

DRAFT INTERNATIONAL ISO/IEC  
STANDARD **WD 10967-3**

WD for the First edition  
Date: 2000-11-26

---

---

**Information technology —**  
**Language independent arithmetic —**

Part 3: Complex floating point arithmetic and  
complex elementary numerical functions

*Technologies de l'information —*  
*Arithmétique indépendante des langues —*

*Partie 3: Arithmétique des nombres en virgule flottante complexe et*  
*fonctions numériques élémentaires complexe*

**EDITOR'S WORKING DRAFT**  
**November 27, 2000 9:52**

**Editor:**  
**Kent Karlsson**  
**IMI, Industri-Matematik International**  
**Kungsgatan 12**  
**SE-411 19 Göteborg**  
**SWEDEN**  
**Telephone: +46-31 10 22 44**  
**Facsimile: +46-31 13 13 25**  
**E-mail: keka@im.se**

## Copyright notice

This ISO document is a Draft International Standard and is copyright-protected by ISO. Except as permitted under the applicable laws of the user's country, neither this ISO draft nor any extract from it may be reproduced, stored in a retrieval system or transmitted in any form or by any means, electronic, photocopying, recording or otherwise, without prior written permission being secured.

Requests for permission to reproduce should be addressed to ISO at the address below or ISO's member body in the country of the requester.

*Copyright Manager  
ISO Central Secretariat  
1 rue de Varembé  
CH-1211 Genève 20  
Switzerland*

*tel. +41 22 749 0111  
fax. +41 22 734 1079  
e-mail: iso@iso.ch*

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

Violators may be prosecuted.

## Contents

Foreword . . . . .	vii
Introduction . . . . .	viii
<b>1 Scope</b>	<b>1</b>
1.1 Inclusions . . . . .	1
1.2 Exclusions . . . . .	2
<b>2 Conformity</b>	<b>2</b>
<b>3 Normative references</b>	<b>3</b>
<b>4 Symbols and definitions</b>	<b>4</b>
4.1 Symbols . . . . .	4
4.1.1 Sets and intervals . . . . .	4
4.1.2 Operators and relations . . . . .	4
4.1.3 Mathematical functions . . . . .	4
4.1.4 Datatypes and exceptional values . . . . .	4
4.2 Definitions of terms . . . . .	5
<b>5 Specifications for complex datatypes and operations</b>	<b>8</b>
5.1 Complex integer datatypes and operations . . . . .	8
5.1.1 The complex integer <i>result</i> helper function . . . . .	9
5.1.2 Complex integer operations . . . . .	9
5.2 Complex floating point datatypes and operations . . . . .	11
5.2.1 The complex floating point <i>result</i> helper functions . . . . .	11
5.2.2 Basic arithmetic for complex floating point . . . . .	12
5.2.3 Complex multiplication and division . . . . .	14
5.3 Elementary transcendental complex floating point operations . . . . .	15
5.3.1 Operations for exponentiations and logarithms . . . . .	15
5.3.1.1 Natural exponentiation . . . . .	15
5.3.1.2 Complex exponentiation of argument base . . . . .	16
5.3.1.3 Complex square root . . . . .	16
5.3.1.4 Natural logarithm . . . . .	17
5.3.1.5 Argument base logarithm . . . . .	18
5.3.2 Operations for radian trigonometric elementary functions . . . . .	18
5.3.2.1 Radian angle normalisation . . . . .	18
5.3.2.2 Radian sine . . . . .	19
5.3.2.3 Radian cosine . . . . .	19
5.3.2.4 Radian tangent . . . . .	20
5.3.2.5 Radian cotangent . . . . .	21
5.3.2.6 Radian secant . . . . .	21
5.3.2.7 Radian cosecant . . . . .	22
5.3.2.8 Radian arc sine . . . . .	23
5.3.2.9 Radian arc cosine . . . . .	24
5.3.2.10 Radian arc tangent . . . . .	25
5.3.2.11 Radian arc cotangent . . . . .	26
5.3.2.12 Radian arc secant . . . . .	27
5.3.2.13 Radian arc cosecant . . . . .	27
5.3.3 Operations for hyperbolic elementary functions . . . . .	28
5.3.3.1 Hyperbolic normalisation . . . . .	28
5.3.3.2 Hyperbolic sine . . . . .	29

5.3.3.3	Hyperbolic cosine . . . . .	29
5.3.3.4	Hyperbolic tangent . . . . .	29
5.3.3.5	Hyperbolic cotangent . . . . .	29
5.3.3.6	Hyperbolic secant . . . . .	30
5.3.3.7	Hyperbolic cosecant . . . . .	30
5.3.3.8	Inverse hyperbolic sine . . . . .	30
5.3.3.9	Inverse hyperbolic cosine . . . . .	30
5.3.3.10	Inverse hyperbolic tangent . . . . .	31
5.3.3.11	Inverse hyperbolic cotangent . . . . .	31
5.3.3.12	Inverse hyperbolic secant . . . . .	31
5.3.3.13	Inverse hyperbolic cosecant . . . . .	31
5.4	Operations for conversion between numeric datatypes . . . . .	32
5.4.1	Integer to complex integer conversions . . . . .	32
5.4.2	Floating point to complex floating point conversions . . . . .	32
<b>6</b>	<b>Notification</b>	<b>32</b>
6.1	Continuation values . . . . .	33
<b>7</b>	<b>Relationship with language standards</b>	<b>33</b>
<b>8</b>	<b>Documentation requirements</b>	<b>34</b>
<b>Annex A</b>	<b>(normative) Partial conformity</b>	<b>37</b>
A.1	Maximum error relaxation . . . . .	37
A.2	Extra accuracy requirements relaxation . . . . .	37
A.3	Partial conformity to part 1 or to part 2 . . . . .	37
<b>Annex B</b>	<b>(informative) Rationale</b>	<b>39</b>
B.1	Scope . . . . .	39
B.1.1	Inclusions . . . . .	39
B.1.2	Exclusions . . . . .	39
B.2	Conformity . . . . .	40
B.3	Normative references . . . . .	40
B.4	Symbols and definitions . . . . .	40
B.4.1	Symbols . . . . .	40
B.4.1.1	Sets and intervals . . . . .	40
B.4.1.2	Operators and relations . . . . .	40
B.4.1.3	Mathematical functions . . . . .	40
B.4.1.4	Datatypes and exceptional values . . . . .	41
B.4.2	Definitions of terms . . . . .	41
B.5	Specifications for the complex datatypes and operations . . . . .	42
<b>Annex C</b>	<b>(informative) Example bindings for specific languages</b>	<b>43</b>
C.1	Ada . . . . .	44
C.2	C . . . . .	48
C.3	C++ . . . . .	52
C.4	Fortran . . . . .	54
C.5	Haskell . . . . .	56
C.6	Java . . . . .	58
C.7	Common Lisp . . . . .	60
C.8	ISLisp . . . . .	62
C.9	Modula-2 . . . . .	67
C.10	PL/I . . . . .	69

C.11 SML . . . . . 71

**Annex D (informative) Bibliography** **75**



## Foreword

ISO (the International Organization for Standardization) and IEC (the International Electrotechnical Commission) form the specialised system for world-wide 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 organisations, governmental and non-governmental, in liaison with ISO and IEC, also take part in the work.

International Standards are drafted in accordance with the rules in the ISO/IEC Directives, Part 3.

In the field of information technology, ISO and IEC have established a joint technical committee, ISO/IEC JTC 1, *Implementation of information technology*. Draft International Standards adopted by the joint technical committee are circulated to national bodies for voting. Publication as an International Standard requires approval by at least 75 % of the national bodies casting a vote.

Attention is drawn to the possibility that some of the elements of this part of ISO/IEC 10967 may be the subject of patent rights. ISO and IEC shall not be held responsible for identifying any or all such patent rights.

International Standard ISO/IEC 10967-3 was prepared by Joint Technical Committee ISO/IEC JTC 1, *Implementation of information technology*, Subcommittee SC 22, *Programming languages, their environments and system software interfaces*.

ISO/IEC 10967 consists of the following parts, under the general title *Information technology* — *Language independent arithmetic*:

- *Part 1: Integer and floating point arithmetic*
- *Part 2: Elementary numerical functions*
- *Part 3: Complex floating point arithmetic and complex elementary numerical functions*

Additional parts will specify other arithmetic datatypes or arithmetic operations.

**Annex B is informative and is intended to be read in parallel with the main normative text of the standard.**

**Notes and Annexes B, C, and D as well as notes are informative.**

## Introduction

### The aims

Portability is a key issue for scientific and numerical software in today's heterogeneous computing environment. Such software may be required to run on systems ranging from personal computers to high performance pipelined vector processors and massively parallel systems, and the source code may be ported between several programming languages.

Part 1 of ISO/IEC 10967 specifies the basic properties of integer and floating point types that can be relied upon in writing portable software.

Part 2 of ISO/IEC 10967 specifies a number of additional operations for integer and floating point types, in particular specifications for numerical approximations to elementary functions on reals.

### The content

The content of this part is based on part 1 and part 2, and extends part 1's and part 2's specifications to specifications for operations approximating integer complex (Gaussian integers) arithmetic, real complex arithmetic, and real complex elementary functions.

The numerical functions covered by this part are computer approximations to mathematical functions of one or more complex arguments. Accuracy versus performance requirements often vary with the application at hand. This is recognised by recommending that implementors support more than one library of these numerical functions. Various documentation and (program available) parameters requirements are specified to assist programmers in the selection of the library best suited to the application at hand.

### The benefits

Adoption and proper use of this part can lead to the following benefits.

Language standards will be able to define their arithmetic semantics more precisely without preventing the efficient implementation of their language on a wide range of machine architectures.

Programmers of numeric software will be able to assess the portability of their programs in advance. Programmers will be able to trade off program design requirements for portability in the resulting program.

Programs will be able to determine (at run time) the crucial numeric properties of the implementation. They will be able to reject unsuitable implementations, and (possibly) to correctly characterize the accuracy of their own results. Programs will be able to extract apparently implementation dependent data (such as the exponent of a floating point number) in an implementation independent way. Programs will be able to detect (and possibly correct for) exceptions in arithmetic processing.

End users will find it easier to determine whether a (properly documented) application program is likely to execute satisfactorily on their platform. This can be done by comparing the documented requirements of the program against the documented properties of the platform.

