

James W Thomas 3/11/13 7:08 PM
Deleted: 4

ISO/IEC JTC 1/SC 22/WG 14 N1685

Date: yyyy-mm-dd

Reference number of document: ISO/IEC TS 18661

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

Secretariat: ANSI

James W Thomas 3/4/13 6:55 PM
Deleted: WG 14 Nnnnn .

James W Thomas 3/4/13 6:55 PM
Deleted: 0000

James W Thomas 3/14/13 8:17 PM
Formatted: Font:12 pt, Bold, Font color: Light Blue

James W Thomas 3/4/13 6:55 PM
Deleted: nnn-n

5

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.

Document type: Technical Specification
Document subtype:
Document stage: (20) Preparation
Document language: E

Copyright notice

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:

*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
Web www.iso.org*

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

Violators may be prosecuted.

James W Thomas 3/14/13 8:16 PM
 Deleted: nnnn
 James W Thomas 3/14/13 8:16 PM
 Formatted: Font:11 pt, Bold, Font color: Auto
 James W Thomas 3/11/13 7:08 PM
 Deleted: 4

Contents

Page

	Introduction.....	v
	Background.....	v
	IEC 60559 floating-point standard	v
5	C support for IEC 60559	vi
	Purpose	vii
	More about formats	vii
1	1 Scope	1
2	2 Conformance	1
10	3 Normative references	1
	4 Terms and definitions	2
5	5 C standard conformance	2
	5.1 Freestanding implementations	2
	5.2 Predefined macros	2
15	5.3 Standard headers	2
6	6 Types	3
7	7 Characteristics	7
8	8 Conversions	9
9	9 Constants.....	10
20	10 Expressions.....	11
	11 Mathematics <math.h>	12
	12 Numeric conversion functions <stdlib.h>	22
	13 Complex arithmetic <complex.h>	24
	14 Type-generic macros <tgmath.h>	26
25	Bibliography.....	30

John Benito 3/26/13 12:08 PM
 Deleted: 23

James W Thomas 3/14/13 9:02 PM
 Deleted: Introduction . v .
 Background . v .
 IEC 60559 floating-point standard . v .
 C support for IEC 60559 . vi .
 Purpose . vii .
 More about formats . vii .
 1 . Scope . 1 .
 2 . Conformance . 1 .
 3 . Normative references . 1 .
 4 . Terms and definitions . 2 .
 5 . C standard conformance . 2 .
 5.1 . Freestanding implementations . 2 .
 5.2 . Predefined macros . 2 .
 5.3 . Standard headers . 2 .
 6 . Types . 3 .
 7 . Characteristics . 5 .
 8 . Conversions . 7 .
 9 . Constants . 8 .
 10 . Expressions . 9 .
 11 . Mathematics <math.h> . 9 .
 12 . Numeric conversion functions <stdlib.h> . 18 .
 Bibliography . 20 .

James W Thomas 3/11/13 7:08 PM

Deleted: 4

James W Thomas 3/14/13 8:16 PM

Deleted: Nnnnn

James W Thomas 3/14/13 8:16 PM

Formatted: Font:11 pt, Bold, Font color: Auto

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: *Supplemental functions*
- Part 5: *Supplemental 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.

James W Thomas 3/4/13 6:55 PM

Deleted: nnn-n

James W Thomas 3/4/13 6:55 PM

Deleted: nnn

James W Thomas 3/4/13 6:55 PM

Deleted: 2008

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:

- 30 • *formats* - for binary floating-point data, including representations for Not-a-Number (NaN) and signed infinities and zeros
- *operations* - basic arithmetic operations (addition, multiplication, etc.) on the format data to compose a well-defined, closed arithmetic system (It also specified conversions between floating-point formats and decimal character sequences, and a few auxiliary operations.)
- 35 • *context* - status flags for detecting exceptional conditions (invalid operation, division by zero, overflow, underflow, and inexact) and controls for choosing different rounding methods

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.

5 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.

10 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 inquiries 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.

15 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.

20 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.

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.

35 C support for IEC 60559

40 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.

45 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.

James W Thomas 3/14/13 8:16 PM

Deleted: nnnn

James W Thomas 3/14/13 8:16 PM

Formatted: Font:11 pt, Bold, Font color: Auto

James W Thomas 3/11/13 7:08 PM

Deleted: 4

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.

- 5 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".

10 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

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.

- 15 Part 1 provides suggested 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 suggested changes in Part 1.

20 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 (this document), 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.

25 More about 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)
- 30 • 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
- extended formats, for each basic format, with minimum range and precision specified

35 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.

40 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.

45 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

James W Thomas 3/4/13 6:55 PM

Deleted: format

James W Thomas 3/4/13 6:55 PM

Deleted: where

James W Thomas 3/4/13 6:55 PM

Deleted: may be done

James W Thomas 3/11/13 7:08 PM

Deleted: 4

James W Thomas 3/14/13 8:16 PM

Deleted: Nnnnn

James W Thomas 3/14/13 8:16 PM

Formatted: Font:11 pt, Bold, Font color: Auto

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 this Technical Specification.

5 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.

James W Thomas 3/4/13 6:55 PM

Deleted: C support for IEC 60559 .

James W Thomas 3/4/13 6:55 PM

Deleted: this

10 Part 2 of Technical Specification 18661 defines types `_Decimal32`, `_Decimal64`, and `_Decimal128` with IEC 60559 formats `decimal32`, `decimal64`, and `decimal128`. Although IEC 60559 regards `decimal32` as an interchange format, not a basic format, and does not require `decimal32` arithmetic (other than conversions), Part 2 of Technical Specification 18661 has full arithmetic and library support for `_Decimal32`, just like for `_Decimal64` and `_Decimal128`.

James W Thomas 3/4/13 6:55 PM

Deleted: this

James W Thomas 3/4/13 6:55 PM

Deleted: this

15 The C language provides just three "generic" floating types (`float`, `double`, and `long double`), which Annex F of the C Standard requires 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.

James W Thomas 3/4/13 6:55 PM

Deleted: formats

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

James W Thomas 3/4/13 6:55 PM

Deleted: Purpose .
The purpose

James W Thomas 3/4/13 6:55 PM

Deleted: specification is to provide

James W Thomas 3/4/13 6:55 PM

Deleted: .
This specification

- 20 · names
- characteristics
- conversions
- constants
- function suffixes
- 25 · character sequence conversion interfaces

30 This specification includes interchange and extended nomenclatures for types that, in some cases, already have C nomenclatures. For example, a type with the IEC 60559 double format may be referred to as `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 interchange and extended formats and C generic floating types.

James W Thomas 3/4/13 6:55 PM

Deleted: proposed

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 nomenclature for types with the interchange and extended floating-point formats specified in ISO/IEC/IEEE 60559:2011.

10 This document proposes nomenclature for all applicable types in Parts 1 and 2 of Technical Specification 18661 and for any other types with IEC 60559 interchange or extended formats supported by the implementation.

2 Conformance

An implementation conforms to Part 3 of Technical Specification 18661 if

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

3 Normative references

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

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

25 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)

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*

30 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*

Suggested changes proposed by Part 3 of Technical Specification 18661 are relative to ISO/IEC 9899:2011 (C11). The actual specification is given by a synthesis with the suggested changes from Parts 1 or 2 or both, depending on which Parts the implementation supports.

35 **ISSUE 1: Will the approach of the previous paragraph lead to a specification that is intelligible? Is there a better approach?**

James W Thomas 3/4/13 6:55 PM
Deleted: 00000

James W Thomas 3/4/13 6:55 PM
Deleted: TS 00000-3

James W Thomas 3/4/13 6:55 PM
Deleted: this

James W Thomas 3/4/13 6:55 PM
Deleted: this

James W Thomas 3/4/13 6:55 PM
Deleted: this

James W Thomas 3/4/13 6:55 PM
Deleted: this

James W Thomas 3/4/13 6:55 PM
Deleted: this

James W Thomas 3/4/13 6:55 PM
Deleted: this

James W Thomas 3/4/13 6:55 PM
Deleted: .
Suggested change to C11: .
 Replace the third sentence of 4#6: .
 A conforming freestanding implementation shall accept any strictly conforming program that does not use complex types and in which the use of the features specified in the library clause (clause 7) is confined to the contents of the standard headers `<float.h>`, `<iso646.h>`, `<limits.h>`, `<stdalign.h>`, `<stdarg.h>`, `<stdbool.h>`, `<stddef.h>`, and `<stdnoreturn.h>`.
 with: .
 A conforming freestanding implementation shall accept any strictly conforming program that does not use complex types and in which the use of the features specified in the library clause (clause 7) is confined to the contents of the standard headers `<float.h>`, `<iso646.h>`, `<limits.h>`, `<math.h>`, `<stdalign.h>`, `<stdarg.h>`, `<stdbool.h>`, `<stddef.h>`, `<stdint.h>`, and `<stdnoreturn.h>` and the numeric conversion functions (7.22.1) of the standard header `<stdlib.h>`.
 The library functions, macros, and types defined in Part 3 of this Technical Specification are defined by their respective headers if the macro `__STDC_WANT_IEC_00000_EXT3` is defined at the point in the source file `[...]` [1]

