

Programming languages, their

interfaces — Floating-point

environments, and system software

FINAL DRAFT Technical Specification

ISO/IEC DTS 18661-4

ISO/IEC JTC 1/SC 22

Secretariat: ANSI

Voting begins on: **2024-12-06**

extensions for C — Part 4: Supplementary functions

ISO/IEC DTS 18661-4

https://standards.iteh.ai/catalog/standards/iso/498b4ed2-60a1-434a-93f7-57e92230b09c/iso-iec-dts-18661-4

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.

IN ADDITION TO THEIR EVALUATION AS BEING ACCEPTABLE FOR INDUSTRIAL, TECHNO-LOGICAL, COMMERCIAL AND USER PURPOSES, DRAFT INTERNATIONAL STANDARDS MAY ON OCCASION HAVE TO BE CONSIDERED IN THE LIGHT OF THEIR POTENTIAL TO BECOME STANDARDS TO WHICH REFERENCE MAY BE MADE IN NATIONAL REGULATIONS.

iTeh Standards (https://standards.iteh.ai) Document Preview

ISO/IEC DTS 18661-4

https://standards.iteh.ai/catalog/standards/iso/498b4ed2-60a1-434a-93f7-57e92230b09c/iso-iec-dts-18661-4



© ISO/IEC 2024

All rights reserved. Unless otherwise specified, or required in the context of its implementation, no part of this publication may be reproduced or utilized otherwise in any form or by any means, electronic or mechanical, including photocopying, or posting on the internet or an intranet, without prior written permission. Permission can be requested from either ISO at the address below or ISO's member body in the country of the requester.

ISO copyright office CP 401 • Ch. de Blandonnet 8 CH-1214 Vernier, Geneva Phone: +41 22 749 01 11 Email: copyright@iso.org Website: www.iso.org Published in Switzerland

ISO/IEC DTS 18661-4:2024(en)

Contents

Foreword			
Introduction			v
1	Scone	٥	1
1	New		
Z	ΝΟΓΠ	lative references	I
3	Term	is and definitions	
4	Confo	Conformance	
5	C standard extensions		
	5.1	Predefined macros	
	5.2	Freestanding implementations	2
	5.3	Headers	2
	5.4	Future directions	2
6	Reduction functions <reduc.h></reduc.h>		2
	6.1	General	
	6.2	The reduc sum functions	
	6.3	The reduc sumabs functions	4
	6.4	The reduc sumsq functions	4
	6.5	The reduc_sumprod functions	5
	6.6	The scaled_prod functions	5
	6.7	The scaled_prodsum functions	6
	6.8	The scaled_proddiff functions	7
7	Augmented arithmetic functions <augarith.h></augarith.h>		
	7.1	General	9
	7.2	The aug add functions 3.4/SUALIDEALUS.IUCIL.AL	
	7.3	The aug sub functions	
	7.4	The aug_mul functions ocument Preview	
Bibliography			13
Dibuogi ahui			10

SO/IEC DTS 18661-4

https://standards.iteh.ai/catalog/standards/iso/498b4ed2-60a1-434a-93f7-57e92230b09c/iso-iec-dts-18661-4

Foreword

ISO (the International Organization for Standardization) and IEC (the International Electrotechnical Commission) form the specialized system for worldwide standardization. National bodies that are members of ISO or IEC participate in the development of International Standards through technical committees established by the respective organization to deal with particular fields of technical activity. ISO and IEC technical committees collaborate in fields of mutual interest. Other international organizations, governmental and non-governmental, in liaison with ISO and IEC, also take part in the work.

The procedures used to develop this document and those intended for its further maintenance are described in the ISO/IEC Directives, Part 1. In particular, the different approval criteria needed for the different types of document should be noted. This document was drafted in accordance with the editorial rules of the ISO/IEC Directives, Part 2 (see www.iso.org/directives or www.iso.org/directiv

ISO and IEC draw attention to the possibility that the implementation of this document may involve the use of (a) patent(s). ISO and IEC take no position concerning the evidence, validity or applicability of any claimed patent rights in respect thereof. As of the date of publication of this document, ISO and IEC had not received notice of (a) patent(s) which may be required to implement this document. However, implementers are cautioned that this may not represent the latest information, which may be obtained from the patent database available at www.iso.org/patents and https://patents.iec.ch. ISO and IEC shall not be held responsible for identifying any or all such patent rights.