Finally, end users of numeric application packages will be able to rely on the correct execution of those packages. That is, for correctly programmed algorithms, the results are reliable if and only if there is no notification.



# Information technology — Language independent arithmetic —

## Part 3: Complex arithmetic and complex elementary numerical functions

### 1 Scope

This part of ISO/IEC 10967 defines the properties of numerical approximations for complex arithmetic operations and many of the complex elementary numerical functions available in standard libraries for a variety of programming languages in common use for mathematical and numerical applications.

An implementor may choose any combination of hardware and software support to meet the specifications of this part. It is the computing environment, as seen by the programmer/user, that does or does not conform to the specifications.

The term *implementation* of this part denotes the total computing environment pertinent to this part, including hardware, language processors, subroutine libraries, exception handling facilities, other software, and documentation.

#### 1.1 Inclusions

The specifications of part 1 and part 2 are included by reference in this part.

This part provides specifications for numerical functions for which operand or result values are of complex integer or complex floating point datatypes constructed from integer and floating point datatypes satisfying the requirements of part 1. Boundaries for the occurrence of exceptions and the maximum error allowed are prescribed for each specified operation. Also the result produced by giving a special value operand, such as an infinity, or a **NaN**, is prescribed for each specified floating point operation.

This part provides specifications for

- a) basic complex integer (Gaussian integer) operations,
- b) non-transcendental Cartesian complex floating point operations,
- c) exponentiations, logarithms, hyperbolics, and
- d) trigonometric operations for Cartesian complex floating point.

This part also provides specifications for

- e) the results produced by an included floating point operation when one or more operand values include IEC 60559 special values, and
- f) program-visible parameters that characterise certain aspects of the operations.

## 1.2 Exclusions

This part provides no specifications for:

- a) Datatypes and operations for polar complex floating point. This standard neither requires nor excludes the presence of such polar complex datatypes and operations.
- b) Numerical functions whose operands are of more than one datatype. This standard neither requires nor excludes the presence of such “mixed operand” operations.
- c) A complex interval datatype, or the operations on such data. This standard neither requires nor excludes such data or operations.
- d) A complex fixed point datatype, or the operations on such data. This standard neither requires nor excludes such data or operations.
- e) A complex rational datatype, or the operations on such data. This standard neither requires nor excludes such data or operations.
- f) Matrix, statistical, or symbolic operations. This standard neither requires nor excludes such data or operations.
- g) The properties of complex arithmetic datatypes that are not related to the numerical process, such as the representation of values on physical media.
- h) The properties of integer and floating point datatypes that properly belong in language standards or other specifications. Examples include
  - 1) the syntax of numerals and expressions in the programming language,
  - 2) the syntax used for parsed (input) or generated (output) character string forms for numerals by any specific programming language or library,
  - 3) the precedence of operators,
  - 4) the consequences of applying an operation to values of improper datatype, or to uninitialised data,
  - 5) the rules for assignment, parameter passing, and returning value,
  - 6) the presence or absence of automatic datatype coercions.

Furthermore, this part does not provide specifications for:

- i) how numerical functions should be implemented,
- j) which algorithms are to be used for the various operations.

## 2 Conformity

It is expected that the provisions of this part of ISO/IEC 10967 will be incorporated by reference and further defined in other International Standards; specifically in language standards and in language binding standards.

A binding standard specifies the correspondence between one or more datatypes, operations, and parameters specified in this part and the concrete language syntax of some programming language. More generally, a binding standard specifies the correspondence between certain operations and the elements of some arbitrary computing entity. A language standard that explicitly provides such binding information can serve as a binding standard.

Conformity to this part is always with respect to a specified set of datatypes and set of operations. Conformity to this part implies conformity to part 1 and part 2 for the integer and floating point datatypes and operations used.

When a binding standard for a language exists, an implementation shall be said to conform to this part if and only if it conforms to the binding standard. In case of conflict between a binding standard and this part, the specifications of the binding standard takes precedence.

When a binding standard covers only a subset of the datatypes and operations defined in this part, an implementation remains free to conform to this part with respect to other datatypes or operations independently of that binding standard.

When no binding standard for a language and some operations specified in this part exists, an implementation conforms to this part if and only if it provides one or more datatypes and one or more operations that together satisfy all the requirements of clauses 5 through 8 that are relevant to those datatypes and operations. The implementation shall then document the binding.

An implementation is free to provide datatypes or operations that do not conform to this part, or that are beyond the scope of this part. The implementation shall not claim or imply conformity to this part with respect to such datatypes or operations.

An implementation is permitted to have modes of operation that do not conform to this part. A conforming implementation shall specify how to select the modes of operation that ensure conformity.

#### NOTES

- 1 Language bindings are essential. Clause 8 requires an implementation to supply a binding if no binding standard exists. See annex C for suggested language bindings.
- 2 A complete binding for this part will include (explicitly or by reference) a binding for part 2 and part 1 as well, which in turn may include (explicitly or by reference) a binding for IEC 60559 as well.
- 3 It is not possible to conform to this part without specifying to which set of datatypes and set of operations conformity is claimed.

### 3 Normative references

The following standards contain provisions which, through reference in this text, constitute provisions of this part. At the time of publication, the editions indicated were valid. All standards are subject to revision, and parties to agreements based on this part are encouraged to investigate the possibility of applying the most recent edition of the standards indicated below. Members of IEC and ISO maintain registers of currently valid International Standards.

IEC 60559:1989, *Binary floating-point arithmetic for microprocessor systems*.

ISO/IEC 10967-1:2002, *Information technology — Language independent arithmetic — Part 1: Integer and floating point arithmetic*.

ISO/IEC 10967-2:2000, *Information technology — Language independent arithmetic — Part 2: Elementary numerical functions*.

## 4 Symbols and definitions

### 4.1 Symbols

#### 4.1.1 Sets and intervals

In this part,  $\mathcal{Z}$  denotes the set of mathematical integers,  $\mathcal{G}$  denotes the set of Gaussian integers (complex integers),  $\mathcal{R}$  denotes the set of classical real numbers, and  $\mathcal{C}$  denotes the set of complex numbers over  $\mathcal{R}$ . Note that  $\mathcal{Z} \subset \mathcal{R} \subset \mathcal{C}$ .

$[x, z]$  designates the interval  $\{y \in \mathcal{R} \mid x \leq y \leq z\}$ ,  
 $]x, z]$  designates the interval  $\{y \in \mathcal{R} \mid x < y \leq z\}$ ,  
 $[x, z[$  designates the interval  $\{y \in \mathcal{R} \mid x \leq y < z\}$ , and  
 $]x, z[$  designates the interval  $\{y \in \mathcal{R} \mid x < y < z\}$ .

NOTE – The notation using a round bracket for an open end of an interval is not used, for the risk of confusion with the notation for pairs.

#### 4.1.2 Operators and relations

All prefix and infix operators have their conventional (exact) mathematical meaning. The conventional notation for set definition and manipulation is also used. In particular this part uses

$\Rightarrow$  and  $\Leftrightarrow$  for logical implication and equivalence  
 $+$ ,  $-$ ,  $/$ , and  $|x|$  on complex values  
 $\cdot$  for multiplication on complex values  
 $<$ ,  $\leq$ ,  $=$ ,  $\neq$ ,  $\geq$ , and  $>$  between reals  
 $\cup$ ,  $\cap$ ,  $\times$ ,  $\in$ ,  $\notin$ ,  $\subset$ ,  $\subseteq$ ,  $\not\subseteq$ ,  $\neq$ , and  $=$  with sets  
 $\times$  for the Cartesian product of sets  
 $\rightarrow$  for a mapping between sets

#### 4.1.3 Mathematical functions

This part specifies properties for a number of operations numerically approximating some of the elementary functions. The following ideal mathematical functions are defined in Chapter 4 of the *Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables* [43] ( $e$  is the Napierian base):

$e^x$ ,  $x^y$ ,  $\sqrt{x}$ ,  $\ln$ ,  $\log_b$ ,  
 $\sinh$ ,  $\cosh$ ,  $\tanh$ ,  $\coth$ ,  $\operatorname{sech}$ ,  $\operatorname{csch}$ ,  $\operatorname{arcsinh}$ ,  $\operatorname{arccosh}$ ,  $\operatorname{arctanh}$ ,  $\operatorname{arcoth}$ ,  $\operatorname{arcsech}$ ,  $\operatorname{arccsch}$ ,  
 $\sin$ ,  $\cos$ ,  $\tan$ ,  $\cot$ ,  $\sec$ ,  $\csc$ ,  $\arcsin$ ,  $\arccos$ ,  $\arctan$ ,  $\operatorname{arccot}$ ,  $\operatorname{arcsec}$ ,  $\operatorname{arccsc}$ .

Many of the inverses above are multi-valued. The selection of which value to return, the principal value, so as to make the inverses into functions, is done in the conventional way. E.g.,  $\sqrt{x} \in [0, \infty[$  when  $x \in [0, \infty[$ .

#### 4.1.4 Datatypes and exceptional values

The datatype **Boolean** consists of the two values **true** and **false**.

Integer datatypes and floating point datatypes are defined in part 1.

Let  $I$  be the non-special value set for an integer datatype conforming to part 1. Let  $F$  be the non-special value set for a floating point datatype conforming to part 1. The following symbols are defined in part 1 or part 2, and used in this part.

Exceptional values:

**overflow**, **underflow**, **invalid**, **infinitary**, and **absolute\_precision\_underflow**.

Integer operations:

$neg_I$ ,  $add_I$ ,  $sub_I$ , and  $mul_I$ .

Floating point helper functions:

$e_F$ ,  $u_F$ ,  $result_F$ , and  $rnd_F$ .

Floating point operations from part 1:

$sign_F$ ,  $neg_F$ ,  $add_F$ ,  $sub_F$ ,  $mul_F$ , and  $div_F$ .

Floating point operations from part 2:

$sqrt_F$ ,  $hypot_F$ ,  $exp_F$ ,  $power_F$ ,  $ln_F$ ,  $logbase_F$

$sinh_F$ ,  $cosh_F$ ,  $tanh_F$ ,  $coth_F$ ,  $sech_F$ ,  $csch_F$ ,

$arcsinh_F$ ,  $arccosh_F$ ,  $arctanh_F$ ,  $arccoth_F$ ,  $arcsech_F$ ,  $arccsch_F$ ,

$rad_F$ ,  $sin_F$ ,  $cos_F$ ,  $tan_F$ ,  $cot_F$ ,  $sec_F$ ,  $csc_F$ ,

$arcsin_F$ ,  $arccos_F$ ,  $arctan_F$ ,  $arccot_F$ ,  $arcsec_F$ ,  $arccsc_F$ ,  $arcF$ .