James W Thomas 3/4/13 6:55 PM
Deleted: nnnnn

James W Thomas 3/4/13 6:55 PM
Deleted: nnnnn

James W Thomas 3/4/13 6:55 PM
Deleted: this

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*

5 C standard conformance

5.1 Freestanding implementations

The following suggested change to C11 expands the conformance requirements for freestanding implements so that they might conform to this Part of Technical Specification 18661

Suggested change to C11:

Replace the third sentence of 4#6:

A conforming freestanding implementation shall accept any strictly conforming program that does not use complex types and in which the use of the features specified in the library clause (clause 7) is confined to the contents of the standard headers `<float.h>`, `<iso646.h>`, `<limits.h>`, `<stdalign.h>`, `<stdarg.h>`, `<stdbool.h>`, `<stddef.h>`, `<stdint.h>`, and `<stdnoreturn.h>`.

with:

A conforming freestanding implementation shall accept any strictly conforming program that does not use complex types and in which the use of the features specified in the library clause (clause 7) is confined to the contents of the standard headers `<fenv.h>`, `<float.h>`, `<iso646.h>`, `<limits.h>`, `<math.h>`, `<stdalign.h>`, `<stdarg.h>`, `<stdbool.h>`, `<stddef.h>`, `<stdint.h>`, and `<stdnoreturn.h>` and the numeric conversion functions (7.22.1) of the standard header `<stdlib.h>`.

5.2 Predefined macros

Suggested change to C11:

In 6.10.8.3#1, add:

`__STDC_IEC_60559_TYPES__` The integer constant `201ymmL`, intended to indicate support of interchange and extended types according to IEC 60559.

5.3 Standard headers

The library functions, macros, and types defined in this Part of Technical Specification 18661 are defined by their respective headers if the macro `__STDC_WANT_IEC_18661_EXT3__` is defined at the point in the source file where the appropriate header is first included.

ISSUE 2: In Part 1 the content in 5.3 above is at the end of 5.1 - which should be fixed.

6 Types

This clause recommends changes to C11 to include the interchange and extended types specified in IEC 60559.

Suggested change to C11:

- 5 Change the first sentence of 6.2.5#10 from:

[10] There are three *real floating types*, designated as `float`, `double`, and `long double`

to:

[10] There are three *generic floating types*, designated as `float`, `double`, and `long double`.

After 6.2.5#10, insert:

- 10 [10a] IEC 60559 specifies interchange formats, identified by their width, which can be used for the exchange of floating-point data between implementations. Tables 1 and 2 give parameters for the IEC 60559 interchange formats.

Table 1 – Binary interchange format parameters

Parameter	binary16	binary32	binary64	binary128	binaryN ($N \geq 128$)
N , storage width in bits	16	32	64	128	multiple of 32
p , precision in bits	11	24	53	113	$N - \text{round}(4 \times \log_2(N)) + 13$
e_{\max} , maximum exponent e	15	127	1023	16383	$2^{(N-p-1)} - 1$
<i>Encoding parameters</i>					
$bias$, $E-e$	15	127	1023	16383	e_{\max}
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	$1 + w + t$

15

The function `round()` in Table 1 rounds to the nearest integer.

For example, `binary256` would have $p = 237$ and $e_{\max} = 262143$.

James W Thomas 3/11/13 7:08 PM
Deleted: 4
 James W Thomas 3/14/13 8:16 PM
Deleted: Nnnnn
 James W Thomas 3/14/13 8:16 PM
Formatted: Font:11 pt, Bold, Font color: Auto

Table 2 – Decimal interchange format parameters

Parameter	decimal32	decimal64	decimal128	decimalN (N ≥ 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 $emax = 1572864$.

Types designated

`_FloatN`, where N is 16, 32, 64, or ≥ 128 and a multiple of 32

`_DecimalN`, where N ≥ 32 and a multiple of 32

support the corresponding IEC 60559 interchange formats and are collectively called the *data-interchange types*. Each data-interchange type has the IEC 60559 interchange format corresponding to its width and radix. Data-interchange types that are supported by all applicable floating-point operations are collectively called the *interchange floating types*. Data-interchange types (including interchange floating types) are not compatible with any other types. An implementation that defines `__STDC_IEC_60559_BFP__` shall provide

`_Float32` and `_Float64` as interchange floating types with the same representation and alignment requirements as `float` and `double`, respectively, and

`_Float16` as a data-interchange type

and may provide:

`_Float16` as an interchange floating type

If the implementation's `long double` type has an IEC 60559 interchange format of width N, then the implementation shall provide the type `_FloatN` as an interchange floating type with the same representation and alignment requirements as `long double`. An implementation may provide any of the other data-interchange types and may provide any of its data-interchange types as interchange floating types.

James W Thomas 3/4/13 6:55 PM
Deleted: The implementation may provide any of the
 James W Thomas 3/4/13 6:55 PM
Deleted: .
 [10b] The implementation may provide an interchange type as an *interchange floating type*. An interchange floating type is
 James W Thomas 3/4/13 6:55 PM
Deleted: and adheres to IEC 60559. .
 [10c]
 James W Thomas 3/12/13 11:22 AM
Deleted:
 James W Thomas 3/12/13 11:20 AM
Deleted: equivalent
 James W Thomas 3/12/13 11:24 AM
Formatted: Indent: Left: 0.75"
 James W Thomas 3/12/13 11:21 AM
Deleted: to
 James W Thomas 3/4/13 6:55 PM
Deleted: and shall provide
 James W Thomas 3/12/13 11:24 AM
Formatted: Indent: Left: 0.79"
 James W Thomas 3/4/13 6:55 PM
Deleted: type. The `_Float16` type need not be an interchange
 James W Thomas 3/12/13 11:21 AM
Deleted: .
 James W Thomas 3/12/13 11:23 AM
Formatted: Font:(Default) Courier New, Bold

5

10

15

20

[10b] For each of its basic formats, IEC 60559 specifies as extended format whose maximum exponent and precision exceed those of the basic format it is associated with. Table 3 below gives the minimum values of these parameters:

Table 3 – Extended format parameters for floating-point numbers

Parameter	Extended formats associated with:				
	binary32	binary64	binary128	decimal64	decimal128
<i>p</i> digits \geq	32	64	128	22	40
<i>emax</i> \geq	1023	16383	65535	6144	24576

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` shall provide `_Float32x`, which may be equivalent to `double`, and may provide any of the other two binary extended floating types. An implementation that defines `_STDC_IEC_60559_DFP` shall provide: `_Decimal64x`, which may be equivalent to `_Decimal128`, and may provide `_Decimal128x`.

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

Replace 6.2.5#11:

[11] There are three *complex types*, designated as `float _Complex`, `double _Complex`, and `long double _Complex`. 43) (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:

[11] For the generic 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`. 43) (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*.

In 6.2.5#14, change the first sentence from:

[14] The type `char`, the signed and unsigned integer types, and the floating types are collectively called the *basic types*. ...

to:

[14] The type `char`, the signed and unsigned integer types, the floating types, and the data-interchange types are collectively called the *basic types*. ...

In 6.2.5#21, change the first sentence from:

[21] Arithmetic types and pointer types are collectively called *scalar types*. ...

to:

- James W Thomas 3/14/13 8:16 PM Deleted: nnnn
- James W Thomas 3/11/13 7:08 PM Deleted: 4
- James W Thomas 3/14/13 8:16 PM Formatted: Font:11 pt, Bold, Font color: Auto
- James W Thomas 3/4/13 6:55 PM Deleted: 10d
- James W Thomas 3/4/13 6:55 PM Deleted: Table 3 goes here

- James W Thomas 3/4/13 6:55 PM Deleted: 10e
- James W Thomas 3/4/13 6:55 PM Deleted: ln
- James W Thomas 3/4/13 6:55 PM Deleted: , after the first sentence, insert
- James W Thomas 3/4/13 6:55 PM Deleted: For each supported interchange
- James W Thomas 3/4/13 6:55 PM Deleted: type
- James W Thomas 3/12/13 9:59 AM Formatted: Font:(Default) Courier New, Bold, Not Italic
- James W Thomas 3/4/13 6:55 PM Deleted: is a
- James W Thomas 3/4/13 6:55 PM Deleted: type
- James W Thomas 3/11/13 7:11 PM Formatted: Normal
- James W Thomas 3/12/13 11:29 AM Formatted: Font:(Default) Courier New
- James W Thomas 3/11/13 7:11 PM Formatted: Normal

[21] Arithmetic types, data-interchange types, and pointer types are collectively called *scalar types*. ...