Any trade name used in this document is information given for the convenience of users and does not constitute an endorsement.

For an explanation of the voluntary nature of standards, the meaning of ISO specific terms and expressions related to conformity assessment, as well as information about ISO's adherence to the World Trade Organization (WTO) principles in the Technical Barriers to Trade (TBT) see www.iso.org/iso/foreword.html. In the IEC, see www.iso.org/iso/foreword.html.

This document was prepared by Joint Technical Committee ISO/IEC JTC 1, *Information technology*, Subcommittee SC 22, *Programming languages, their environments and system software interfaces*.

This second edition cancels and replaces the first edition (ISO/IEC TS 18661-4:2015), which has been technically revised.

https://standards.iteh.ai/catalog/standards/iso/498b4ed2-60a1-434a-93f7-57e92230b09c/iso-iec-dts-18661-4 The main changes are as follows:

- The specification has been updated to extend ISO/IEC 9899:2024.
- The mathematical functions and constant rounding modes have been removed. These features are now incorporated into ISO/IEC 9899:2024.
- Functions to support the augmented arithmetic operations specified in IEEE 754-2019 have been added.
- New headers have been added, and all extensions to the <math.h> header have been removed.

A list of all parts in the ISO 18661 series can be found on the ISO website.

ISO/IEC DTS 18661-4:2024(en)

Introduction

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 reliable programs, debugging and moving programs between systems exceedingly difficult. Now the great majority of systems provide data formats and arithmetic operations according to IEEE 754. Corresponding versions of IEEE 754 and ISO/IEC 60559 have equivalent content.

Support for IEEE 754-1985 was added in ISO/IEC 9899:1999 (also referred to as C99), and ISO/IEC 9899:2018 is still based on IEEE 754-1985. However, IEEE 754 underwent a major revision in 2008 and a minor revision in 2019, which added several new features.

The purpose of the ISO/IEC TS 18661 series (first published 2014 through 2016) has been to specify C language support for the new features introduced into IEEE 754 since 1985. Most of the ISO/IEC TS 18661 series has been incorporated into ISO/IEC 9899:2024 (also referred to as C23 because major work on this revision was completed in 2023), which supports all required and most recommended features in IEEE 754-2019.

To supplement the IEEE 754 support in C23, this document specifies two C headers with functions corresponding to the reduction and augmented arithmetic operations recommended by IEEE 754, but not included in C23.

The reduction operations perform widely used vector computations involving sums and products, including scaled products. These operations are allowed to associate in any order, and to evaluate in any wider format.

The augmented arithmetic operations, added in IEEE 754-2019, are versions of operations commonly called twoSum and twoProduct. These operations can be used to implement arithmetic with extra precision, for example, for double-double format. In theory, they can also be used to implement efficient reproducible dot products.

(https://standards.iteh.ai) Document Preview

ISO/IEC DTS 18661-4

https://standards.iteh.ai/catalog/standards/iso/498b4ed2-60a1-434a-93f7-57e92230b09c/iso-iec-dts-18661-4

iTeh Standards (https://standards.iteh.ai) Document Preview

ISO/IEC DTS 18661-4

https://standards.iteh.ai/catalog/standards/iso/498b4ed2-60a1-434a-93f7-57e92230b09c/iso-iec-dts-18661-4

Programming languages, their environments, and system software interfaces — Floating-point extensions for C —

Part 4: **Supplementary functions**

1 Scope

This document specifies extensions to programming language C to include functions corresponding to operations specified and recommended in ISO/IEC 60559, but not supported in ISO/IEC 9899:2024 (also referred to as C23).

2 Normative references

The following documents, in whole or in part, are normatively referenced in this document and are indispensable for its application. For dated references, only the edition cited applies. For undated references, the latest edition of the referenced document (including any amendments) applies.

ISO/IEC 9899:2024, Information technology — Programming languages — C

ISO/IEC 60559:2020, Information technology — Microprocessor Systems — Floating-point arithmetic

3 Terms and definitions **Document Preview**

For the purposes of this document, the terms and definitions given in ISO/IEC 9899:2024 and ISO/IEC 60559:2020 apply.