Floating point datatypes that conform to part 1 shall, for use with this part, have a value for the parameter  $p_F$  such that  $p_F \geq 2 \cdot \max\{1, \log_{r_F}(2 \cdot \pi)\}$ , and have a value for the parameter  $emin_F$  such that  $emin_F \leq -p_F - 1$ .

NOTES

- 1 This implies that  $fminN_F < 0.5 \cdot \epsilon_F / r_F$  in this part, rather than just  $fminN_F \leq \epsilon_F$ .
- 2 These extra requirements, which do not limit the use of any existing floating point datatype, are made so that angles in radians are not too degenerate within the first two cycles, plus and minus, when represented in  $F$ .
- 3  $F$  should also be such that  $p_F \geq 2 + \log_{r_F}(1000)$ , to allow for a not too coarse angle resolution anywhere in the interval  $[-big\_angle_{r_F}, big\_angle_{r_F}]$ . See clause 5.3.9 of part 2.

The following symbols represent floating point values defined in IEC 60559 and used in this part:

**-0**, **+∞**, **-∞**, **qNaN**, and **sNaN**.

These floating point values are not part of the set  $F$ , but if  $iec\_559_F$  has the value **true**, these values are included in the floating point datatype corresponding to  $F$ .

NOTE 4 – This part uses the above five special values for compatibility with IEC 60559. In particular, the symbol **-0** (in bold) is not the application of (mathematical) unary  $-$  to the value 0, and is a value logically distinct from 0.

The specifications cover the results to be returned by an operation if given one or more of the IEC 60559 special values **-0**, **+∞**, **-∞**, or **NaNs** as input values. These specifications apply only to systems which provide and support these special values. If an implementation is not capable of representing a **-0** result or continuation value, the actual result or continuation value shall be 0. If an implementation is not capable of representing a prescribed result or continuation value of the IEC 60559 special values **+∞**, **-∞**, or **qNaN**, the actual result or continuation value is binding or implementation defined.

## 4.2 Definitions of terms

For the purposes of this part, the following definitions apply:

**accuracy**: The closeness between the true mathematical result and a computed result.

**arithmetic datatype**: A datatype whose non-special values are members of  $\mathcal{Z}$ ,  $\mathcal{G}$ ,  $\mathcal{R}$ , or  $\mathcal{C}$ .

**continuation value**: A computational value used as the result of an arithmetic operation when an exception occurs. Continuation values are intended to be used in subsequent arithmetic

processing. A continuation value can be a (in the datatype representable) value in  $\mathcal{R}$  or an IEC 60559 special value. (Contrast with *exceptional value*. See 6.1.2 of part 1.)

**denormalisation loss:** A larger than normal rounding error caused by the fact that subnormal values have less than full precision. (See 5.2.5 of part 1 for a full definition.)

**error:** (1) The difference between a computed value and the correct value. (Used in phrases like “rounding error” or “error bound”.)

(2) A synonym for *exception* in phrases like “error message” or “error output”. Error and exception are not synonyms in any other context.

**exception:** The inability of an operation to return a suitable finite numeric result from finite arguments. This might arise because no such finite result exists mathematically, or because the mathematical result cannot be represented with sufficient accuracy.

**exceptional value:** A non-numeric value produced by an arithmetic operation to indicate the occurrence of an exception. Exceptional values are not used in subsequent arithmetic processing. (See clause 5 of part 1.)

#### NOTES

- 1 Exceptional values are used as part of the defining formalism only. With respect to this part, they do not represent values of any of the datatypes described. There is no requirement that they be represented or stored in the computing system.
- 2 Exceptional values are not to be confused with the NaNs and infinities defined in IEC 60559. Contrast this definition with that of *continuation value* above.

**helper function:** A function used solely to aid in the expression of a requirement. Helper functions are not visible to the programmer, and are not required to be part of an implementation.

**implementation** (of this part): The total arithmetic environment presented to a programmer, including hardware, language processors, exception handling facilities, subroutine libraries, other software, and all pertinent documentation.

**literal:** A syntactic entity denoting a constant value without having proper sub-entities that are expressions.

**monotonic approximation:** An floating point operation  $op_F : \dots \times F \times \dots \rightarrow F$ , for a floating point datatype with non-special value set  $F$ , where the other arguments are kept constant, is a monotonic approximation of a predetermined mathematical function  $h : \mathcal{R} \rightarrow \mathcal{R}$  if, for every  $a \in F$  and  $b \in F$ ,  $a < b$ ,

- a)  $h$  is monotonic non-decreasing on  $[a, b]$  implies  $op_F(\dots, a, \dots) \leq op_F(\dots, b, \dots)$ ,
- b)  $h$  is monotonic non-increasing on  $[a, b]$  implies  $op_F(\dots, a, \dots) \geq op_F(\dots, b, \dots)$ .

**monotonic non-decreasing:** A function  $h : \mathcal{R} \rightarrow \mathcal{R}$  is monotonic non-decreasing on a real interval  $[a, b]$  if for every  $x$  and  $y$  such that  $a \leq x \leq y \leq b$ ,  $h(x)$  and  $h(y)$  are well-defined and  $h(x) \leq h(y)$ .

**monotonic non-increasing:** A function  $h : \mathcal{R} \rightarrow \mathcal{R}$  is monotonic non-increasing on a real interval  $[a, b]$  if for every  $x$  and  $y$  such that  $a \leq x \leq y \leq b$ ,  $h(x)$  and  $h(y)$  are well-defined and  $h(x) \geq h(y)$ .

**normalised:** The non-zero values of a floating point type  $F$  that provide the full precision allowed by that type. (See  $F_N$  in 5.2 of part 1 for a full definition.)

**notification:** The process by which a program (or that program’s end user) is informed that an arithmetic exception has occurred. For example, dividing 2 by 0 results in a notification. (See clause 6 of part 1 for details.)

**numeral:** A numeric literal. It may denote a value in  $\mathcal{Z}$  or  $\mathcal{R}$ ,  $-\mathbf{0}$ , an infinity, or a NaN.

**numerical function:** A computer routine or other mechanism for the approximate evaluation of a mathematical function.

**operation:** A function directly available to the programmer, as opposed to helper functions or theoretical mathematical functions.

**pole:** A mathematical function  $f$  has a pole at  $x_0$  if  $x_0$  is finite,  $f$  is defined, finite, monotone, and continuous in at least one side of the neighbourhood of  $x_0$ , and  $\lim_{x \rightarrow x_0} f(x)$  is infinite.

**precision:** The number of digits in the fraction of a floating point number. (See clause 5.2 of part 1.)

**rounding:** The act of computing a representable final result for an operation that is close to the exact (but unrepresentable) result for that operation. Note that a suitable representable result may not exist (see 5.2.6 of part 1). (See also A.5.2.6 of part 1 for some examples.)

**rounding function:** Any function  $rnd: \mathcal{R} \rightarrow X$  (where  $X$  is a given discrete and unlimited subset of  $\mathcal{R}$ ) that maps each element of  $X$  to itself, and is monotonic non-decreasing. Formally, if  $x$  and  $y$  are in  $\mathcal{R}$ ,

$$\begin{aligned} x \in X &\Rightarrow rnd(x) = x \\ x < y &\Rightarrow rnd(x) \leq rnd(y) \end{aligned}$$

Note that if  $u \in \mathcal{R}$  is between two adjacent values in  $X$ ,  $rnd(u)$  selects one of those adjacent values.

**round to nearest:** The property of a rounding function  $rnd$  that when  $u \in \mathcal{R}$  is between two adjacent values in  $X$ ,  $rnd(u)$  selects the one nearest  $u$ . If the adjacent values are equidistant from  $u$ , either may be chosen deterministically.

**round toward minus infinity:** The property of a rounding function  $rnd$  that when  $u \in \mathcal{R}$  is between two adjacent values in  $X$ ,  $rnd(u)$  selects the one less than  $u$ .

**round toward plus infinity:** The property of a rounding function  $rnd$  that when  $u \in \mathcal{R}$  is between two adjacent values in  $X$ ,  $rnd(u)$  selects the one greater than  $u$ .

**shall:** A verbal form used to indicate requirements strictly to be followed in order to conform to the standard and from which no deviation is permitted. (Quoted from the directives [1].)

**should:** A verbal form used to indicate that among several possibilities one is recommended as particularly suitable, without mentioning or excluding others; or that (in the negative form) a certain possibility is deprecated but not prohibited. (Quoted from the directives [1].)

**signature** (of a function or operation): A summary of information about an operation or function. A signature includes the function or operation name; a subset of allowed argument values to the operation; and a superset of results from the function or operation (including exceptional values if any), if the argument is in the subset of argument values given in the signature.

The signature

$$add_I : I \times I \rightarrow I \cup \{\mathbf{overflow}\}$$

states that the operation named  $add_I$  shall accept any pair of  $I$  values as input, and (when given such input) shall return either a single  $I$  value as its output or the exceptional value **overflow**.

A signature for an operation or function does not forbid the operation from accepting a wider range of arguments, nor does it guarantee that every value in the result range will actually be returned for some input. An operation given an argument outside the stipulated argument domain may produce a result outside the stipulated result range.

**subnormal:** The non-zero values of a floating point type  $F$  that provide less than the full precision allowed by that type. (See  $F_D$  in 5.2 of part 1 for a full definition.)

**ulp:** The value of one “unit in the last place” of a floating point number. This value depends on the exponent, the radix, and the precision used in representing the number. Thus, the ulp of a normalised value  $x$  (in  $F$ ), with exponent  $t$ , precision  $p$ , and radix  $r$ , is  $r^{t-p}$ , and the ulp of a subnormal value is  $fminD_F$ . (See 5.2 of part 1.)

## 5 Specifications for complex datatypes and operations

This clause specifies Gaussian integer (complex integer) datatypes, complex floating point datatypes and a number of helper functions and operations for complex integer and complex floating point datatypes.

