



Technical  
Specification

**ISO/IEC TS 18661-4**

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

Part 4:  
**Supplementary functions**

*Langages de programmation, leurs environnements et interfaces  
du logiciel système — Extensions à virgule flottante pour C —*

*Partie 4: Fonctions supplémentaires*

Second edition  
2025-03

IECNORM.COM : Click to view the full PDF of ISO/IEC TS 18661-4:2025



**COPYRIGHT PROTECTED DOCUMENT**

© ISO/IEC 2025

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](mailto:copyright@iso.org)  
Website: [www.iso.org](http://www.iso.org)

Published in Switzerland

# Contents

	Page
<b>Foreword</b> .....	<b>iv</b>
<b>Introduction</b> .....	<b>v</b>
<b>1 Scope</b> .....	<b>1</b>
<b>2 Normative references</b> .....	<b>1</b>
<b>3 Terms and definitions</b> .....	<b>1</b>
<b>4 Conformance</b> .....	<b>1</b>
<b>5 C standard extensions</b> .....	<b>2</b>
5.1 Predefined macros.....	2
5.2 Freestanding implementations.....	2
5.3 Headers.....	2
5.4 Future directions.....	2
<b>6 Reduction functions</b> <reduc.h>.....	<b>2</b>
6.1 General.....	2
6.2 The <code>reduc_sum</code> functions.....	3
6.3 The <code>reduc_sumabs</code> functions.....	4
6.4 The <code>reduc_sumsq</code> functions.....	4
6.5 The <code>reduc_sumprod</code> functions.....	5
6.6 The <code>scaled_prod</code> functions.....	5
6.7 The <code>scaled_prodsum</code> functions.....	6
6.8 The <code>scaled_proddiff</code> functions.....	7
<b>7 Augmented arithmetic functions</b> <augarith.h>.....	<b>9</b>
7.1 General.....	9
7.2 The <code>aug_add</code> functions.....	9
7.3 The <code>aug_sub</code> functions.....	11
7.4 The <code>aug_mul</code> functions.....	11
<b>Bibliography</b> .....	<b>13</b>

## 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](http://www.iso.org/directives) or [www.iec.ch/members\\_experts/refdocs](http://www.iec.ch/members_experts/refdocs)).

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](http://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](http://www.iso.org/iso/foreword.html). In the IEC, see [www.iec.ch/understanding-standards](http://www.iec.ch/understanding-standards).

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.

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.

Any feedback or questions on this document should be directed to the user's national standards body. A complete listing of these bodies can be found at [www.iso.org/members.html](http://www.iso.org/members.html) and [www.iec.ch/national-committees](http://www.iec.ch/national-committees).

## 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 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 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.

IECNORM.COM : Click to view the full PDF of ISO/IEC TS 18661-4:2025

IECNORM.COM : Click to view the full PDF of ISO/IEC TS 18661-4:2025

# 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

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 <https://www.iso.org/obp>
- IEC Electropedia: available at <https://www.electropedia.org/>

### 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 (Clause 6).

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 `reduc_` or `scaled_` are potentially reserved identifiers and may be added to the declarations in the `<reduc.h>` header.

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

See C23, 7.33.

## 6 Reduction functions `<reduc.h>`

### 6.1 General

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

The type declared is `size_t` (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 `reduc_sum`, `reduc_sumabs`, `reduc_sumsq`, and `reduc_sumprod` functions, the “overflow” or “underflow” floating-point exception is raised and a range error occurs, if and only if the final result overflows or underflows. The `scaled_prod`, `scaled_prodsum`, and `scaled_proddiff` functions do not raise the “overflow” or “underflow” floating-point exceptions and do not cause a range error.

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.

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.

For each of the following synopses, an implementation shall declare the functions suffixed with  $fN$  or  $fNx$  only if it supports the corresponding binary floating type and the macro `__STDC_WANT_IEC_60559_TYPES_EXT__` is defined at the point in the code where `<reduc.h>` is first included. An implementation shall declare the functions suffixed with  $dN$  for  $N \neq 32, 64$  or  $128$  or with  $dNx$  only if it supports the corresponding decimal floating type and the macro `__STDC_WANT_IEC_60559_TYPES_EXT__` is defined at the point in the code where `<reduc.h>` is first included (see C23, Annex H.)