ISO and IEC maintain terminological databases for use in standardization at the following addresses:

- ISO Online browsing platform: available at <u>https://www.iso.org/obp</u>
- IEC Electropedia: available at <u>https://www.electropedia.org/</u>

4 Conformance

An implementation that meets the requirements for a conforming implementation of C23 may conform to one or both feature sets (reduction functions and augmented arithmetic) in this document. The implementation conforms to the reduction functions feature if:

- a) it defines __STDC_IEC_60559_BFP__ or __STDC_IEC_60559_DFP__ or both, indicating support for ISO/IEC 60559 binary or decimal floating-point arithmetic, as specified in C23, Annex F;
- b) it defines __STDC_IEC_60559_FUNCS_REDUCTION__ to 202401L and provides the <reduc.h> header specified in this document (<u>Clause 6</u>).

The implementation conforms to the augmented arithmetic feature if:

- c) it defines __stdc_iec_60559_bfp_, indicating support for ISO/IEC 60559 binary floating-point arithmetic, as specified in C23, Annex F;
- d) it defines __stDc_iec_60559_funcs_AUGMENTED_ARITHMETIC__ to 202401L and provides the <augarith.h> header specified in this document (Clause 7).

5 C standard extensions

5.1 Predefined macros

The implementation defines one or both of the following macros to indicate conformance to the specification in this document for support of the corresponding features specified and recommended in ISO/IEC 60559.

___STDC_IEC_60559_FUNCS_REDUCTION___ The integer constant 202401L. ___STDC_IEC_60559_FUNCS_AUGMENTED_ARITHMETIC___ The integer constant 202401L.

5.2 Freestanding implementations

The strictly conforming programs that shall be accepted by a conforming freestanding implementation that defines one of the feature macros in 5.1 may also use features in the corresponding header specified in this document. See C23, Clause 4.

5.3 Headers

If the implementation defines one of the feature macros in 5.1 then the implementation provides the corresponding header specified in this document. The header and its use follow the general specification in C23, 7.1 for the C Library as though the header were a subclause of C23, Clause 7 for a conditional feature.

5.4 Future directions

For implementations that define __stDc_iec_60559_funcs_Reduction__, function names that begin with <code>reduc_</code> or <code>scaled_</code> are potentially reserved identifiers and may be added to the declarations in the <code><reduc.h></code> header.

For implementations that define __STDC_IEC_60559_FUNCS_AUGMENTED_ARITHMETIC__, tag names that end with <code>aug_t</code> and function names that begin with <code>aug_</code> are potentially reserved identifiers and may be added to the declarations in the <code><augarith.h></code> header.

See C23, 7.33. <u>ISO/IEC DTS 18661-4</u> https://standards.iteh.ai/catalog/standards/iso/498b4ed2-60a1-434a-93f7-57e92230b09c/iso-iec-dts-18661-4

6 Reduction functions <reduc.h>

6.1 General

The header <reduc.h> declares the type and functions in this clause.

The type declared is <code>size_t</code> (described in C23, 7.21.1).

Each function in this clause is declared in <reduc.h> if and only if the corresponding type is supported according to C23, Annex F or Annex H.

The functions in this clause shall be implemented so that intermediate computations do not overflow or underflow. For the <code>reduc_sum</code>, <code>reduc_sum</code>,

The reduction functions do not raise the "divide-by-zero" floating-point exception.

With ISO/IEC 60559 default exception handling, these functions raise the "inexact" floating-point exception in response to "overflow" and "underflow" exceptions; otherwise, whether they raise the "inexact" floating-point exception is unspecified.

ISO/IEC DTS 18661-4:2024(en)

Numerical results and exceptional behavior, including the "invalid" floating-point exception, can differ due to the precision of intermediates and the order of evaluation. However, only one floating-point exception is raised (other than "inexact" in response to "overflow" or "underflow") per reduction function invocation; exceptions are not raised for each exceptional intermediate operand or result. Reduction functions may raise the "invalid" floating-point exception if an element of an array argument is a signaling NaN (see C23, F.2.2). Once an invalid floating-point exception is raised, due to signaling NaN, $\infty - \infty$, or $0 \times \infty$, processing of array elements may stop.

Whether and how rounding direction modes affect functions in this clause are implementation defined and may be indeterminate. This applies to constant as well as dynamic rounding modes, C23, 7.6.3 notwithstanding.