Each operation is given a signature and is further specified by a number of cases. These cases may refer to other operations (specified in this part, in part 1, or in part 2), to mathematical functions, and to helper functions (specified in this part, in part 1, or in part 2). They also use special abstract values ( $-\infty$ ,  $+\infty$ ,  $-0$ , **qNaN**, **sNaN**). For each datatype, two of these abstract values may represent several actual values each: **qNaN** and **sNaN**. Finally, the specifications may refer to exceptional values.

The signatures in the specifications in this clause specify only all non-special values as input values, and indicate as output values a superset of all non-special, special, and exceptional values that may result from these (non-special) input values. Exceptional and special values that can never result from non-special input values are not included in the signatures given. Also, signatures that, for example, include IEC 60559 special values as arguments are not given in the specifications below. This does not exclude such signatures from being valid for these operations.

### 5.1 Complex integer datatypes and operations

Clause 5.1 of part 1 and clause 5.1 of part 2 specify integer datatypes and a number of operations on values of an integer datatype. In this clause complex integer (Gaussian integer) datatypes and operations on values of a complex integer datatype are specified.

A complex integer datatype is constructed from an integer datatype. For each integer datatype, there is one complex integer datatype.

$I$  is the set of non-special values,  $I \subseteq \mathcal{Z}$ , for an integer datatype conforming to part 1. Integer datatypes conforming to part 1 often do not contain any **NaN** or infinity values, even though they may do so. Therefore this clause has no specifications for such values as arguments or results.

$i(I)$  is the set of non-special values in a imaginary integer datatype, constructed from the datatype corresponding to  $I$ .

$$i(I) = \{\hat{\mathbf{i}} \cdot x' \mid x' \in I\}$$

$c(I)$  is the set of non-special values in a complex integer datatype, constructed from the datatype corresponding to  $I$ .

$$c(I) = \{x + \hat{\mathbf{i}} \cdot x' \mid x, x' \in I\}$$



### 5.1.1 The complex integer *result* helper function

The  $result_{c(I)}$  helper function:

$$result_{c(I)} : \mathcal{G} \rightarrow c(I) \cup \{\mathbf{overflow}\}$$

$$result_{c(I)}(z) = result_I(\Re(z)) + \hat{\mathbf{i}} \cdot result_I(\Im(z))$$

NOTE – If one or both of the  $result_I$  function applications on the right side returns **overflow**, then the  $result_{c(I)}$  application returns **overflow**. Similarly below for the specifications that do not use  $result_{c(I)}$  but specify the result parts directly.

### 5.1.2 Complex integer operations

$$itimes_I : I \rightarrow i(I) \cup \{\mathbf{overflow}\}$$

$$itimes_I(x) = \hat{\mathbf{i}} \cdot x$$

$$itimes_{c(I)} : c(I) \rightarrow c(I) \cup \{\mathbf{overflow}\}$$

$$\begin{aligned} itimes_{c(I)}(x + \hat{\mathbf{i}} \cdot x') \\ = neg_I(x') + \hat{\mathbf{i}} \cdot x \end{aligned}$$

$$re_I : I \rightarrow I$$

$$re_I(x) = x$$

$$re_{i(I)} : i(I) \rightarrow I$$

$$re_{i(I)}(\hat{\mathbf{i}} \cdot x') = 0$$

$$re_{c(I)} : c(I) \rightarrow I$$

$$re_{c(I)}(x + \hat{\mathbf{i}} \cdot x') = x$$

$$im_I : I \rightarrow I$$

$$im_I(x) = 0$$

$$im_{i(I)} : i(I) \rightarrow I$$

$$im_{i(I)}(\hat{\mathbf{i}} \cdot x') = x'$$

$$im_{c(I)} : c(I) \rightarrow I$$

$$im_{c(I)}(x + \hat{\mathbf{i}} \cdot x') = x'$$

$$plusitimes_{c(I)} : I \times I \rightarrow c(I)$$

$$\begin{aligned} plusitimes_{c(I)}(x, y) \\ = x + \hat{\mathbf{i}} \cdot y \end{aligned}$$

$$neg_{i(I)} : i(I) \rightarrow i(I) \cup \{\mathbf{overflow}\}$$

$$neg_{i(I)}(\hat{\mathbf{i}} \cdot x') = \hat{\mathbf{i}} \cdot neg_I(x')$$

$$neg_{c(I)} : c(I) \rightarrow c(I) \cup \{\mathbf{overflow}\}$$

$$\begin{aligned} neg_{c(I)}(x + \hat{\mathbf{i}} \cdot x') \\ &= neg_I(x) + \hat{\mathbf{i}} \cdot neg_I(x') \end{aligned}$$

$$\begin{aligned} conj_I : I &\rightarrow I \cup \{\mathbf{overflow}\} \\ conj_I(x) &= x \end{aligned}$$

$$\begin{aligned} conj_{i(I)} : i(I) &\rightarrow i(I) \cup \{\mathbf{overflow}\} \\ conj_{i(I)}(\hat{\mathbf{i}} \cdot x') &= \hat{\mathbf{i}} \cdot neg_I(x') \end{aligned}$$

$$\begin{aligned} conj_{c(I)} : c(I) &\rightarrow c(I) \cup \{\mathbf{overflow}\} \\ conj_{c(I)}(x + \hat{\mathbf{i}} \cdot x') \\ &= x + \hat{\mathbf{i}} \cdot neg_I(x') \end{aligned}$$

$$\begin{aligned} add_{i(I)} : i(I) \times i(I) &\rightarrow i(I) \cup \{\mathbf{overflow}\} \\ add_{i(I)}(\hat{\mathbf{i}} \cdot x', \hat{\mathbf{i}} \cdot y') \\ &= \hat{\mathbf{i}} \cdot add_I(x', y') \end{aligned}$$

$$\begin{aligned} add_{c(I)} : c(I) \times c(I) &\rightarrow c(I) \cup \{\mathbf{overflow}\} \\ add_{c(I)}(x + \hat{\mathbf{i}} \cdot x', y + \hat{\mathbf{i}} \cdot y') \\ &= add_F(x, y) + \hat{\mathbf{i}} \cdot add_F(x', y') \end{aligned}$$

$$\begin{aligned} sub_{i(I)} : i(I) \times i(I) &\rightarrow i(I) \cup \{\mathbf{overflow}\} \\ sub_{i(I)}(x, y) &= add_{i(I)}(x, neg_{i(I)}(y)) \end{aligned}$$

$$\begin{aligned} sub_{c(I)} : c(I) \times c(I) &\rightarrow c(I) \cup \{\mathbf{overflow}\} \\ sub_{c(I)}(x, y) &= add_{c(I)}(x, neg_{c(I)}(y)) \end{aligned}$$

$$\begin{aligned} mul_{c(I)} : c(I) \times c(I) &\rightarrow c(I) \cup \{\mathbf{overflow}\} \\ mul_{c(I)}(x + \hat{\mathbf{i}} \cdot x', y + \hat{\mathbf{i}} \cdot y') \\ &= result_{c(I)}((x + \hat{\mathbf{i}} \cdot x') \cdot (y + \hat{\mathbf{i}} \cdot y')) \\ &\quad \text{if } x, y \in c(I) \end{aligned}$$

$$\begin{aligned} eq_{c(I)} : c(I) \times c(I) &\rightarrow \mathbf{Boolean} \\ eq_{c(I)}(x, y) &= \mathbf{true} && \text{if } x, y \in c(I) \text{ and } x = y \\ &= \mathbf{false} && \text{if } x, y \in c(I) \text{ and } x \neq y \end{aligned}$$

$$\begin{aligned} neq_{c(I)} : c(I) \times c(I) &\rightarrow \mathbf{Boolean} \\ neq_{c(I)}(x, y) &= \mathbf{true} && \text{if } x, y \in c(I) \text{ and } x \neq y \\ &= \mathbf{false} && \text{if } x, y \in c(I) \text{ and } x = y \end{aligned}$$

## 5.2 Complex floating point datatypes and operations

Clause 5.2 of part 1 and clause 5.2 of part 2 specify floating point datatypes and a number of operations on values of a floating point datatype. In this clause complex floating point datatypes and operations on values of a complex floating point datatype are specified.

NOTE – Further operations on values of a complex floating point datatype, for elementary complex floating point numerical functions, are specified in clause 5.3.

$F$  is the non-special value set,  $F \subset \mathcal{R}$ , for a floating point datatype conforming to part 1. Floating point datatypes conforming to part 1 often do contain  $-0$ , infinity, and **NaN** values. Therefore, in this clause there are specifications for such values as arguments.

$i(F)$  is the set of non-special values in a imaginary floating point datatype, constructed from the datatype corresponding to  $F$ .

$$i(F) = \{\hat{\mathbf{i}} \cdot y \mid y \in F\}$$

$c(F)$  is the set of non-special values in a complex floating point datatype, constructed from the datatype corresponding to  $F$ .

$$c(F) = \{x + \hat{\mathbf{i}} \cdot x' \mid x, x' \in F\}$$

### 5.2.1 The complex floating point *result* helper functions

$$result_{c(F)}^* : \mathcal{C} \times (\mathcal{R} \rightarrow F^*) \rightarrow F \cup \{\mathbf{underflow}, \mathbf{overflow}\}$$

$$result_{c(F)}^*(z, rnd) \\ = result_F^*(\Re(z), rnd) + \hat{\mathbf{i}} \cdot result_F^*(\Im(z), rnd)$$

overflow, underflow, ... combination...!, (inexact)????

$result_F^*$  is defined in part 2.