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>

#ifdef __STDC_IEC_60559_BFP__
double reduc_sum(size_t n, const double p[static n]);
float reduc_sumf(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_sumfNx(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_sumdNx(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  $p$  is a NaN, the functions return a quiet NaN. If any two elements of array  $p$  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  $p$  is a NaN and no two elements of  $p$  are infinities with different signs), if any element of array  $p$  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  $+\infty$ ; 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>

#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]);
long double reduc_sumsql(size_t n, const long double p[static n]);
_FloatN reduc_sumsqfN(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  $+\infty$ ; otherwise, if any element of array  $p$  is a NaN, the functions return a quiet NaN.

#### Returns

The `reduc_sumsq` functions return the computed sum.

## 6.5 The `reduc_sumprod` functions

### Synopsis

```
#include <reduc.h>

#ifdef __STDC_IEC_60559_BFP__
double reduc_sumprod(size_t n,
    const double p[static n],
    const double q[static n]);
float reduc_sumprodf(size_t n,
    const float p[static n],
    const float q[static n]);
long double reduc_sumprodl(size_t n,
    const long double p[static n],
    const long double q[static n]);
_FloatN reduc_sumprodfN(size_t n,
    const _FloatN p[static n],
    const _FloatN q[static n]);
_FloatNx reduc_sumprodfNx(size_t n,
    const _FloatNx p[static n],
    const _FloatNx q[static n]);
#endif
#ifdef __STDC_IEC_60559_DFP__
_DecimalN reduc_sumproddN(size_t n,
    const _DecimalN p[static n],
    const _DecimalN q[static n]);
_DecimalNx reduc_sumproddNx(size_t n,
    const _DecimalNx p[static n],
    const _DecimalNx q[static n]);
#endif
```

### Description

The `reduc_sumprod` functions compute the dot product of the sequences of elements of the arrays `p` and `q`:  $\sum_{i=0}^{n-1} (p[i] \times q[i])$ . If the length `n = 0`, the functions return the value `+0`. If any element of array `p` or `q` is a NaN, the functions return a quiet NaN. If a product is  $0 \times \infty$ , the functions return a quiet NaN and raise the “invalid” floating-point exception and a domain error occurs. If a sum is of infinities of different signs, the functions return a quiet NaN and raise the “invalid” floating-point exception and a domain error occurs. Otherwise (if no array element is a NaN, no product is  $0 \times \infty$ , and no sum is of infinities of different signs), if a term in the summation is an infinity, the functions return that same infinity.

### Returns

The `reduc_sumprod` functions return the computed dot product.

## 6.6 The `scaled_prod` functions

### Synopsis

```
#include <reduc.h>

#ifdef __STDC_IEC_60559_BFP__
double scaled_prod(size_t n,
    const double p[static restrict n],
    long int * restrict sfptr);
float scaled_prodf(size_t n,
    const float p[static restrict n],
    long int * restrict sfptr);
long double scaled_prodl(size_t n,
    const long double p[static restrict n],
    long int * restrict sfptr);
_FloatN scaled_prodfN(size_t n,
    const _FloatN p[static restrict n],
    long int * restrict sfptr);
_FloatNx scaled_prodfNx(size_t n,
```

```

    const _FloatNx p[static restrict n],
    long int * restrict sfptr);
#endif
#ifdef __STDC_IEC_60559_DFP__
    _DecimalN scaled_proddN(size_t n,
    const _DecimalN p[static restrict n],
    long int * restrict sfptr);
    _DecimalNx scaled_proddNx(size_t n,
    const _DecimalNx p[static restrict n],
    long int * restrict sfptr);
#endif
#endif

```

## Description

The `scaled_prod` functions compute a scaled product  $pr$  of the  $n$  elements of the array  $p$  and a scale factor  $sf$ , such that

$$pr \times b^{sf} = \prod_{i=0}^{n-1} p[i]$$

where  $b$  is the radix of the type. These functions store the scale factor  $sf$  in the object pointed to by `sfptr`. If the length  $n = 0$ , the functions return the value +1 and store 0 in the object pointed to by `sfptr`. If any element of array  $p$  is a NaN, the functions return a quiet NaN. If any two elements of array  $p$  are a zero and an infinity, the functions return a quiet NaN and raise the “invalid” floating-point exception and a domain error occurs. Otherwise, if any element of array  $p$  is an infinity, the functions return an infinity. Otherwise, if any element of array  $p$  is a zero, the functions return a zero. Otherwise, if the scale factor is outside the range of the `long int` type, the functions return a quiet NaN and raise the “invalid” floating-point exception. If a zero, infinity, or NaN is returned, the functions store 0 in the object pointed to by `sfptr`.