Add the following to 6.4.1 Keywords:

keyword:

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

Add the following to 6.7.2 Type specifiers:

type-specifier:

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

Add the following bullets in 6.7.2#2 Constraints:

- `_FloatN`, where N is 16, 32, 64, or ≥ 128 and a multiple of 32
- `_Float32x`
- `_Float64x`
- `_Float128x`
- `_DecimalN`, where $N \geq 32$ and a multiple of 32
- `_Decimal64x`
- `_Decimal128x`
- `_FloatN_Complex`, where N is 16, 32, 64, or ≥ 128 and a multiple of 32
- `_Float32x_Complex`
- `_Float64x_Complex`
- `_Float128x_Complex`

Add the following after 6.7.2#3:

[3a] The type specifiers `_FloatN` (where N is 16, 32, 64, or ≥ 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 interchange and extended types (see 6.10.8.3).

Add the following after 6.5#8:

[8a] Expressions involving operands of interchange or extended 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.

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.

with:

[2] For the generic real 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 suggests new `<float.h>` macros, analogous to the macros for generic floating types, that characterize the data-interchange types and the extended floating types. It also suggests macros to indicate which data-interchange types are provided as interchange floating types.

Suggested changes to C11:

In 5.2.4.2.2#7, change the sentence:

All except `DECIMAL_DIG`, `FLT_EVAL_METHOD`, `FLT_RADIX`, and `FLT_ROUNDS` have separate names for all three floating-point types.

to:

All except `DECIMAL_DIG`, `FLT_EVAL_METHOD`, `FLT_RADIX`, `FLT_ROUNDS`, `FLT_N_IS_ARITH`, and `DEC_N_IS_ARITH` have separate names for all floating-point types.

After 5.2.4.2.2#7, add the paragraph

[7a] Some of the macros in `<float.h>` provide characteristics of data-interchange types and extended floating types, as specified in IEC 60559. The prefixes `FLT_N_` and `DEC_N_` are used for binary and decimal data-interchange types of width *N*. The prefixes `FLT_Nx_` and `DEC_Nx_` are used for binary and decimal extended floating types that extend a basic format of width *N*. For each data-interchange or extended floating type that the implementation provides, `<float.h>` shall define the associated macros. Conversely, for each such type that the implementation does not provide, `<float.h>` shall not define the associated macros.

In 5.2.4.2.2#11, add the following in the bullet defining `FLT_DECIMAL_DIG`, etc.:

`FLT_N_DECIMAL_DIG`
`FLT_Nx_DECIMAL_DIG`

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:

James W Thomas 3/14/13 8:16 PM
Deleted: nnnn

James W Thomas 3/14/13 8:16 PM
Formatted ... [2]

James W Thomas 3/11/13 7:08 PM
Deleted: 4

James W Thomas 3/12/13 9:18 AM
Formatted ... [3]

James W Thomas 3/4/13 6:55 PM
Deleted: ln

James W Thomas 3/4/13 6:55 PM
Deleted: , after the first sentence, insert

James W Thomas 3/4/13 6:55 PM
Deleted: each supported

James W Thomas 3/4/13 6:55 PM
Deleted: type

James W Thomas 3/4/13 6:55 PM
Deleted: is an

James W Thomas 3/4/13 6:55 PM
Deleted: type

James W Thomas 3/12/13 9:57 AM
Deleted: FloatN

James W Thomas 3/4/13 6:55 PM
Deleted: [4]

James W Thomas 3/14/13 8:37 PM
Formatted ... [5]

James W Thomas 3/14/13 8:37 PM
Formatted ... [6]

James W Thomas 3/14/13 8:37 PM
Formatted ... [7]

James W Thomas 3/14/13 8:38 PM
Formatted ... [8]

James W Thomas 3/14/13 8:38 PM
Formatted ... [9]

James W Thomas 3/14/13 8:35 PM
Formatted ... [10]

James W Thomas 3/14/13 8:41 PM
Formatted ... [11]

James W Thomas 3/14/13 8:41 PM
Formatted ... [12]

James W Thomas 3/12/13 3:02 PM
Formatted ... [13]

James W Thomas 3/12/13 3:02 PM
Formatted ... [14]

James W Thomas 3/12/13 1:44 PM
Formatted ... [15]

James W Thomas 3/12/13 2:43 PM
Formatted ... [16]

James W Thomas 3/12/13 1:40 PM
Formatted ... [17]

James W Thomas 3/12/13 1:40 PM
Formatted ... [18]

James W Thomas 3/11/13 7:08 PM

Deleted: 4

James W Thomas 3/14/13 8:16 PM

Deleted: Nnnnn

James W Thomas 3/14/13 8:16 PM

Formatted: Font:11 pt, Bold, Font color: Auto

James W Thomas 3/12/13 1:26 PM

Formatted: English (US)

James W Thomas 3/12/13 1:26 PM

Formatted: std text, Indent: Left: 0.25", Hanging: 0.19"

James W Thomas 3/12/13 2:57 PM

Deleted: Macros in <float.h> provide characteristics of interchange types and extended floating types, as specified in IEC 60559, for all the interchange types and extended floating types supported by the implementation. The prefixes `FLT_N_` and `DEC_N_` are used for binary and decimal interchange types of width N . The prefixes `FLT_Nx_` and `DEC_Nx_` are used for binary and decimal extended floating types that extend a basic format of width N . For extended floating types, the type parameters p , e_{max} , and e_{min} are for the extended floating type itself, not for the basic format that it extends.

James W Thomas 3/4/13 6:55 PM

Deleted: with the interger constant 1

James W Thomas 3/4/13 6:55 PM

Deleted: respectfully

James W Thomas 3/4/13 6:55 PM

Deleted: ISSUE 4: Why is "coefficient" used here instead of "significant"? It is in Part 2 too.

— number of decimal digits, n , such that any floating-point number in the widest of the supported floating and data-interchange types with ...

Add the following after 5.2.4.2.2#13:

[13a] The macros

```
FLT_N_IS_ARITH
DEC_N_IS_ARITH
```

are defined if the implementation supports the data-interchange types `_FloatN` and `_DecimalN`, respectively, as interchange floating types. If defined the macros have the integer value 1.

[13b] In the following lists, the type parameters p , e_{max} , and e_{min} for extended floating types are for the extended floating type itself, not for the basic format that it extends.

[13c] The integer values given in the following lists shall be expressed by constant expressions suitable for use in `#if` preprocessing directives:

— number of digits in the floating-point significand, p

```
FLT_N_MANT_DIG
FLT_Nx_MANT_DIG
```

— number of digits in the coefficient, p

```
DEC_N_MANT_DIG
DEC_Nx_MANT_DIG
```

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

```
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 binary digits and back again without change to the q decimal digits, $\text{floor}((p - 1) \log_{10} 2)$

```
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}

```
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, $\text{ceiling}(\log_{10} 2^{e_{min}-1})$