Define the  $no\_result_{c(F)}$  and  $no\_result2_{c(F)}$  helper functions:

$$no\_result_{c(F)} : c(F) \rightarrow \{\mathbf{invalid}\}$$

$$no\_result_{c(F)}(x + \hat{\mathbf{i}} \cdot x') \\ = \mathbf{qNaN} + \hat{\mathbf{i}} \cdot \mathbf{qNaN} \quad \text{if } x \text{ is a quiet NaN and } x' \text{ is not a signalling NaN} \\ = \mathbf{qNaN} + \hat{\mathbf{i}} \cdot \mathbf{qNaN} \quad \text{if } x' \text{ is a quiet NaN and } x \text{ is not a signalling NaN} \\ = \mathbf{invalid}(\mathbf{qNaN} + \hat{\mathbf{i}} \cdot \mathbf{qNaN}) \\ \quad \text{if } x, x' \in F \cup \{-\infty, -0, +\infty\} \\ = \mathbf{invalid}(\mathbf{qNaN} + \hat{\mathbf{i}} \cdot \mathbf{qNaN}) \\ \quad \text{if } x \text{ is a signalling NaN or } x' \text{ is a signalling NaN}$$

$$no\_result2_{c(F)} : c(F) \times c(F) \rightarrow \{\mathbf{invalid}\}$$

$$no\_result2_{c(F)}(x + \hat{\mathbf{i}} \cdot x', y + \hat{\mathbf{i}} \cdot y') \\ = \mathbf{qNaN} + \hat{\mathbf{i}} \cdot \mathbf{qNaN} \quad \text{if } x \text{ is a quiet NaN and neither } x', y, \text{ nor } y' \text{ is a signalling NaN} \\ = \mathbf{qNaN} + \hat{\mathbf{i}} \cdot \mathbf{qNaN} \quad \text{if } x' \text{ is a quiet NaN and neither } x, y, \text{ nor } y' \text{ is not a signalling NaN} \\ = \mathbf{qNaN} + \hat{\mathbf{i}} \cdot \mathbf{qNaN} \quad \text{if } y \text{ is a quiet NaN and neither } x, x', \text{ nor } y' \text{ is not a signalling NaN} \\ = \mathbf{qNaN} + \hat{\mathbf{i}} \cdot \mathbf{qNaN} \quad \text{if } y' \text{ is a quiet NaN and neither } x, x', \text{ nor } y \text{ is not a signalling NaN} \\ = \mathbf{invalid}(\mathbf{qNaN} + \hat{\mathbf{i}} \cdot \mathbf{qNaN}) \\ \quad \text{if } x, x', y, y' \in F \cup \{-\infty, -0, +\infty\} \\ = \mathbf{invalid}(\mathbf{qNaN} + \hat{\mathbf{i}} \cdot \mathbf{qNaN}) \\ \quad \text{if at least one of } x, x', y, \text{ or } y' \text{ is a signalling NaN}$$

These helper functions are used to specify both NaN argument handling and to handle non-NaN-argument cases where **invalid(qNaN + î · qNaN)** is the appropriate result.

NOTE – The handling of other special values, if available, is left unspecified by this part.

### 5.2.2 Basic arithmetic for complex floating point

$$itimes_F : F \rightarrow i(F) \cup \{\mathbf{overflow}\}$$

$$itimes_F(x) = \hat{\mathbf{i}} \cdot x$$

$$itimes_{c(F)} : c(F) \rightarrow c(F) \cup \{\mathbf{overflow}\}$$

$$itimes_{c(F)}(x + \hat{\mathbf{i}} \cdot x') = neg_F(x') + \hat{\mathbf{i}} \cdot x$$

$$re_F : F \rightarrow F$$

$$re_F(x) = x$$

$$re_{i(F)} : i(F) \rightarrow F$$

$$re_{i(F)}(\hat{\mathbf{i}} \cdot x') = -\mathbf{0}$$

$$re_{c(F)} : c(F) \rightarrow F$$

$$re_{c(F)}(x + \hat{\mathbf{i}} \cdot x') = x$$

$$im_F : F \rightarrow F$$

$$im_F(x) = -\mathbf{0}$$

$$im_{i(F)} : i(F) \rightarrow F$$

$$im_{i(F)}(\hat{\mathbf{i}} \cdot x') = x'$$

$$im_{c(F)} : c(F) \rightarrow F$$

$$im_{c(F)}(x + \hat{\mathbf{i}} \cdot x') = x'$$

$$plusitimes_{c(F)} : F \times F \rightarrow c(F)$$

$$plusitimes_{c(F)}(x, y) = x + \hat{\mathbf{i}} \cdot y$$

$$neg_{i(F)} : i(F) \rightarrow i(F)$$

$$neg_{i(F)}(\hat{\mathbf{i}} \cdot x') = \hat{\mathbf{i}} \cdot neg_F(x')$$

$$neg_{c(F)} : c(F) \rightarrow c(F)$$

$$neg_{c(F)}(x + \hat{\mathbf{i}} \cdot x') = neg_F(x) + \hat{\mathbf{i}} \cdot neg_F(x')$$

$$conj_F : F \rightarrow F$$

$$\text{conj}_F(x) = x$$

$$\text{conj}_{i(F)} : i(F) \rightarrow i(F)$$

$$\text{conj}_{i(F)}(\hat{\mathbf{i}} \cdot x') = \hat{\mathbf{i}} \cdot \text{neg}_F(x')$$

$$\text{conj}_{c(F)} : c(F) \rightarrow c(F)$$

$$\begin{aligned} \text{conj}_{c(F)}(x + \hat{\mathbf{i}} \cdot x') \\ = x + \hat{\mathbf{i}} \cdot \text{neg}_F(x') \end{aligned}$$

$$\text{add}_{i(F)} : i(F) \times i(F) \rightarrow i(F) \cup \{\mathbf{overflow}\}$$

$$\begin{aligned} \text{add}_{i(F)}(\hat{\mathbf{i}} \cdot x', \hat{\mathbf{i}} \cdot y') \\ = \hat{\mathbf{i}} \cdot \text{add}_F(x', y') \end{aligned}$$

$$\text{add}_{c(F)} : c(F) \times c(F) \rightarrow c(F) \cup \{\mathbf{overflow}\}$$

$$\begin{aligned} \text{add}_{c(F)}(x + \hat{\mathbf{i}} \cdot x', y + \hat{\mathbf{i}} \cdot y') \\ = \text{add}_F(x, x') + \hat{\mathbf{i}} \cdot \text{add}_F(y, y') \end{aligned}$$

$$\text{sub}_{i(F)} : i(F) \times i(F) \rightarrow i(F) \cup \{\mathbf{overflow}\}$$

$$\begin{aligned} \text{sub}_{i(F)}(\hat{\mathbf{i}} \cdot x', \hat{\mathbf{i}} \cdot y') \\ = \hat{\mathbf{i}} \cdot \text{sub}_F(x', y') \end{aligned}$$

$$\text{sub}_{c(F)} : c(F) \times c(F) \rightarrow c(F) \cup \{\mathbf{overflow}\}$$

$$\text{sub}_{c(F)}(x, y) = \text{add}_{c(F)}(x, \text{neg}_{c(F)}(y))$$

$$\text{mul}_{i(F)} : i(F) \times i(F) \rightarrow F \cup \{\mathbf{overflow}\}$$

$$\begin{aligned} \text{mul}_{i(F)}(\hat{\mathbf{i}} \cdot x, \hat{\mathbf{i}} \cdot x') \\ = \text{neg}_F(\text{mul}_F(x, x')) \end{aligned}$$

NOTE 1 –  $\text{mul}_{c(F)}$  is specified in clause 5.2.3

$$\text{div}_{i(F)} : i(F) \times i(F) \rightarrow F \cup \{\mathbf{overflow}\}$$

$$\begin{aligned} \text{div}_{i(F)}(\hat{\mathbf{i}} \cdot x, \hat{\mathbf{i}} \cdot x') \\ = \text{div}_F(x, x') \end{aligned}$$

NOTE 2 –  $\text{div}_{c(F)}$  is specified in clause 5.2.3

$$\text{eq}_{c(F)} : c(F) \times c(F) \rightarrow \mathbf{Boolean}$$

$$\begin{aligned} \text{eq}_{c(F)}(x, y) &= \mathbf{true} && \text{if } x, y \in c(F) \text{ and } x = y \\ &= \mathbf{false} && \text{if } x, y \in c(F) \text{ and } x \neq y \end{aligned}$$

$$\text{neq}_{c(F)} : c(F) \times c(F) \rightarrow \mathbf{Boolean}$$

$$\begin{aligned} \text{neq}_{c(F)}(x, y) &= \mathbf{true} && \text{if } x, y \in c(F) \text{ and } x \neq y \\ &= \mathbf{false} && \text{if } x, y \in c(F) \text{ and } x = y \end{aligned}$$

$$\text{abs}_{i(F)} : i(F) \rightarrow F$$

$$\text{abs}_{i(F)}(\hat{\mathbf{i}} \cdot x') = \text{abs}_F(x')$$

$$abs_{c(F)} : c(F) \rightarrow F \cup \{\mathbf{underflow}, \mathbf{overflow}\}$$

$$abs_{c(F)}(x + \hat{\mathbf{i}} \cdot x') = hypot_F(x, x')$$

$$phase_F : F \rightarrow F$$

$$phase_F(x) = arc_F(x, -\mathbf{0})$$

$$phase_{i(F)} : i(F) \rightarrow F$$

$$phase_{i(F)}(\hat{\mathbf{i}} \cdot x') = arc_F(-\mathbf{0}, x')$$

$$phase_{c(F)} : c(F) \rightarrow F \cup \{\mathbf{underflow}\}$$

$$phase_{c(F)}(x + \hat{\mathbf{i}} \cdot x') = arc_F(x, x')$$

$$sign_{i(F)} : i(F) \rightarrow i(F)$$

$$sign_{i(F)}(\hat{\mathbf{i}} \cdot x') = \hat{\mathbf{i}} \cdot sign_F(x')$$

$$sign_{c(F)} : c(F) \rightarrow c(F) \cup \{\mathbf{underflow}\}$$

$$sign_{c(F)}(x + \hat{\mathbf{i}} \cdot x') = sin_F(arc_F(x, x')) + \hat{\mathbf{i}} \cdot cos_F(arc_F(x, x'))$$

### 5.2.3 Complex multiplication and division

There shall be two maximum error parameters for complex multiplication and division.

$$max\_error\_mul_{c(F)} \in F$$

$$max\_error\_div_{c(F)} \in F$$

no monotonicity requirements? no sign requirements? dependency on argument values?

The  $mul_{c(F)}^*$  approximation helper function:

$$mul_{c(F)}^* : \mathcal{C} \times \mathcal{C} \rightarrow \mathcal{C}$$

$mul_{c(F)}^*(x, y)$  returns a close approximation to  $x \cdot y$  in  $\mathcal{C}$  with maximum error  $max\_error\_mul_{c(F)}$ .