## Returns

The `scaled_prod` functions return the computed scaled product  $pr$ .

## 6.7 The `scaled_proddsum` functions

### Synopsis

```

#include <reduc.h>

#ifdef __STDC_IEC_60559_BFP__
double scaled_proddsum(size_t n,
    const double p[static restrict n],
    const double q[static restrict n],
    long int * restrict sfptr);
float scaled_proddsumf(size_t n,
    const float p[static restrict n],
    const float q[static restrict n],
    long int * restrict sfptr);
long double scaled_proddsuml(size_t n,
    const long double p[static restrict n],
    const long double q[static restrict n],
    long int * restrict sfptr);
_FloatN scaled_proddsumfN(size_t n,
    const _FloatN p[static restrict n],
    const _FloatN q[static restrict n],
    long int * restrict sfptr);
_FloatNx scaled_proddsumfNx(size_t n,
    const _FloatNx p[static restrict n],
    const _FloatNx q[static restrict n],
    long int * restrict sfptr);
#endif
#ifdef __STDC_IEC_60559_DFP__
    _DecimalN scaled_proddsumdN(size_t n,
    const _DecimalN p[static restrict n],
    const _DecimalN q[static restrict n],
    long int * restrict sfptr);

```

```

_DecimalNx scaled_prodsimdNx(size_t n,
    const _DecimalNx p[static restrict n],
    const _DecimalNx q[static restrict n],
    long int * restrict sfptr);
#endif

```

## Description

The `scaled_prodsim` functions compute a scaled product  $pr$  of the sums of the corresponding elements of the arrays  $p$  and  $q$  and a scale factor  $sf$ , such that

$$pr \times b^{sf} = \prod_{i=0}^{n-1} (p[i] + q[i])$$

where  $b$  is the radix of the type. These functions store the scale factor  $sf$  in the object pointed to by `sfptr`. If the length  $n = 0$ , the functions return the value +1 and store 0 in the object pointed to by `sfptr`. If any element of array  $p$  or  $q$  is a NaN, the functions return a quiet NaN. If any sum is of infinities with different signs or if any two factors in the product are a zero and an infinity, the functions return a quiet NaN and raise the “invalid” floating-point exception and a domain error occurs. Otherwise, if any factor in the product is an infinity, the functions return an infinity. Otherwise, if any factor in the product is a zero, the functions return a zero. Otherwise, if the scale factor is outside the range of the `long int` type, the functions return a quiet NaN and raise the “invalid” floating-point exception. If a zero, infinity, or NaN is returned, the functions store 0 in the object pointed to by `sfptr`.

## Returns

The `scaled_prodsim` functions return the computed scaled product  $pr$ .

## 6.8 The `scaled_proddiff` functions

### Synopsis

```

#include <reduc.h>

#ifdef __STDC_IEC_60559_BFP__
double scaled_proddiff(size_t n,
    const double p[static restrict n],
    const double q[static restrict n],
    long int * restrict sfptr);
float scaled_proddiff(float p[static restrict n],
    const float q[static restrict n],
    long int * restrict sfptr);
long double scaled_proddiff(long double p[static restrict n],
    const long double q[static restrict n],
    long int * restrict sfptr);
_FloatN scaled_proddiffN(size_t n,
    const _FloatN p[static restrict n],
    const _FloatN q[static restrict n],
    long int * restrict sfptr);
_FloatNx scaled_proddiffNx(size_t n,
    const _FloatNx p[static restrict n],
    const _FloatNx q[static restrict n],
    long int * restrict sfptr);
#endif
#ifdef __STDC_IEC_60559_DFP__
_DecimalN scaled_proddiffdN(size_t n,
    const _DecimalN p[static restrict n],
    const _DecimalN q[static restrict n],
    long int * restrict sfptr);
_DecimalNx scaled_proddiffdNx(size_t n,
    const _DecimalNx p[static restrict n],
    const _DecimalNx q[static restrict n],
    long int * restrict sfptr);
#endif