```
FLT_N_MIN_10_EXP
FLT_Nx_MIN_10_EXP
```


James W Thomas 3/14/13 8:16 PM

Deleted: nnnn

James W Thomas 3/14/13 8:16 PM

Formatted: Font:11 pt, Bold, Font color: Auto

James W Thomas 3/11/13 7:08 PM

Deleted: 4

- 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_NX_MAX_EXP`
`DEC_N_MAX_EXP`
`DEC_NX_MAX_EXP`

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

10
`FLT_N_MAX_10_EXP`
`FLT_NX_MAX_10_EXP`

[13d] The values given in the following list shall be replaced by constant expressions:

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

15
`FLT_N_MAX`
`FLT_NX_MAX`
`DEC_N_MAX`
`DEC_NX_MAX`

- 20
 — the difference between 1 and the least value greater than 1 that is representable in the given floating-point type, b^{1-p}

25
`FLT_N_EPSILON`
`FLT_NX_EPSILON`
`DEC_N_EPSILON`
`DEC_NX_EPSILON`

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

30
`FLT_N_MIN`
`FLT_NX_MIN`
`DEC_N_MIN`
`DEC_NX_MIN`

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

35
`FLT_N_TRUE_MIN`
`FLT_NX_TRUE_MIN`
`DEC_N_TRUE_MIN`
`DEC_NX_TRUE_MIN`

8 Conversions

- 40 The following suggested change to C11 supports the IEC 60559 restrictions against operands whose types do not have the same radix or such that neither type is a subset of (or equivalent to) the other.

Suggested change to C11:

In 6.3.1.8#1, replace the first 3 items after “This pattern is called the *usual arithmetic conversions*”:

James W Thomas 3/4/13 6:55 PM

Deleted: with implementation-defined (positive) values

James W Thomas 3/4/13 6:55 PM

Deleted: After 6.3.1.8#2, add . [2a]

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`.

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:

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

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

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.

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

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.

The following suggested change to C11 provides conversions between data-interchange types and other data-interchange types and real floating types.

After 6.3.2.3, add the subclause:

6.3.2.3a Data-interchange types

[1] Any supported data-interchange type can be converted to and from any supported data-interchange type and any real floating type, with rounding to IEC 60559 formats as specified in IEC 60559.

9 Constants

The following suggested changes to C11 provide suffixes that designate constants of data-interchange types and extended floating types.

Suggested changes to C11:

- James W Thomas 3/11/13 7:08 PM Deleted: 4
- James W Thomas 3/14/13 8:16 PM Deleted: Nnnnn
- James W Thomas 3/14/13 8:16 PM Formatted: Font:11 pt, Bold, Font color: Auto
- James W Thomas 3/14/13 8:56 PM Formatted: Indent: Left: 0.56"
- James W Thomas 3/14/13 8:54 PM Formatted: Font:(Default) Courier New
- James W Thomas 3/14/13 8:55 PM Formatted: Font:(Default) Courier New
- James W Thomas 3/14/13 8:55 PM Formatted: Font:(Default) Courier New
- James W Thomas 3/14/13 8:55 PM Formatted: Font:(Default) Courier New
- James W Thomas 3/14/13 8:56 PM Formatted: std text, Indent: Left: 0.56"
- James W Thomas 3/14/13 8:56 PM Formatted: Font:(Default) Courier New
- James W Thomas 3/14/13 8:56 PM Formatted: Font:(Default) Courier New
- James W Thomas 3/14/13 8:57 PM Formatted: Indent: Left: 0.56"
- James W Thomas 3/4/13 6:55 PM Deleted: is
- James W Thomas 3/4/13 6:55 PM Deleted: be
- James W Thomas 3/4/13 6:55 PM Deleted: of
- James W Thomas 3/4/13 6:55 PM Deleted: [2b]
- James W Thomas 3/4/13 6:55 PM Deleted: the numerical values represent ... [19]
- James W Thomas 3/4/13 6:55 PM Deleted: type A compose a strict subset ... [20]
- James W Thomas 3/14/13 8:57 PM Formatted ... [21]
- James W Thomas 3/12/13 2:12 PM Deleted: both operands are real floatir ... [22]
- James W Thomas 3/14/13 8:57 PM Formatted ... [23]
- James W Thomas 3/4/13 6:55 PM Deleted: with narrower
- James W Thomas 3/4/13 6:55 PM Deleted: types of the operands are not ... [24]
- James W Thomas 3/4/13 6:55 PM Deleted: the
- James W Thomas 3/4/13 6:55 PM Deleted: wider corresponding real type. If
- James W Thomas 3/4/13 6:55 PM Deleted: ISSUE 6: This isn't an adequate ... [25]

James W Thomas 3/14/13 8:16 PM
Formatted: Font:11 pt, Bold, Font color: Auto

James W Thomas 3/14/13 8:16 PM

Deleted: nnnn

James W Thomas 3/11/13 7:08 PM

Deleted: 4

Change *floating-suffix* in 6.4.4.2 from:

floating-suffix: one of
f l F L

to:

- 5 *floating-suffix*: one of
f l F L
fN FN - for each supported binary data-interchange type `_FloatN`
fNx FNx - for each supported binary extended floating type `_FloatNx`
10 dN DN - for each supported decimal data-interchange type `_DecimalN`
dNx DNx - for each supported decimal extended floating type `_DecimalNx`

Add the following paragraph after 6.4.4.2#2:

Constraints

[2a] A *floating-suffix* dN, DN, dNx, or DNx shall not be used in a *hexadecimal-floating-constant*.

Add the following paragraph after 6.4.4.2#4:

- 15 [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 dN or DN, it has type `_DecimalN`. If suffixed by dNx or DNx, it has type `_DecimalNx`.

Add the following paragraph after 6.4.4.2#5:

- 20 [5a] Decimal floating-point constants that have the same numerical value but different quantum exponents have distinguishable internal representations. The quantum exponent is specified to be the same as for the corresponding `strtodN` or `strtodNx` function for the same numeric string.

10 Expressions

- 25 The following suggested change to C11 is intended to ensure that data-interchange types can be converted to and from real floating types and data-interchange types, by assignment, cast, argument passing, and function return.

Suggested changes to C11:

At the end of 6.5.16.1#1, append the bullet:

- the left operand has atomic, qualified, or unqualified real floating type or data-interchange type, and the right has real floating type or data-interchange type;

- 30 The following suggested changes to C11 specify that certain arithmetic operators need not handle operands of data-interchange type.

Suggested changes to C11:

Change 6.5.3.3#1 from:

- 35 [1] The operand of the unary + or - operator shall have arithmetic type; of the ~ operator, integer type; of the ! operator, scalar type.

to:

James W Thomas 3/4/13 6:55 PM

Deleted: [2a]

James W Thomas 3/11/13 7:17 PM

Deleted: allow

James W Thomas 3/11/13 7:17 PM

Deleted: to

James W Thomas 3/11/13 7:16 PM

Formatted: std text, Indent: Left: 0.25", Widow/Orphan control, Adjust space between Latin and Asian text, Adjust space between Asian text and numbers, Tabs:Not at 0.15" + 0.5"

James W Thomas 3/11/13 7:20 PM

Formatted: Justified, Indent: Left: 0", First line: 0", Line spacing: at least 11.5 pt, Widow/Orphan control, Adjust space between Latin and Asian text, Adjust space between Asian text and numbers, Tabs:Not at 0.15" + 0.5"

James W Thomas 3/11/13 7:08 PM
 Deleted: 4
 James W Thomas 3/14/13 8:16 PM
 Deleted: Nnnnn
 James W Thomas 3/14/13 8:16 PM
 Formatted: Font:11 pt, Bold, Font color: Auto

[1] The operand of the unary + or - operator shall have arithmetic type; of the ~ operator, integer type; of the ! operator, floating or pointer type.

Change 6.5.13#2 from:

[2] Each of the operands shall have scalar type.

to:

[2] Each of the operands shall have floating or pointer type.

Change 6.5.14#2 from:

[2] Each of the operands shall have scalar type.

to:

[2] Each of the operands shall have floating or pointer type.

Change 6.5.15#2 from:

[2] The first operand shall have scalar type.

to:

[2] The first operand shall have floating or pointer type.

In 7.2.1.1#1, change:

```
void assert(scalar expression);
```

to:

```
void assert(floating or pointer expression);
```

In 7.2.1.1#2, change the second sentence from:

When it is executed, if *expression* (which shall have a scalar type) is false ...

to:

When it is executed, if *expression* (which shall have a floating or pointer type) is false ...

11 Mathematics <math.h>

This clause suggests changes to C11 to include functions and macros for interchange and extended floating types. Binary interchange floating types and binary extended floating types are supported by functions and macros corresponding to all those specified for generic floating types (`float`, `double`, and `long double`) in C11 and Part 1 of Technical Specification 18661. Decimal interchange floating types and decimal extended floating types are supported by functions and macros corresponding to all those for the decimal types in Part 2 of Technical Specification 18661. Data-interchange types (including ones that are not interchange floating types) are supported by the classification macros in C11 and Parts 1 and 2 of Technical Specification 18661, and by the `totalorder` and `totalordermag` functions in Parts 1 and 2.

The list of elementary functions specified in the mathematics library is extended to handle interchange floating types and extended floating types. These include functions specified in C11 (7.12.4, 7.12.5, 7.12.6, 7.12.7, 7.12.8, 7.12.9, 7.12.10, 7.12.11, 7.12.12, and 7.12.13) and in Part 1 of Technical Specification 18661 (14.1, 14.2, 14.3, 14.4, 14.5, 14.8, 14.9, and 14.10). Macros analogous to the `HUGE_VAL`, `INFINITY`, `NAN`, and the

James W Thomas 3/11/13 7:39 PM
 Formatted: Normal
 James W Thomas 3/11/13 7:43 PM
 Formatted: Font:(Default) Courier New
 James W Thomas 3/11/13 7:44 PM
 Formatted: Font:(Default) Courier New
 James W Thomas 3/11/13 7:44 PM
 Formatted: Indent: Left: 0.53", First line: 0.03"
 James W Thomas 3/11/13 7:44 PM
 Formatted: Font:(Default) Courier New, Bold, Not Italic
 James W Thomas 3/11/13 7:39 PM
 Formatted: Normal
 James W Thomas 3/11/13 7:50 PM
 Formatted: Indent: Left: 0.25", No bullets or numbering
 James W Thomas 3/12/13 9:21 AM
 Formatted: Font:(Default) Courier New
 James W Thomas 3/12/13 9:21 AM
 Formatted: Font:(Default) Courier New
 James W Thomas 3/11/13 7:33 PM
 Deleted: .
 James W Thomas 3/4/13 6:55 PM
 Deleted: this
 James W Thomas 3/4/13 6:55 PM
 Deleted: this
 James W Thomas 3/4/13 6:55 PM
 Deleted: .
 James W Thomas 3/4/13 6:55 PM
 Deleted: this
 James W Thomas 3/4/13 6:55 PM
 Deleted: this
 James W Thomas 3/4/13 6:55 PM
 Deleted: .

SNAN macros are defined for data-interchange types and extended floating types. With the exception of the floating-point functions listed in 11.2, which have accuracy as specified in IEC 60559, the accuracy of floating-point results is implementation-defined. The implementation may state that the accuracy is unknown. All comparison macros specified in C11 (7.12.14) and in Part 1 of Technical Specification 18661 (14.6) are extended to handle interchange floating types and extended floating types. All classification macros specified in C11 (7.12.3) and in Part 1 of Technical Specification 18661 (14.7) are extended to handle data-interchange types and extended floating types.

Suggested changes to C11:

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 type **float**, **long double**, **_FloatN**, **_FloatNx**, **_DecimalN**, and **_DecimalNx**, respectively.

Add after 7.12#1:

[1a] 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.

Add after 7.12#2:

[2a] For each decimal data-interchange type that the implementation provides, one of the types