Further requirement on the  $mul_{c(F)}^*$  approximation helper function are:

$$mul_{c(F)}^*(z, z') = mul_{c(F)}^*(z', z) \quad \text{if } z, z' \in \mathcal{C}$$

$$mul_{c(F)}^*(-z, z') = -mul_{c(F)}^*(z, z') \quad \text{if } z, z' \in \mathcal{C}$$

$$mul_{c(F)}^*(conj(z), conj(z')) = conj(mul_{c(F)}^*(z, z')) \quad \text{if } z, z' \in \mathcal{C}$$

The  $mul_{c(F)}$  operation:

$$mul_{c(F)} : c(F) \times c(F) \rightarrow c(F) \cup \{-\mathbf{0}, \dots, \mathbf{underflow}, \mathbf{overflow}\}$$

$$mul_{c(F)}(x + \hat{\mathbf{i}} \cdot x', y + \hat{\mathbf{i}} \cdot y') = result_{c(F)}^*(mul_{c(F)}^*(x + \hat{\mathbf{i}} \cdot x', y + \hat{\mathbf{i}} \cdot y'))$$

$$= sub_F(mul_F(x, y), mul_F(x', y')) + \hat{\mathbf{i}} \cdot add_F(mul_F(x, y'), mul_F(x', y))$$

$$\text{otherwise}$$

The  $div_{c(F)}^*$  approximation helper function:

$$\mathit{div}_{\mathcal{C}(F)}^* : \mathcal{C} \times \mathcal{C} \rightarrow \mathcal{C}$$

$\mathit{div}_{\mathcal{C}(F)}^*(x, y)$  returns a close approximation to  $x/y$  in  $\mathcal{C}$  with maximum error  $\mathit{max\_error\_div}_{\mathcal{C}(F)}$ .

Further requirement on the  $\mathit{div}_{\mathcal{C}(F)}^*$  approximation helper function are:

$$\begin{aligned} \mathit{div}_{\mathcal{C}(F)}^*(-z, z') &= -\mathit{div}_{\mathcal{C}(F)}^*(z, z') && \text{if } z, z' \in \mathcal{C} \text{ and } z' \neq 0 \\ \mathit{div}_{\mathcal{C}(F)}^*(z, -z') &= -\mathit{div}_{\mathcal{C}(F)}^*(z, z') && \text{if } z, z' \in \mathcal{C} \text{ and } z' \neq 0 \\ \mathit{div}_{\mathcal{C}(F)}^*(\mathit{conj}(z), \mathit{conj}(z')) &= \mathit{conj}(\mathit{div}_{\mathcal{C}(F)}^*(z, z')) && \text{if } z, z' \in \mathcal{C} \text{ and } z' \neq 0 \end{aligned}$$

The  $\mathit{div}_{\mathcal{C}(F)}$  operation:

$$\begin{aligned} \mathit{div}_{\mathcal{C}(F)} : \mathcal{C}(F) \times \mathcal{C}(F) &\rightarrow \mathcal{C}(F) \cup \{-\mathbf{0}\dots, \mathbf{underflow}, \mathbf{overflow}, \mathbf{infinitary}, \mathbf{invalid}\} \\ \mathit{div}_{\mathcal{C}(F)}(x + \hat{\mathbf{i}} \cdot x', y + \hat{\mathbf{i}} \cdot y') &= \mathit{result}_{\mathcal{C}(F)}^*(\mathit{div}_{\mathcal{C}(F)}^*(x + \hat{\mathbf{i}} \cdot x', y + \hat{\mathbf{i}} \cdot y')) \\ &= ??? && \text{if } x + \hat{\mathbf{i}} \cdot x', y + \hat{\mathbf{i}} \cdot y' \in \mathcal{C}(F) \text{ and } |y| \neq |y'| \\ & && \text{otherwise} \end{aligned}$$

### 5.3 Elementary transcendental complex floating point operations

#### 5.3.1 Operations for exponentiations and logarithms

There shall be two maximum error parameters for complex exponentiations and logarithms.

$$\begin{aligned} \mathit{max\_error\_exp}_{\mathcal{C}(F)} &\in F \\ \mathit{max\_error\_power}_{\mathcal{C}(F)} &\in F \end{aligned}$$

no monotonicity requirements? no sign requirements? dependency on argument values?

##### 5.3.1.1 Natural exponentiation

The  $\mathit{exp}_{\mathcal{C}(F)}^*$  approximation helper function:

$$\mathit{exp}_{\mathcal{C}(F)}^* : \mathcal{C} \rightarrow \mathcal{C}$$

$\mathit{exp}_{\mathcal{C}(F)}^*(z)$  returns a close approximation to  $e^z$  in  $\mathcal{C}$  with maximum error  $\mathit{max\_error\_exp}_{\mathcal{C}(F)}$ .

A further requirement on the  $\mathit{exp}_{\mathcal{C}(F)}^*$  approximation helper function is:

$$\mathit{exp}_{\mathcal{C}(F)}^*(\mathit{conj}(z)) = \mathit{conj}(\mathit{exp}_{\mathcal{C}(F)}^*(z)) \quad \text{if } z \in \mathcal{C}$$

The relationship to the  $\mathit{cos}_F^*$ ,  $\mathit{sin}_F^*$ , and  $\mathit{exp}_F^*$  approximation helper functions in an associated library for real-valued operations shall be:

$$\begin{aligned} \mathit{exp}_{\mathcal{C}(F)}^*(\hat{\mathbf{i}} \cdot x') &= \mathit{cos}_F^*(x') + \hat{\mathbf{i}} \cdot \mathit{sin}_F^*(x') && \text{if } x' \in \mathcal{R} \\ \mathit{exp}_{\mathcal{C}(F)}^*(x) &= \mathit{exp}_F^*(x) && \text{if } x \in \mathcal{R} \\ &\dots\text{cyclic rep. in general...} \end{aligned}$$

The  $\mathit{exp}_{\mathcal{C}(F)}$  operation:

$$\begin{aligned} \mathit{exp}_{\mathcal{C}(F)} : \mathcal{C}(F) &\rightarrow \mathcal{C}(F) \cup \{\mathbf{underflow}, \mathbf{overflow}, \mathbf{absolute\_precision\_underflow}\} \\ \mathit{exp}_{\mathcal{C}(F)}(x + \hat{\mathbf{i}} \cdot x') &= \mathit{result}_{\mathcal{C}(F)}^*(\mathit{exp}_{\mathcal{C}(F)}^*(x + \hat{\mathbf{i}} \cdot x'), \mathit{nearest}_F) \\ &= \mathit{exp}_{\mathcal{C}(F)}(0 + \hat{\mathbf{i}} \cdot x') && \text{if } x + \hat{\mathbf{i}} \cdot x' \in \mathcal{C}(F) \text{ and } |x'| \leq \mathit{big\_angle\_r}_F \\ &= \mathit{exp}_{\mathcal{C}(F)}(0 + \hat{\mathbf{i}} \cdot 0) && \text{if } x = -\mathbf{0} \\ &= \mathit{mul}_{\mathcal{C}(F)}(\mathit{exp}_{\mathcal{C}(F)}(x + \hat{\mathbf{i}} \cdot 0)) && \text{if } x' = -\mathbf{0} \text{ and } x \neq -\mathbf{0} \\ &= \mathit{mul}_F(0, \mathit{cos}_F(x')) + \hat{\mathbf{i}} \cdot \mathit{mul}_F(0, \mathit{sin}_F(x')) \end{aligned}$$

$$\begin{aligned}
& \text{if } x = -\infty \text{ and } x' \in F \text{ and } |x'| \leq \text{big\_angle\_}x_F \\
& = \text{mul}_F(+\infty, \cos_F(x')) + \hat{\mathbf{i}} \cdot \text{mul}_F(+\infty, \sin_F(x')) \\
& \text{if } x = +\infty \text{ and } x' \in F \text{ and } |x'| \leq \text{big\_angle\_}x_F \text{ and } x' \neq 0 \\
& = +\infty + \hat{\mathbf{i}} \cdot 0 \\
& \text{if } x = +\infty \text{ and } x' \in F \text{ and } x' = 0 \\
& = \text{radh}_{c(F)}(x + \hat{\mathbf{i}} \cdot x') \quad \text{otherwise}
\end{aligned}$$

NOTE –  $\text{radh}_{c(F)}$  is specified in clause 5.3.3.1.

### 5.3.1.2 Complex exponentiation of argument base

The  $\text{power}_{c(F)}^*$  approximation helper function:

$$\text{power}_{c(F)}^* : \mathcal{C} \times \mathcal{C} \rightarrow \mathcal{C}$$

$\text{power}_{c(F)}^*(b, z)$  returns a close approximation to  $b^z$  in  $\mathcal{C}$  with maximum error  $\text{max\_error\_power}_{c(F)}$ .

A further requirement on the  $\text{power}_{c(F)}^*$  approximation helper function is:

$$\begin{aligned}
\text{power}_{c(F)}^*(\text{conj}(b), \text{conj}(z)) &= \text{conj}(\text{power}_{c(F)}^*(b, z)) \\
&\text{if } b, z \in \mathcal{C}
\end{aligned}$$

The  $\text{power}_{c(F)}$  operation:

$$\begin{aligned}
\text{power}_{c(F)} : c(F) \times c(F) &\rightarrow c(F) \cup \{\text{underflow, overflow, absolute\_precision\_underflow, invalid}\} \\
\text{power}_{c(F)}(x + \hat{\mathbf{i}} \cdot x', y + \hat{\mathbf{i}} \cdot y') &= \text{result}_{c(F)}^*(\text{power}_{c(F)}^*(x + \hat{\mathbf{i}} \cdot x', y + \hat{\mathbf{i}} \cdot y'), \text{nearest}_F) \\
&\text{if } x + \hat{\mathbf{i}} \cdot x', y + \hat{\mathbf{i}} \cdot y' \in c(F) \text{ and } x \neq 0 \text{ and } \\
&\quad |y' \cdot \ln(|x|)| \text{ is not too large...?} \\
&= \text{power}_{c(F)}(0 + \hat{\mathbf{i}} \cdot x', y + \hat{\mathbf{i}} \cdot y') \\
&\text{if } x = -\mathbf{0} \text{ (?) } \\
&= \text{power}_{c(F)}(x + \hat{\mathbf{i}} \cdot x', 0 + \hat{\mathbf{i}} \cdot y') \\
&\text{if } x' = -\mathbf{0} \text{ and } x \neq -\mathbf{0} \\
&= \text{power}_{c(F)}(x + \hat{\mathbf{i}} \cdot x', 0 + \hat{\mathbf{i}} \cdot y') \\
&\text{if } y = -\mathbf{0} \\
&= \text{conj}_{c(F)}(\text{power}_{c(F)}(x + \hat{\mathbf{i}} \cdot x', y + \hat{\mathbf{i}} \cdot \mathbf{0})) \\
&\text{if } y' = -\mathbf{0} \text{ (?) } \\
&= \text{exp}_{c(F)}(\text{mul}_{c(F)}(\ln_{c(F)}(x + \hat{\mathbf{i}} \cdot x'), y + \hat{\mathbf{i}} \cdot y')) \\
&\text{otherwise}
\end{aligned}$$

NOTE – Complex raising to a power is multi-valued. The principal result is given by  $b^q = e^{q \cdot \ln(b)}$ . The  $b^q$  function branch cuts at  $\{x \mid x \in \mathcal{R} \text{ and } x < 0\} \times \mathcal{C}$  (except when  $q$  is in  $\mathcal{Z}$ ). Thus  $\text{power}_{c(F)}(x + \hat{\mathbf{i}} \cdot \mathbf{0}, y) \neq \text{power}_{c(F)}(x + \hat{\mathbf{i}} \cdot (-\mathbf{0}), y)$ .