```

## Description

The `scaled_proddiff` functions compute a scaled product  $pr$  of the differences of the corresponding elements of the arrays  $p$  and  $q$  and a scale factor  $sf$ , such that

$$pr \times b^{sf} = \prod_{i=0}^{n-1} (p[i] - q[i])$$

where  $b$  is the radix of the type. These functions store the scale factor  $sf$  in the object pointed to by `sfptr`. If the length  $n = 0$ , the functions return the value +1 and store 0 in the object pointed to by `sfptr`. If any element of array  $p$  or  $q$  is a NaN, the functions return a quiet NaN. If any difference is of infinities with the same signs or if any two factors in the product are a zero and an infinity, the functions return a quiet NaN and raise the “invalid” floating-point exception and a domain error occurs. Otherwise, if any factor in the product is an infinity, the functions return an infinity. Otherwise, if any factor in the product is a zero, the functions return a zero. Otherwise, if the scale factor is outside the range of the `long int` type, the functions return a quiet NaN and raise the “invalid” floating-point exception. If a zero, infinity, or NaN is returned, the functions store 0 in the object pointed to by `sfptr`.

## Returns

The `scaled_proddiff` functions return the computed scaled product  $pr$ .

**EXAMPLE** The scaled reduction functions support computing quantities of modest magnitudes whose intermediate results can overflow and underflow. One example is the computation of Clebsch-Gordan coefficients or Wigner 3-j symbols for quantum physics. Expressions for these quantities involve quotients of products of factorials, and so are prone to intermediate overflow. As a simplified example, consider computing a fragment of the Clebsch-Gordan calculation.

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

// compute quot = n1! * n2! / n3!

int n1 = 140, n2 = 160, n3 = 200; // factorial magnitudes
// 1e241, 1e284, 1e374
// quot magnitude 1e151

// products scaled to avoid overflow
double num1, num2, den;
// scale factors
long int num1e, num2e, dene;

// products scaled again to avoid intermediate overflow
// in final computation
double num1s, num2s, dens;
// scale factors
long int num1es, num2es, denes;

// result
double quot;

// arrays { 2, 3, 4, ... }
double num1p[n1-1], num2p[n2-1], denp[n3-1];

// n1! scaled to avoid overflow
for (int i = 2; i <= n1; i++) {
    num1p[i-2] = i;
}
num1 = scaled_prod(n1-1, num1p, &num1e);

// n2! scaled to avoid overflow
for (int i = 2; i <= n2; i++) {
    num2p[i-2] = i;
}
num2 = scaled_prod(n2-1, num2p, &num2e);
```

```

// n3! scaled to avoid overflow
for (int i = 2; i <= n3; i++) {
    denp[i-2] = i;
}
den = scaled_prod(n3-1, denp, &dene);

// re-scale to avoid subsequent intermediate overflow
num1es = llogb(num1);
num1s = scalbln(num1, -num1es);
num2es = llogb(num2);
num2s = scalbln(num2, -num2es);
denes = llogb(den);
dens = scalbln(den, -denes);

// compute result without intermediate overflow
quot = scalbln(num1s * num2s / dens,
    num1e + num2e - dene + num1es + num2es - denes);

```

## 7 Augmented arithmetic functions <augarith.h>

### 7.1 General

This clause supports augmented arithmetic, as recommended by ISO/IEC 60559 for its binary formats. Each type and each function in this clause are declared in <augarith.h> if and only if the corresponding type is an ISO/IEC 60559 floating type supported according to Annex F of C23 or an interchange or extended type supported according to Annex H of C23.

The types are structures for returning two floating-point values:

```

struct daug_t { double h; double t; };
struct faug_t { float h; float t; };
struct ldaug_t { long double h; long double t; };
struct _fNaug_t { _FloatN h; _FloatN t; };
struct _fNxaug_t { _FloatNx h; _FloatNx t; };

```

The functions in this clause use these structures to return a “head” value  $h$  and “tail” value  $t$  to represent an extra-precise result value given by  $h + t$ . See the example in [7.2](#).

The functions in this clause round to nearest with ties toward zero, a rounding direction specified by ISO/IEC 60559 for use by augmented arithmetic operations. Thus, results are independent of dynamic and constant rounding direction modes. Like other ISO/IEC 60559 operations, rounding is done with gradual underflow.

NOTE Reference [\[11\]](#) shows how to use currently available ISO/IEC 60559 operations and to-nearest, ties-to-even rounding to implement the augmented arithmetic operations with their special to-nearest, ties-toward-zero rounding.

For each of the following synopses, an implementation shall declare the functions suffixed with  $fN$  or  $fNx$  only if it supports the corresponding binary floating type and the macro `__STDC_WANT_IEC_60559_TYPES_EXT__` is defined at the point in the code where <augarith.h> is first included.

### 7.2 The `aug_add` functions

#### Synopsis

```

#include <augarith.h>

struct daug_t aug_add(double x, double y);
struct faug_t aug_addf(float x, float y);
struct ldaug_t aug_addl(long double x, long double y);
struct _fNaug_t aug_addfN(_FloatN x, _FloatN y);
struct _fNxaug_t aug_addfNx(_FloatNx x, _FloatNx y);

```

## Description

The `aug_add` functions compute two result values:

`h`: the sum  $x + y$  rounded to the type using round-to-nearest with ties toward zero;

`t`: the error in `h` as a computation of  $x + y$ .

If `h` is a non-zero finite number, `t` has the value  $x + y - h$  (which is exactly representable in the type), where if `t` is zero its sign is the sign of `h`. If `h` is zero, `t` has the value of `h` (and both have the same sign). If `h` is infinite, `t` has the value of `h`. If `h` is a NaN, `t` is the same NaN.

These functions raise floating-point exceptions like the computation of `h`, except that under ISO/IEC 60559 default exception handling they raise the “inexact” floating-point exception only when the computation of `h` overflows.

A range error occurs when the computation of `h` overflows. A domain error occurs when the arguments are infinities with different signs.

## Returns

These functions return the sum and error in a structure.

**EXAMPLE** The augmented arithmetic operations are useful for extending precision, particularly for implementing “double-double” arithmetic, which provides a faster, though less precise, and less predictable, alternative to ISO/IEC 60559 binary128 on systems which lack hardware support for binary128.

Double-double represents numbers as a pair of doubles, the second no larger in magnitude than the first, and usually much smaller, where the pair  $(h, t)$  represents the number  $h + t$  (to infinite precision). Ideally, in the pair  $(h, t)$ , `h` equals the correctly rounded result of computed  $(h+t)$ , and `t` equals the correctly rounded result of  $h + t - \text{correctly-rounded-computed}(h+t)$ . Performance considerations often compromise this ideal. There is no standard specification for double-double.

The code below uses augmented addition to compute the double-double sum  $s = a + b$ , where  $s$ ,  $a$  and  $b$  are represented by  $(sh, st)$ ,  $(ah, at)$  and  $(bh, bt)$ , respectively, and  $s = sh + st$  closely approximates  $a + b = ah + at + bh + bt$ .

```
#include <augarith.h>

// components of double-double values a = 1/3, b = 2/3, and s
double ah = 0x0.AAAAAAAAAAAAAA8p-1, at = 0x0.AAAAAAAAAAAAAA8p-55;
double bh = 0x0.AAAAAAAAAAAAAA8p0, bt = 0x0.AAAAAAAAAAAAAA8p-54;
double sh, st;

struct daug_t u, v, w, y, z;

// compute components of s = a + b
// exact sum is ah + at + bh + bt
u = aug_add(ah, bh); // exact sum is u.h + u.t + at + bt
v = aug_add(at, bt); // exact sum is u.h + u.t + v.h + v.t
w = aug_add(u.t, v.t); // exact sum is u.h + v.h + w.h + w.t
y = aug_add(v.h, w.h); // exact sum is u.h + y.h + y.t + w.t
z = aug_add(u.h, y.h); // exact sum is z.h + z.t + y.t + w.t

sh = z.h;
st = z.t;
```

The code gives a good approximation to the ideal result, with absolute error  $y.t + w.t$ , and it is commutative and without conditional branches.

The steps for `w` and `y` can use regular addition (+) rather than `aug_add`, because `w.t` and `y.t` are not used in the calculation. The code above gives a name to `w.t` and `y.t` for the didactic purpose of the error formula, and it also assures a consistent result regardless of the evaluation method.