```
decencodingdN_t
binencodingdN_t
```

is declared and represents values of the type in one of the two alternative encodings allowed for decimal formats by the IEC 60559 standard: the encoding (indicated by the prefix **dec**) based on decimal encoding of the significand, or the encoding (indicated by the prefix **bin**) based on binary encoding of the significand. These types are used by the decimal re-encoding functions (7.12.11).

ISSUE 3: Do we need types corresponding to float_t and double_t for the interchange floating types? Is this more appropriate for Part 5 which might deal with expression evaluation methods?

Add at the end of 7.12#3 the following macros:

[3] ... For each data-interchange type (**_FloatN** and **_DecimalN**) that the implementation provides, the corresponding one of the macros

```
HUGE_VAL_FN
HUGE_VAL_DN
```

is defined and expands to a constant expression of the type representing positive infinity. The macros

James W Thomas 3/14/13 8:16 PM
Formatted: Font:11 pt, Bold, Font color: Auto

James W Thomas 3/14/13 8:16 PM
Deleted: nnnn...4 ... [26]

James W Thomas 3/14/13 8:34 AM
Deleted: N...this ...this ... [27]

James W Thomas 3/14/13 8:34 AM
Formatted: Font:Not Italic

James W Thomas 3/4/13 6:55 PM
Deleted: s

James W Thomas 3/13/13 10:36 AM
Deleted:

James W Thomas 3/13/13 10:37 AM
Formatted: Indent: Left: 0.56"

James W Thomas 3/13/13 10:30 AM
Formatted ... [28]

James W Thomas 3/13/13 10:29 AM
Deleted: .

James W Thomas 3/13/13 11:41 AM
Deleted: The implementation provides all the functions for its supported interchange floating types and extended floating types.

James W Thomas 3/13/13 10:44 AM
Formatted: Normal, Indent: Left: 0.25"

James W Thomas 3/13/13 10:44 AM
Formatted: Font:Not Bold

James W Thomas 3/13/13 10:43 AM
Deleted: paragraph [29]

James W Thomas 3/4/13 6:55 PM
Deleted: ISSUE: _t types?The ... [30]

James W Thomas 3/13/13 10:48 AM
Deleted: dpd

James W Thomas 3/13/13 10:48 AM
Deleted: d

James W Thomas 3/13/13 11:20 AM
Deleted: represent values of decimal interchange types in one of the two alternative encodings required by the IEC 60559 standard: the Densely Packed Decimal (DPD) encoding or the Binary Integer Decimal (BID) enc(... [31]

James W Thomas 3/13/13 10:45 AM
Deleted: paragraph... [32]

James W Thomas 3/13/13 3:00 PM
Deleted: supported

James W Thomas 3/13/13 3:13 PM
Formatted ... [33]

James W Thomas 3/13/13 3:02 PM
Deleted: _FloatN and ... [34]

James W Thomas 3/13/13 3:15 PM
Deleted: s...s... _FloatN and ... [35]

James W Thomas 3/11/13 7:08 PM
Deleted: 4
 James W Thomas 3/14/13 8:16 PM
Deleted: Nnnnn
 James W Thomas 3/14/13 8:16 PM
Formatted: Font:11 pt, Bold, Font color: Auto

HUGE_VAL_DNx
 HUGE_VAL_DNx

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

5 After 7.12#5, add the following:

[5b] For each data-interchange type (`_FloatN` and `_DecimalN`) that the implementation provides, the corresponding one of the signaling NaN macros

SNANFN
 SNANDN

is defined and expands into a constant expression of the type representing a signaling NaN. The signaling NaN macros

SNANFNx
 SNANDNx

expand into constant expressions of types `_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.

Add at the end of 7.12 paragraph 7 the following macros.

[7] ... The macros

FP_FAST_FMAFN
 FP_FAST_FMADN

are, respectively, `_FloatN` and `_DecimalN` analogs of `FP_FAST_FMA`. The macros

FP_FAST_FMAFNx
 FP_FAST_FMADNx

are, respectively, `_FloatNx` and `_DecimalNx` analogs of `FP_FAST_FMA`.

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

7.12.4 Trigonometric functions

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

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

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

James W Thomas 3/13/13 3:21 PM
Deleted: Add at the end of
 James W Thomas 3/13/13 3:21 PM
Deleted: paragraph
 James W Thomas 3/13/13 3:21 PM
Deleted:
 James W Thomas 3/13/13 3:22 PM
Deleted: macros.
 James W Thomas 3/13/13 3:16 PM
Deleted: supported interface
 James W Thomas 3/13/13 3:16 PM
Deleted: type
 James W Thomas 3/13/13 3:16 PM
Deleted: and
 James W Thomas 3/13/13 3:17 PM
Deleted: ,
 James W Thomas 3/13/13 3:18 PM
Deleted: s
 James W Thomas 3/13/13 3:19 PM
Deleted: type `_FloatN` and `_DecimalN`, respectively,
 James W Thomas 3/13/13 3:22 PM
Deleted: 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. For each supported extended floating type `_FloatNx` and `_DecimalNx`,
 James W Thomas 3/13/13 3:20 PM
Deleted: t
 James W Thomas 3/13/13 3:23 PM
Deleted: For each supported interface floating type `_FloatN` and `_DecimalN`, t
 James W Thomas 3/13/13 3:24 PM
Deleted: For each supported extended floating type `_FloatN` and `_DecimalN`, t

```

_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);

```

20

7.12.5 Hyperbolic functions

```

_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 sinhfN (_FloatN x);
_FloatNx sinhfNx(_FloatNx x);
_DecimalN sinhndN (_DecimalN x);
_DecimalNx sinhndNx(_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);

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

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

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

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

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

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

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

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

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

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

55  _FloatN logbfN (_FloatN x);
    _FloatNx logbfNx(_FloatNx x);
```



```
DecimalN logbdN (DecimalN x);
DecimalNx logbdNx(DecimalNx x);
```

```
5 FloatN modffN (FloatN x, FloatN *iptr);
FloatNx modffNx(FloatNx x, FloatNx *iptr);
DecimalN modfdN (DecimalN x, DecimalN *iptr);
DecimalNx modfdNx(DecimalNx x, DecimalNx *iptr);
```

```
10 FloatN scalbnfN (FloatN value, int exp);
FloatNx scalbnfN (FloatNx value, int exp);
DecimalN scalbndN (DecimalN value, int exp);
DecimalNx scalbndNx (DecimalNx value, int exp);
```

```
15 FloatN scalblnfN (FloatN value, long int exp);
FloatNx scalblnfN (FloatNx value, long int exp);
DecimalN scalblndN (DecimalN value, long int exp);
DecimalNx scalblndNx (DecimalNx value, long int exp);
```

7.12.7 Power and absolute-value functions

```
20 FloatN cbrtfN (FloatN x);
FloatNx cbrtfNx(FloatNx x);
DecimalN cbrtdN (DecimalN x);
DecimalNx cbrtdNx(DecimalNx x);
```

```
25 FloatN fabsfN (FloatN x);
FloatNx fabsfNx(FloatNx x);
DecimalN fabsdN (DecimalN x);
DecimalNx fabsdNx(DecimalNx x);
```

```
30 FloatN hypotfN (FloatN x, FloatN y);
FloatNx hypotfNx(FloatNx x, FloatNx y);
DecimalN hypotdN (DecimalN x, DecimalN y);
DecimalNx hypotdNx(DecimalNx x, DecimalNx y);
```

```
35 FloatN powfN (FloatN x, FloatN y);
FloatNx powfNx(FloatNx x, FloatNx y);
DecimalN powdN (DecimalN x, DecimalN y);
DecimalNx powdNx(DecimalNx x, DecimalNx y);
```