### 5.3.1.3 Complex square root

The  $\text{sqrt}_{c(F)}^*$  approximation helper function:

$$\text{srt}_{c(F)}^* : \mathcal{C} \rightarrow \mathcal{C}$$

$\text{srt}_{c(F)}^*(z)$  returns a close approximation to  $\sqrt{z}$  in  $\mathcal{C}$  with maximum error  $\text{max\_error\_exp}_{c(F)}$ .

Further requirements on the  $\text{srt}_{c(F)}^*$  approximation helper function are:

$$\begin{aligned}
\text{srt}_{c(F)}^*(\text{conj}(z)) &= \text{conj}(\text{srt}_{c(F)}^*(z)) && \text{if } z \in \mathcal{C} \\
\text{srt}_{c(F)}^*(x) &= \sqrt{x} && \text{if } x \in \mathcal{R} \text{ and } x \geq 0 \\
\text{srt}_{c(F)}^*(x) &= \hat{\mathbf{i}} \cdot \text{srt}_{c(F)}^*(-x) && \text{if } x \in \mathcal{R} \text{ and } x < 0 \\
\Re(\text{srt}_{c(F)}^*(\hat{\mathbf{i}} \cdot x')) &= \Im(\text{srt}_{c(F)}^*(\hat{\mathbf{i}} \cdot x')) && \text{if } x' \in \mathcal{R} \text{ and } x' \geq 0
\end{aligned}$$

The  $\text{srt}_{c(F)}$  operation:



$$\begin{aligned}
& \mathit{sqr}t_{c(F)} : c(F) \rightarrow c(F) \\
& \mathit{sqr}t_{c(F)}(x + \hat{\mathbf{i}} \cdot x') \\
& \quad = \mathit{result}_{c(F)}^*(\mathit{sqr}t_{c(F)}^*(x + \hat{\mathbf{i}} \cdot x'), \mathit{nearest}_F) \\
& \quad \quad \quad \text{if } x + \hat{\mathbf{i}} \cdot x' \in c(F) \\
& \quad = \mathit{sqr}t_{c(F)}(0 + \hat{\mathbf{i}} \cdot x') \quad \text{if } x = -\mathbf{0} \text{ and } x' \in F \cup \{-\infty, +\infty\} \\
& \quad = -\mathbf{0} + \hat{\mathbf{i}} \cdot (-\mathbf{0}) \quad \text{if } x = \mathit{negz} \text{ and } x' = -\mathbf{0} \\
& \quad = \mathit{conj}_{c(F)}(\mathit{sqr}t_{c(F)}(x + \hat{\mathbf{i}} \cdot 0)) \\
& \quad \quad \quad \text{if } x \in F \cup \{-\infty, +\infty\} \text{ and } x' = -\mathbf{0} \\
& \quad = +\infty + \hat{\mathbf{i}} \cdot (+\infty) \quad \text{if } x \in F \cup \{-\infty, -\mathbf{0}, +\infty\} \text{ and } x' = +\infty \\
& \quad = +\infty + \hat{\mathbf{i}} \cdot 0 \quad \text{if } x = +\infty \text{ and } x' \in F \text{ and } x \geq 0 \\
& \quad = +\infty + \hat{\mathbf{i}} \cdot (-\mathbf{0}) \quad \text{if } x = +\infty \text{ and } ((x' \in F \text{ and } x' < 0 \text{ or } x' = -\mathbf{0}) \\
& \quad = +\infty + \hat{\mathbf{i}} \cdot (-\infty) \quad \text{if } x \in F \cup \{-\infty, -\mathbf{0}, +\infty\} \text{ and } x' = -\infty \\
& \quad = 0 + \hat{\mathbf{i}} \cdot (+\infty) \quad \text{if } x = -\infty \text{ and } x' \in F \text{ and } x \geq 0 \\
& \quad = 0 + \hat{\mathbf{i}} \cdot (-\infty) \quad \text{if } x = -\infty \text{ and } ((x' \in F \text{ and } x' < 0 \text{ or } x' = -\mathbf{0}) \\
& \quad = \mathit{no\_result}_{c(F)}(x + \hat{\mathbf{i}} \cdot x') \\
& \quad \quad \quad \text{otherwise}
\end{aligned}$$

NOTE – The inverse of complex square is multi-valued. The principal result is given by  $\sqrt{b} = e^{0.5 \cdot \ln(b)}$ . The  $\sqrt{\phantom{x}}$  function branch cuts at  $\{x \mid x \in \mathcal{R} \text{ and } x < 0\}$ . Thus  $\mathit{sqr}t_{c(F)}(x + \hat{\mathbf{i}} \cdot 0) \neq \mathit{sqr}t_{c(F)}(x + \hat{\mathbf{i}} \cdot (-\mathbf{0}))$  when  $x < 0$ .

### 5.3.1.4 Natural logarithm

The  $\mathit{ln}_{c(F)}^*$  approximation helper function:

$$\mathit{ln}_{c(F)}^* : \mathcal{C} \rightarrow \mathcal{C}$$

$\mathit{ln}_{c(F)}^*(z)$  returns a close approximation to  $\ln(z)$  in  $\mathcal{C}$  with maximum error  $\mathit{max\_error\_exp}_{c(F)}$ .

A further requirement on the  $\mathit{ln}_{c(F)}^*$  approximation helper function is:

$$\mathit{ln}_{c(F)}^*(\mathit{conj}(z)) = \mathit{conj}(\mathit{ln}_{c(F)}^*(z)) \quad \text{if } z \in \mathcal{C}$$

The relationship to the  $\mathit{arc}_F^*$  and  $\mathit{ln}_F^*$  approximation helper functions in an associated library for real-valued operations shall be:

$$\begin{aligned}
\Im(\mathit{ln}_{c(F)}^*(x + \hat{\mathbf{i}} \cdot x')) &= \mathit{arc}_F^*(x, x') & \text{if } x, x' \in \mathcal{R} \\
\Re(\mathit{ln}_{c(F)}^*(x + \hat{\mathbf{i}} \cdot x')) &= \mathit{ln}_F^*(|x + \hat{\mathbf{i}} \cdot x'|) & \text{if } x, x' \in \mathcal{R} \text{ and } (x = 0 \text{ or } x' = 0)
\end{aligned}$$

The  $\mathit{ln}_{c(F)}$  operation:

$$\begin{aligned}
& \mathit{ln}_{c(F)} : c(F) \rightarrow c(F) \cup \{\mathbf{infinitary}\} \\
& \mathit{ln}_{c(F)}(x + \hat{\mathbf{i}} \cdot x') = \mathit{result}_{c(F)}^*(\mathit{ln}_{c(F)}^*(x + \hat{\mathbf{i}} \cdot x'), \mathit{nearest}_F) \\
& \quad \quad \quad \text{if } x + \hat{\mathbf{i}} \cdot x' \in c(F) \text{ and } (x \neq 0 \text{ or } x' \neq 0) \\
& \quad = \mathbf{infinitary}(-\infty + \hat{\mathbf{i}} \cdot \mathit{arc}_F(x, x')) \\
& \quad \quad \quad \text{if } x, x' \in \{-\mathbf{0}, 0\} \\
& \quad = \mathit{conj}_{c(F)}(\mathit{ln}_{c(F)}(x + \hat{\mathbf{i}} \cdot 0)) \\
& \quad \quad \quad \text{if } x' = -\mathbf{0} \\
& \quad = \mathit{ln}_{c(F)}(0 + \hat{\mathbf{i}} \cdot x') \quad \text{if } x = -\mathbf{0} \text{ and } x \in F \text{ and } x' \neq 0 \\
& \quad = +\infty + \hat{\mathbf{i}} \cdot \mathit{arc}_F(x, x') \quad \text{if } x \in \{-\infty, +\infty\} \text{ and } x' \in F \cup \{-\infty, +\infty\} \\
& \quad = +\infty + \hat{\mathbf{i}} \cdot \mathit{arc}_F(x, x') \quad \text{if } x \in F \text{ and } x' \in \{-\infty, +\infty\} \\
& \quad = \mathit{no\_result}_{c(F)}(x + \hat{\mathbf{i}} \cdot x') \\
& \quad \quad \quad \text{otherwise}
\end{aligned}$$

NOTES

- 1 The inverse of natural exponentiation is multi-valued: the imaginary part may have any integer multiple of  $2 \cdot \pi$  added to it, and the result is also in the solution set. The  $\ln$  function (returning the principle value for the inverse) branch cuts at  $\{x \mid x \in \mathcal{R} \text{ and } x < 0\}$ , is continuous on the rest of  $\mathcal{C}$ , and  $\ln(z) \in \mathcal{R}$  if  $x \in \mathcal{R}$  and  $x > 0$ . Thus  $\ln_{c(F)}(x + \hat{\mathbf{i}} \cdot 0) \neq \ln_{c(F)}(x + \hat{\mathbf{i}} \cdot (-0))$  when  $x < 0$ .
- 2  $re_{c(F)}(\ln_{c(F)}(x + \hat{\mathbf{i}} \cdot x')) \approx \ln_F(\text{hypot}_F(x, x'))$  and  $im_{c(F)}(\ln_{c(F)}(x + \hat{\mathbf{i}} \cdot x')) = \text{arc}_F(x, x')$  when there is no notification.