The preferred quantum exponent for the reduction functions for decimal floating types is unspecified.

NOTE For N = 32, 64 and 128, the functions suffixed with dN are declared if the implementation supports decimal floating types (i.e. defines __STDC_IEC_60559_DFP_), without the requirement that the macro __STDC_WANT_IEC_60559_TYPES_EXT_ be defined.

6.2 The reduc_sum functions

```
Synopsis

#include <reduc.h> https://standards.iteh.ai)

#ifdef _STDC_IEC_60559_BFP_
double reduc_sum(size_t n, const double p[static n]);
float reduc_sum(size_t n, const float p[static n]);
long double reduc_suml(size_t n, const long double p[static n]);
_FloatN reduc_sumfN(size_t n, const _FloatN p[static n]);
_FloatNx reduc_sumfN(size_t n, const _FloatNx p[static n]);
#endif
#ifdef _STDC_IEC_60559_DFP_
_DecimalN reduc_sumdN(size_t n, const _DecimalN p[static n]);
_DecimalNx reduc_sumdN(size_t n, const _DecimalNx p[static n]);
#endif
```

Description

The reduc_sum functions compute the sum of the n elements of array p: $\sum_{i=0}^{n-1} p[i]$. If the length n = 0, the

functions return the value +0. If any element of array $_{\rm p}$ is a NaN, the functions return a quiet NaN. If any two elements of array $_{\rm p}$ are infinities with different signs, the functions return a quiet NaN and raise the "invalid" floating-point exception and a domain error occurs. Otherwise (if no element of $_{\rm p}$ is a NaN and no two elements of $_{\rm p}$ are infinities with different signs), if any element of array $_{\rm p}$ is an infinity, the functions return that same infinity.

Returns

The reduc_sum functions return the computed sum.

6.3 The reduc_sumabs functions

Synopsis

```
#include <reduc.h>
```

```
#ifdef __STDC_IEC_60559_BFP__
double reduc_sumabs(size_t n, const double p[static n]);
float reduc_sumabsf(size_t n, const float p[static n]);
long double reduc_sumabsl(size_t n, const long double p[static n]);
_FloatN reduc_sumabsfN(size_t n, const _FloatN p[static n]);
_FloatNx reduc_sumabsfNx(size_t n, const _FloatNx p[static n]);
#endif
#ifdef __STDC_IEC_60559_DFP____
_DecimalN reduc_sumabsdN(size_t n, const _DecimalN p[static n]);
_DecimalNx reduc_sumabsdNx(size_t n, const _DecimalNx p[static n]);
#endif
```

Description

The reduc_sumabs functions compute the sum of the absolute values of the n elements of array p: $\sum_{i=0}^{n-1} |p[i]|$. If the length n = 0, the functions return the value +0. If any element of array p is an infinity, the functions return + ∞ ; otherwise, if any element of array p is a NaN, the functions return a quiet NaN.

Returns

The reduc_sumabs functions return the computed sum.

6.4 The reduc sumsq functions

Synopsis

```
#include <reduc.h> Document Preview
#ifdef __STDC_IEC_60559_BFP__
double reduc_sumsq(size_t n, const double p[static n]);
float reduc_sumsqf(size_t n, const float p[static n]);
https://long_double_reduc_sumsql(size_t n, const_long_double_p[static n]);
https://long_double_reduc_sumsqfN(size_t n, const_long_double_p[static n]);
FloatN reduc_sumsqfNx(size_t n, const_FloatN p[static n]);
FloatNx_reduc_sumsqfNx(size_t n, const_FloatNx p[static n]);
#endif
#ifdef __STDC_IEC_60559_DFP__
__DecimalN_reduc_sumsqdN(size_t n, const_DecimalN p[static n]);
__DecimalNx_reduc_sumsqdNx(size_t n, const_DecimalNx p[static n]);
#endif
```

Description

The reduc_sumsq functions compute the sum of squares of the values of the n elements of array p: $\sum_{i=0}^{n-1} (p[i] \times p[i])$ If the length n = 0, the functions return the value +0. If any element of array p is an infinity, the functions return + ∞ ; otherwise, if any element of array p is a NaN, the functions return a quiet NaN.

Returns

The ${\tt reduc_sumsq}$ functions return the computed sum.