```
40 FloatN sqrtfN (FloatN x);
FloatNx sqrtfNx(FloatNx x);
DecimalN sqrtdN (DecimalN x);
DecimalNx sqrtdNx(DecimalNx x);
```

7.12.8 Error and gamma functions

```
45 FloatN erffN (FloatN x);
FloatNx erffNx(FloatNx x);
DecimalN erfdN (DecimalN x);
DecimalNx erfdNx(DecimalNx x);
```

```
50 FloatN erfcfN (FloatN x);
FloatNx erfcfNx(FloatNx x);
DecimalN erfcdN (DecimalN x);
DecimalNx erfcdNx(DecimalNx x);
```

```

_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 (_DecimalN x);
```

```
__FloatN truncfN (__FloatN x);
```

```
__FloatNx truncfNx(__FloatNx x);
```

```
__DecimalN truncdN (__DecimalN x);
```

```
__DecimalNx truncdNx(__DecimalNx x);
```

```
__FloatN roundevenfN (__FloatN x);
```

```
__FloatNx roundevenfNx(__FloatNx x);
```

```
__DecimalN roundevendN (__DecimalN x);
```

```
__DecimalNx roundevendNx(__DecimalNx x);
```

```
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);
```

```
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);
```

```
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);
```

```
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);
```

```
__DecimalN fmoddN (__DecimalN x, __DecimalN y);
```

```
__DecimalNx fmoddNx(__DecimalNx x, __DecimalNx y);
```

```
__FloatN remainderfN (__FloatN x, __FloatN y);
```

```
__FloatNx remainderfNx(__FloatNx x, __FloatNx y);
```

```
__DecimalN remainderdN (__DecimalN x, __DecimalN y);
```

```
__DecimalNx remainderdNx(__DecimalNx x, __DecimalNx y);
```

```
__FloatN remquofN (__FloatN x, __FloatN y, int *quo);
```

```
__FloatNx remquofNx(__FloatNx x, __FloatNx y, int *quo);
```

7.12.11 Manipulation functions

```
__FloatN copysignfN (__FloatN x, __FloatN y);
```

```
__FloatNx copysignfNx(__FloatNx x, __FloatNx y);
```

```
__DecimalN copysigndN (__DecimalN x, __DecimalN y);
```

```
__DecimalNx copysigndNx(__DecimalNx x, __DecimalNx y);
```

```
__FloatN nanfN (const char *tagp);
```

```
__FloatNx nanfNx (const char *tagp);
```

```
__DecimalN nandN (const char *tagp);
```

```
__DecimalNx nandNx (const char *tagp);
```

James W Thomas 3/11/13 7:08 PM

Deleted: 4

James W Thomas 3/14/13 8:16 PM

Deleted: Nnnnn

James W Thomas 3/14/13 8:16 PM

Formatted: Font:11 pt, Bold, Font color: Auto

```

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

```

ISSUE 4: Do we need a version of nexttoward that could be used with interchange floating types?

```

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

```

```

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

```

```

_FloatN canonicalizefN (_FloatN x);
_FloatNx canonicalizefNx(_FloatNx x);
_DecimalN canonicalizedN (_DecimalN x);
_DecimalNx canonicalizedNx(_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);

```

```

int quantexpdN (_DecimalN x);
int quantexpdNx(_DecimalNx x);

```

```

decencodingdN_t encodedecdN (_DecimalN x);
_DecimalN decodedecdN (decencodingdN_t x);
binencodingdN_t encodebindN (_DecimalN x);
_DecimalN decodebindN (binencodingdN_t x);

```

James W Thomas 3/4/13 6:55 PM

Deleted: :

James W Thomas 3/13/13 11:31 AM

Formatted: Indent: Left: 0.56"

7.12.12 Maximum, minimum, and positive difference functions

```

_FloatN fdimfN (_FloatN x, _FloatN y);
_FloatNx fdimfNx(_FloatNx x, _FloatNx y);
_DecimalN fdimdN (_DecimalN x, _DecimalN y);
_DecimalNx fdimdNx(_DecimalNx x, _DecimalNx y);
_FloatN fmaxfN (_FloatN x, _FloatN y);
_FloatNx fmaxfNx(_FloatNx x, _FloatNx y);
_DecimalN fmaxdN (_DecimalN x, _DecimalN y);
_DecimalNx fmaxdNx(_DecimalNx x, _DecimalNx y);

```

```

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

```

```

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

```

```

_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

```

_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

```

_FloatM fMaddfN (_FloatN x, _FloatN y); // M < N
_FloatM fMaddfNx (_FloatNx x, _FloatNx y); // M <= N
_FloatMx fMxaddfN (_FloatN x, _FloatN y); // M < N
_FloatMx fMxaddfNx (_FloatNx x, _FloatNx y); // M < N
_DecimalM dMaddN (_DecimalN x, _DecimalN y); // M < N
_DecimalM dMaddNx (_DecimalNx x, _DecimalNx y); // M <= N
_DecimalMx dMxaddN (_DecimalN x, _DecimalN y); // M < N
_DecimalMx dMxaddNx (_DecimalNx x, _DecimalNx y); // M < N

_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
_DecimalM dMsubdN (_DecimalN x, _DecimalN y); // M < N
_DecimalM dMsubdNx (_DecimalNx x, _DecimalNx y); // M <= N
_DecimalMx dMxsubdN (_DecimalN x, _DecimalN y); // M < N
_DecimalMx dMxsubdNx (_DecimalNx x, _DecimalNx y); // M < N

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

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

```

James W Thomas 3/11/13 7:08 PM

Deleted: 4

James W Thomas 3/14/13 8:16 PM

Deleted: Nnnnn

James W Thomas 3/14/13 8:16 PM

Formatted: Font:11 pt, Bold, Font color: Auto

```

_DecimalM dMsqrtdNx (_DecimalNx x); // M <= N
_DecimalMx dMxdivdN (_DecimalN x); // M < N
_DecimalMx dMxdivdNx (_DecimalNx x); // M < N

```

```

_FloatM fMfmafN (_FloatN x, _FloatN y, _FloatN z); // M < N
_FloatM fMfmafNx (_FloatNx x, _FloatNx y, _FloatNx z); // M <= N
_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
_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

```

int totalorderfN (_FloatN x, _FloatN y);
int totalorderdN (_DecimalN x, _DecimalN y);

int totalordermagfN (_FloatN x, _FloatN y);
int totalordermagdN (_DecimalN x, _DecimalN y);

```

F.10.13 Payload functions

```

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

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

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 F.10.12.2 (see Part 1 of Technical Specification 18661), append to paragraph 2:

The `totalorderfN` and `totalordermagfN` functions are declared for each corresponding data-interchange type that the implementation provides.

12 Numeric conversion functions <stdlib.h>

This clause specifies functions to convert between character sequences and the data-interchange types 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 new functions that perform conversions like `snprintf_s` in <stdio.h>.

Suggested changes to C11:

After 7.22.1.4, add:

James W Thomas 3/13/13 5:50 PM

Formatted: Normal, Indent: Left: 0", First line: 0", Add space between paragraphs of the same style

James W Thomas 3/13/13 5:53 PM

Formatted: Font:(Default) Courier New, Bold

James W Thomas 3/13/13 5:54 PM

Formatted: Font:Italic

James W Thomas 3/13/13 5:55 PM

Formatted: Indent: Left: 0.25", First line: 0", Add space between paragraphs of the same style

James W Thomas 3/13/13 5:53 PM

Formatted: Font:(Default) Courier New, Bold

James W Thomas 3/13/13 5:54 PM

Formatted: Font:Italic

James W Thomas 3/14/13 8:16 PM

Deleted: nnnn

James W Thomas 3/11/13 7:08 PM

Deleted: 4

James W Thomas 3/14/13 8:16 PM

Formatted: Font:11 pt, Bold, Font color: Auto

7.22.1.5 The `strtofN`, `strtofNx`, `strtodN`, and `strtodNx` functions

Synopsis

```
[1] #define __STDC_WANT_IEC_18661_EXT3__
#include <stdlib.h>
_FloatN strtodN (const char * restrict nptr, char ** restrict
endptr);
_FloatNx strtodNx (const char * restrict nptr, char ** restrict
endptr);
_DecimalN strtodN (const char * restrict nptr, char ** restrict
endptr);
_DecimalNx strtodNx (const char * restrict nptr, char ** restrict
endptr);
```

James W Thomas 3/13/13 11:46 AM

Deleted: 00000

James W Thomas 3/13/13 11:46 AM

Formatted: Font color: Text 1

Description

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, specified in Part 2 of Technical Specification 18661, except they convert to the types `_DecimalN` and `_DecimalNx` respectively.

James W Thomas 3/4/13 6:55 PM

Deleted: this

Returns

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.

7.22.1.6 The `strfromfN`, `strfromfNx`, `strfromdN`, and `strfromdNx` functions

Synopsis

```
[1] #define __STDC_WANT_IEC_18661_EXT3__
#include <stdlib.h>
int strfromfN (char * restrict s, rsize_t n, const char * restrict
format, _FloatN fp);
int strfromfNx (char * restrict s, rsize_t n, const char * restrict
format, _FloatNx fp);
int strfromdN (char * restrict s, rsize_t n, const char * restrict
format, _DecimalN fp);
int strfromdNx (char * restrict s, rsize_t n, const char * restrict
format, _DecimalNx fp);
```

James W Thomas 3/13/13 11:47 AM

Deleted: 00000

James W Thomas 3/13/13 11:47 AM

Formatted: Font color: Text 1

Description

These functions are equivalent to `snprintf_s(s, n, format, fp)` (K.3.5.3.5), except the type is implied by the function suffix and the `format` string contains no length modifier.

Returns

These functions return the value that would be returned by `snprintf_s(s, n, format, fp)`.

13 Complex arithmetic <complex.h>

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

Suggested changes to C11:

Change 7.3.1#3 from:

[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.

James W Thomas 3/14/13 9:27 AM

Formatted: std text, No bullets or numbering

James W Thomas 3/14/13 9:29 AM

Formatted

... [37]

to:

[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`, `l`, `fN`, and `fNx` suffixes which are corresponding functions whose parameters and return values have corresponding real types `float`, `long double`, `FloatN`, and `FloatNx`.

James W Thomas 3/14/13 9:28 AM

Formatted: std text, No bullets or numbering

James W Thomas 3/14/13 9:32 AM

Formatted

... [38]

James W Thomas 3/14/13 9:33 AM

Formatted

... [39]

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.

James W Thomas 3/14/13 9:29 AM

Formatted: std text

James W Thomas 3/14/13 3:15 PM

Formatted

... [40]

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

7.3.5 Trigonometric functions

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

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

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

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

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

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

7.3.6 Hyperbolic functions

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

```
_FloatN complex casinhfN (_FloatN complex z);
```



```

_FloatNx complex casinhfNx(_FloatNx complex z);
_FloatN complex catanhfN (_FloatN complex z);
_FloatNx complex catanhfNx(_FloatNx complex z);
5
_FloatN complex ccoshfN (_FloatN complex z);
_FloatNx complex acoshfNx(_FloatNx complex z);
_FloatN complex csinhfN (_FloatN complex z);
_FloatNx complex csinhfNx(_FloatNx complex z);
10
_FloatN complex ctanhfN (_FloatN complex z);
_FloatNx complex ctanhfNx(_FloatNx complex z);

```

7.3.7 Exponential and logarithmic functions

```

15
_FloatN complex cexpfN (_FloatN complex z);
_FloatNx complex cexpfNx(_FloatNx complex z);
_FloatN complex clogfN (_FloatN complex z);
20
_FloatNx complex clogfNx(_FloatNx complex z);

```

7.3.8 Power and absolute value functions

```

_FloatN complex cabsfN (_FloatN complex z);
_FloatNx complex cabsfNx(_FloatNx complex z);
25
_FloatN complex cpowfN (_FloatN complex z, _FloatN complex y);
_FloatNx complex cpowfNx(_FloatNx complex z, _FloatNx complex y);
_FloatN complex csqrtfN (_FloatN complex z);
30
_FloatNx complex csqrtfNx(_FloatNx complex z);

```

7.3.9 Manipulation functions

```

_FloatN complex cargfN (_FloatN complex z);
_FloatNx complex cargfNx(_FloatNx complex z);
35
_FloatN complex cimagfN (_FloatN complex z);
_FloatNx complex cimagfNx(_FloatNx complex z);
_FloatN complex CMPLXfN (_FloatN x, _FloatN y);
40
_FloatNx complex CMPLXfNx(_FloatNx x, _FloatNx y);
_FloatN complex conjfN (_FloatN complex z);
_FloatNx complex conjfNx(_FloatNx complex z);
45
_FloatN complex cprojfN (_FloatN complex z);
_FloatNx complex cprojfNx(_FloatNx complex z);
_FloatN crealfN (_FloatN complex z);
_FloatNx crealfNx(_FloatNx complex z);

```

James W Thomas 3/11/13 7:08 PM

Deleted: 4...Nnnnn

... [41]

James W Thomas 3/14/13 8:16 PM

Formatted: Font:11 pt, Bold, Font color: Auto

14 Type-generic macros <tgmath.h>

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

Suggested changes to C11:

In 7.25, replace paragraphs 2 and 3:

[2] Of the <math.h> and <complex.h> functions without an `f(float)` or `l(long double)` suffix, several have one or more parameters whose corresponding real type is `double`. For each such function, except `modf`, there is a corresponding *type-generic macro*.³¹³ The parameters whose corresponding real type is `double` in the function synopsis are *generic parameters*. Use of the macro invokes a function whose corresponding real type and type domain are determined by the arguments for the generic parameters.³¹⁴

[3] Use of the macro invokes a function whose generic parameters have the corresponding real type determined as follows:

- First, if any argument for generic parameters has type `long double`, the type determined is `long double`.
- Otherwise, if any argument for generic parameters has type `double` or is of integer type, the type determined is `double`.
- Otherwise, the type determined is `float`.

with:

[2] This clause specifies a many-to-one correspondence of functions in <math.h> and <complex.h> with a *type-generic macro*.³¹³ Use of the type-generic macro invokes a corresponding function whose type is determined by the types of the arguments for particular parameters called the *generic parameters*.³¹⁴

[3] Of the <math.h> and <complex.h> functions without a type suffix, several have one or more parameters whose corresponding real type is `double`. For each such function, except `modf`, there is a corresponding type-generic macro.³¹³ The parameters whose corresponding real type is `double` in the function synopsis are generic parameters.

[3a] Some of the <math.h> functions for decimal floating types have no unsuffixed counterpart. Of these functions with a `d64` suffix, some have one or more parameters whose type is `Decimal64`. For each such function, except `encodedecd64` and `encodebind64`, there is a corresponding type-generic macro. The parameters whose real type is `Decimal64` in the function synopsis are generic parameters.

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

[3c] Use of a type-generic macro invokes a function whose generic parameters have the corresponding real type determined by the corresponding real types of the arguments as follows:

- 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.
- 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`

James W Thomas 3/14/13 9:05 PM

Formatted: Heading 1

James W Thomas 3/14/13 3:17 PM

Formatted

... [42]

James W Thomas 3/14/13 3:17 PM

Formatted

... [43]

James W Thomas 3/14/13 3:18 PM

Formatted: Font:(Default) Courier New, 10 pt

James W Thomas 3/14/13 3:18 PM

Formatted

... [44]

James W Thomas 3/14/13 3:19 PM

Formatted

... [45]

James W Thomas 3/14/13 3:20 PM

Formatted

... [46]

James W Thomas 3/14/13 3:22 PM

Formatted

... [47]

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`.

5 — Otherwise, if any argument for generic parameters has type `_Decimal132`, the type determined is `_Decimal132`.

10 — Otherwise, if 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 types of these arguments. If `_FloatM` and either `long double` or `_FloatNx` are both widest types (with equivalent sets of values) of these arguments, the type determined is `_FloatM`. Otherwise, if `long double` and `_FloatNx` are both widest types (with equivalent sets of values) of these arguments, the type determined is `long double`.

15 — Otherwise, if any argument for generic parameters has type `double`, `_Float64`, or `_Float32x`, the type determined is the widest of the types of these arguments. If `_Float64` and either `double` or `_Float32x` are both widest types (with equivalent sets of values) of these arguments, the type determined is `_Float64`. Otherwise, if `double` and `_Float32x` are both widest types (with equivalent sets of values) of these arguments, the type determined is `double`.

— Otherwise, if any argument for generic parameters is of integer type, the type determined is `double`.

20 — Otherwise, if any argument for generic parameters has type `_Float32`, the type determined is `_Float32`.

— Otherwise, the type determined is `float`.

If neither `<math.h>` nor `<complex.h>` define a function whose generic parameters have the determined corresponding real type, the behavior is undefined.

25 In 7.25#5, replace the last sentence:

If all arguments for generic parameters are real, then use of the macro invokes a real function; otherwise, use of the macro results in undefined behavior.

with:

30 If all arguments for generic parameters are real, then use of the macro invokes a real function (provided `<math.h>` defines a function of the determined type); otherwise, use of the macro results in undefined behavior.

In 7.25#6, replace the last sentence:

Use of the macro with any real or complex argument invokes a complex function.

with:

35 Use of the macro with any argument of generic floating type, binary floating type, or complex type, invokes a complex function. Use of the macro with an argument of a decimal floating type results in undefined behavior.

After 7.25#6, add:

James W Thomas 3/14/13 8:16 PM
Formatted: Font:11 pt, Bold, Font color: Auto

James W Thomas 3/14/13 8:16 PM
Deleted: nnnn

James W Thomas 3/11/13 7:08 PM
Deleted: 4

James W Thomas 3/14/13 3:26 PM
Formatted: Font:(Default) Courier New, 10 pt, Not Bold

James W Thomas 3/14/13 3:26 PM
Formatted: Font:(Default) Courier New, 10 pt

James W Thomas 3/14/13 3:27 PM
Formatted: Font:(Default) Courier New, 10 pt

James W Thomas 3/14/13 3:27 PM
Formatted: Font:(Default) Courier New, 10 pt

James W Thomas 3/14/13 3:27 PM
Formatted: Font:(Default) Courier New, 10 pt

James W Thomas 3/14/13 3:28 PM
Formatted: Font:(Default) Courier New, 10 pt

James W Thomas 3/14/13 3:28 PM
Formatted: Font:(Default) Courier New, 10 pt

James W Thomas 3/14/13 3:28 PM
Formatted: Font:(Default) Courier New, 10 pt

James W Thomas 3/14/13 3:29 PM
Formatted: Font:10 pt, Not Bold

James W Thomas 3/14/13 3:29 PM
Formatted: Font:10 pt, Not Bold

James W Thomas 3/14/13 3:29 PM
Formatted: Font:(Default) Courier New, 10 pt

James W Thomas 3/14/13 3:29 PM
Formatted: Font:(Default) Courier New, 10 pt, Not Bold, Italic

James W Thomas 3/14/13 3:29 PM
Formatted: ... [48]

James W Thomas 3/14/13 3:31 PM
Formatted: Font:10 pt, Not Bold, Italic

James W Thomas 3/14/13 3:32 PM
Formatted: ... [49]

James W Thomas 3/14/13 3:33 PM
Formatted: ... [50]

James W Thomas 3/14/13 3:33 PM
Formatted: Font:10 pt, Not Bold

James W Thomas 3/14/13 3:34 PM
Formatted: ... [51]

James W Thomas 3/14/13 3:35 PM
Formatted: ... [52]

James W Thomas 3/14/13 3:35 PM
Formatted: ... [53]

James W Thomas 3/11/13 7:08 PM

Deleted: 4

James W Thomas 3/14/13 8:16 PM

Deleted: Nnnnn

James W Thomas 3/14/13 8:16 PM

Formatted: Font:11 pt, Bold, Font color: Auto

[6a] For each `d64`-suffixed function in `<math.h>` (except `encodedecd64` and `encodebind64`) that does not have an unsuffixed counterpart, the corresponding type-generic macro has the name of the function, but without the suffix. These type-generic macros are:

<code><math.h></code> function	type-generic macro
----- <code>quantized64</code> <code>samequantumd64</code> <code>quantexpd64</code>	----- <code>quantize</code> <code>samequantum</code> <code>quantexp</code>

Use of the macro with a generic floating or complex argument or with only integer type arguments results in undefined behavior.

[6b] For an implementation that supports Parts 1 and 3 (but not Part 2) of Technical Specification 18661 and that provides the following types:

type	IEC 60559 format
----- <code>float</code>	----- <code>binary32</code>
<code>double</code>	<code>binary64</code>
<code>long double</code>	<code>binary128</code>
<code>_Float32</code>	<code>binary32</code>
<code>_Float64</code>	<code>binary64</code>
<code>_Float128</code>	<code>binary128</code>
<code>_Float32x</code>	<code>binary64</code>
<code>_Float64x</code>	<code>binary128</code>

a type-generic macro `cbrt` that conforms to the specification in this clause and that is affected by constant rounding modes as specified in Part 1 of Technical Specification 18661 could be implemented as follows:

```
#if defined(__STDC_WANT_IEC_18661_EXT3)
#define cbrt(X) _Generic((X), \
    _Float128: cbrtf128(X), \
    _Float64: cbrtf64(X), \
    _Float32: cbrtf32(X), \
    _Float64x: cbrtf64x(X), \
    _Float32x: cbrtf32x(X), \
    long double: cbrtl(X), \
    default: _Roundwise_cbrt(X), \
    float: cbrtf(X) \
)
#else
#define cbrt(X) _Generic((X), \
    long double: cbrtl(X), \
    default: _Roundwise_cbrt(X), \
    float: cbrtf(X) \
)
#endif
```

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

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

```
#define __STDC_WANT_IEC_18661_EXT3__
```

In 7.25#7, append to the declarations:

James W Thomas 3/14/13 3:38 PM

Formatted: Font:(Default) Courier New

James W Thomas 3/14/13 2:08 PM

Formatted: No Spacing, Indent: Left: 0.28", First line: 0.28"

James W Thomas 3/14/13 3:42 PM

Formatted: Font:Not Bold

James W Thomas 3/14/13 3:38 PM

Formatted: Font:(Default) Courier New

James W Thomas 3/14/13 3:43 PM

Formatted: Font:Not Bold

James W Thomas 3/14/13 3:38 PM

Formatted: Font:(Default) Courier New

James W Thomas 3/14/13 3:43 PM

Formatted: Font:Not Bold

James W Thomas 3/14/13 3:39 PM

Formatted: Font:(Default) Courier New

James W Thomas 3/14/13 3:43 PM

Formatted: Font:Not Bold

James W Thomas 3/14/13 3:39 PM

Formatted: Font:(Default) Courier New

James W Thomas 3/14/13 3:43 PM

Formatted: Font:Not Bold

James W Thomas 3/14/13 3:39 PM

Formatted: Font:(Default) Courier New

James W Thomas 3/14/13 3:43 PM

Formatted: Font:Not Bold

James W Thomas 3/14/13 3:39 PM

Formatted: Font:(Default) Courier New

James W Thomas 3/14/13 3:43 PM

Formatted: Font:Not Bold

James W Thomas 3/14/13 3:39 PM

Formatted: Font:(Default) Courier New

James W Thomas 3/14/13 3:43 PM

Formatted: Font:Not Bold

```

#if __STDC_IEC_60559_TYPES__ >= 201ymmL
_Float32x f32x;
_Float64 f64;
_Float128 f128;
_Float64 complex f64c;
#endif

```

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>
<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 (the sets are not equivalent)
<code>pow(f32x, n)</code>	<code>powf32x(f32x, n)</code>

- James W Thomas 3/14/13 8:16 PM
Deleted: nnnn
- James W Thomas 3/14/13 8:16 PM
Formatted: Font:11 pt, Bold, Font color: Auto
- James W Thomas 3/11/13 7:08 PM
Deleted: 4

- James W Thomas 3/14/13 3:51 PM
Formatted: Font:(Default) Arial
- James W Thomas 3/14/13 3:51 PM
Formatted: Font:(Default) Arial
- James W Thomas 3/14/13 3:56 PM
Formatted: Font:(Default) Arial
- James W Thomas 3/14/13 3:52 PM
Formatted: Font:(Default) Arial
- James W Thomas 3/14/13 3:52 PM
Formatted: Font:(Default) Arial
- James W Thomas 3/14/13 3:52 PM
Formatted: Font:(Default) Arial
- James W Thomas 3/14/13 3:52 PM
Formatted: Font:(Default) Arial
- James W Thomas 3/14/13 3:52 PM
Formatted: Font:(Default) Arial
- James W Thomas 3/14/13 3:52 PM
Formatted: Font:(Default) Arial
- James W Thomas 3/14/13 3:53 PM
Formatted: Font:(Default) Arial
- James W Thomas 3/14/13 10:11 AM
Formatted: Font:(Default) Courier New, Bold, English (US)
- James W Thomas 3/14/13 3:53 PM
Formatted: std text, Don't add space between paragraphs of the same style

James W Thomas 3/11/13 7:08 PM

Deleted: 4

James W Thomas 3/14/13 8:16 PM

Deleted: Nnnnn

James W Thomas 3/14/13 8:16 PM

Formatted: Font:11 pt, Bold, Font color: Auto

Bibliography

- [1] ISO/IEC 9899:2011, *Information technology — Programming languages, their environments and system software interfaces — Programming Language C*
- [2] ISO/IEC/IEEE 60559:2011, *Information technology — Microprocessor Systems — Floating-point arithmetic*
- [3] 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*
- [4] IEC 60559:1989, *Binary floating-point arithmetic for microprocessor systems, second edition*
- [5] IEEE 754-2008, *IEEE Standard for Floating-Point Arithmetic*
- [6] IEEE 754-1985, *IEEE Standard for Binary Floating-Point Arithmetic*
- [7] IEEE 854-1987, *IEEE Standard for Radix-Independent Floating-Point Arithmetic*

James W Thomas 3/4/13 6:55 PM

Deleted: 2008