### 5.3.1.5 Argument base logarithm

The  $\text{logbase}_{c(F)}^*$  approximation helper function:

$$\text{logbase}_{c(F)}^* : \mathcal{C} \times \mathcal{C} \rightarrow \mathcal{C}$$

$\text{logbase}_{c(F)}^*(b, z)$  returns a close approximation to  $\log_b(z)$  in  $\mathcal{C}$  with maximum error  $\text{max\_error\_power}_{c(F)}$ .

A further requirement on the  $\text{logbase}_{c(F)}^*$  approximation helper function is:

$$\begin{aligned} \text{logbase}_{c(F)}^*(\text{conj}(b), \text{conj}(z)) &= \text{conj}(\text{logbase}_{c(F)}^*(b, z)) \\ &\text{if } b, z \in \mathcal{C} \end{aligned}$$

The  $\text{logbase}_{c(F)}$  operation:

$$\begin{aligned} \text{logbase}_{c(F)} : c(F) \times c(F) &\rightarrow c(F) \cup \{\mathbf{infinitary}, \mathbf{invalid}\} \\ \text{logbase}_{c(F)}(x + \hat{\mathbf{i}} \cdot x', y + \hat{\mathbf{i}} \cdot y') &= \text{result}_{c(F)}^*(\text{logbase}_{c(F)}^*(x + \hat{\mathbf{i}} \cdot x', y + \hat{\mathbf{i}} \cdot y'), \text{nearest}_F) \\ &\quad \text{if } x + \hat{\mathbf{i}} \cdot x', y + \hat{\mathbf{i}} \cdot y' \in c(F) \text{ and } x \neq 0 \text{ and} \\ &\quad \quad |\Re(\ln(x + \hat{\mathbf{i}} \cdot x'))| \neq |\Im(\ln(x + \hat{\mathbf{i}} \cdot x'))| \\ &= \text{div}_{c(F)}(\ln_{c(F)}(y + \hat{\mathbf{i}} \cdot y'), \ln_{c(F)}(x + \hat{\mathbf{i}} \cdot x')) \\ &\quad \text{otherwise} \end{aligned}$$

NOTE – Complex logarithm with argument base is multi-valued. The principal result is given by  $\log_b(q) = \ln(q)/\ln(b)$ . Apart from the poles, the  $\log_b(q)$  function branch cuts at  $(\{x \mid x \in \mathcal{R} \text{ and } x < 0\} \times \mathcal{C}) \cup (\mathcal{C} \times \{x \mid x \in \mathcal{R} \text{ and } x < 0\})$ .

### 5.3.2 Operations for radian trigonometric elementary functions

There shall be two maximum error parameters for complex trigonometric operations.

$$\text{max\_error\_sin}_{c(F)} \in F$$

$$\text{max\_error\_tan}_{c(F)} \in F$$

no monotonicity requirements? no sign requirements? dependency on argument values?

#### 5.3.2.1 Radian angle normalisation

$$\text{rad}_{i(F)} : i(F) \rightarrow i(F) \cup \{\mathbf{absolute\_precision\_underflow}\}$$

$$\text{rad}_{i(F)}(\hat{\mathbf{i}} \cdot x') = \hat{\mathbf{i}} \cdot x'$$

$$\text{rad}_{c(F)} : c(F) \rightarrow c(F) \cup \{\mathbf{absolute\_precision\_underflow}\}$$

$$\begin{aligned} \text{rad}_{c(F)}(x + \hat{\mathbf{i}} \cdot x') &= \text{rad}_F(x) + \hat{\mathbf{i}} \cdot x' && \text{if } x' \in F \cup \{-\infty, -0, +\infty\} \text{ and } \text{rad}_F(x) \in F \cup \{-0\} \\ &= \mathbf{absolute\_precision\_underflow}(\text{qNaN} + \hat{\mathbf{i}} \cdot \text{qNaN}) && \text{if } x' \in F \cup \{-\infty, -0, +\infty\} \text{ and} \\ &\quad \text{rad}_F(x) = \mathbf{absolute\_precision\_underflow} \\ &= \text{no\_result}_{c(F)}(x + \hat{\mathbf{i}} \cdot x') \text{ otherwise} \end{aligned}$$

### 5.3.2.2 Radian sine

The  $\sin_{i(F)}$  operation:

$$\begin{aligned}\sin_{i(F)} &: i(F) \rightarrow i(F) \cup \{\mathbf{overflow}\} \\ \sin_{i(F)}(\hat{\mathbf{i}} \cdot x') &= \hat{\mathbf{i}} \cdot \sinh_F(x')\end{aligned}$$

The  $\sin_{c(F)}^*$  approximation helper function:

$$\sin_{c(F)}^* : \mathcal{C} \rightarrow \mathcal{C}$$

$\sin_{c(F)}^*(z)$  returns a close approximation to  $\sin(z)$  in  $\mathcal{C}$  with maximum error  $max\_error\_sin_{c(F)}$ .

Further requirements on the  $\sin_{c(F)}^*$  approximation helper function are:

$$\begin{aligned}\sin_{c(F)}^*(\mathit{conj}(z)) &= \mathit{conj}(\sin_{c(F)}^*(z)) && \text{if } z \in \mathcal{C} \\ \sin_{c(F)}^*(-z) &= -\sin_{c(F)}^*(z) && \text{if } z \in \mathcal{C}\end{aligned}$$

The relationship to the  $\sin_F^*$  and  $\sinh_F^*$  approximation helper functions in an associated library for real-valued operations shall be:

$$\begin{aligned}\sin_{c(F)}^*(x) &= \sin_F^*(x) && \text{if } x \in \mathcal{R} \\ \sin_{c(F)}^*(\hat{\mathbf{i}} \cdot x') &= \hat{\mathbf{i}} \cdot \sinh_F^*(x') && \text{if } x' \in \mathcal{R}\end{aligned}$$

The requirements implied by these relationships and the requirements from part 2 (GENER-ALISE?) remain even if there is no  $\sin_F$  or  $\sinh_F$  operations in any associated library for real-valued operations or there is no associated library for real-valued operations.

The  $\sin_{c(F)}$  operation:

$$\begin{aligned}\sin_{c(F)} &: c(F) \rightarrow c(F) \cup \{\mathbf{underflow}, \mathbf{overflow}, \mathbf{absolute\_precision\_underflow}\} \\ \sin_{c(F)}(x + \hat{\mathbf{i}} \cdot x') &= \mathit{result}_{c(F)}^*(\sin_{c(F)}^*(x + \hat{\mathbf{i}} \cdot x'), \mathit{nearest}_F) && \text{if } x + \hat{\mathbf{i}} \cdot x' \in c(F) \text{ and } |x| \leq \mathit{big\_angle}_F \\ &= \mathit{conj}_{c(F)}(\sin_{c(F)}(x + \hat{\mathbf{i}} \cdot 0)) && \text{if } x' = -\mathbf{0} \\ &= \mathit{neg}_{c(F)}(\sin_{c(F)}(0 + \hat{\mathbf{i}} \cdot \mathit{neg}_F(x'))) && \text{if } x = -\mathbf{0} \text{ and } x' \neq -\mathbf{0} \\ &= 0 + \hat{\mathbf{i}} \cdot (+\infty) && \text{if } x = 0 \text{ and } x' = +\infty \\ &= 0 + \hat{\mathbf{i}} \cdot (-\infty) && \text{if } x = 0 \text{ and } x' = -\infty \\ &= \mathit{mul}_F(+\infty, \sin_F(x)) + \hat{\mathbf{i}} \cdot \mathit{mul}_F(+\infty, \cos_F(x)) && \text{if } x' = +\infty \text{ and } x \notin \{-\mathbf{0}, 0\} \\ &= \mathit{mul}_F(+\infty, \sin_F(x)) + \hat{\mathbf{i}} \cdot \mathit{mul}_F(-\infty, \cos_F(x)) && \text{if } x' = -\infty \text{ and } x \notin \{-\mathbf{0}, 0\} \\ &= \mathit{rad}_{c(F)}(x + \hat{\mathbf{i}} \cdot x') && \text{otherwise}\end{aligned}$$

### 5.3.2.3 Radian cosine

The  $\cos_{i(F)}$  operation:

$$\begin{aligned}\cos_{i(F)} &: i(F) \rightarrow F \cup \{\mathbf{overflow}\} \\ \cos_{i(F)}(\hat{\mathbf{i}} \cdot x') &= \cosh_F(x')\end{aligned}$$

The  $\cos_{c(F)}^*$  approximation helper function:

$$\cos_{c(F)}^* : \mathcal{C} \rightarrow \mathcal{C}$$

$\cos_{c(F)}^*(z)$  returns a close approximation to  $\cos(z)$  in  $\mathcal{C}$  with maximum error  $max\_error\_sin_{c(F)}$ .

Further requirements on the  $\cos_{c(F)}^*$  approximation helper function are:

$$\begin{aligned} \cos_{c(F)}^*(\text{conj}(z)) &= \text{conj}(\cos_{c(F)}^*(z)) && \text{if } z \in \mathcal{C} \\ \cos_{c(F)}^*(-z) &= \cos_{c(F)}^*(z) && \text{if } z \in \mathcal{C} \end{aligned}$$

The relationship to the  $\cos_F^*$  and  $\cosh_F^*$  approximation helper functions in an associated library for real-valued operations shall be:

$$\begin{aligned} \cos_{c(F)}^*(x) &= \cos_F^*(x) && \text{if } x \in \mathcal{R} \\ \cos_{c(F)}^*(\hat{\mathbf{i}} \cdot x') &= \cosh_F^*(x') && \text{if } x' \in \mathcal{R} \end{aligned}$$

The  $\cos_{c(F)}$  operation:

$$\begin{aligned} \cos_{c(F)} : c(F) &\rightarrow c(F) \cup \{\mathbf{underflow}, \mathbf{overflow}, \mathbf{absolute\_precision\_underflow}\} \\ \cos_{c(F)}(x + \hat{\mathbf{i}} \cdot x') &= \text{result}_{c(F)}^*(\cos_{c(F)}^*(x + \hat{\mathbf{i}} \cdot x'), \text{nearest}_F) && \text{if } x + \hat{\mathbf{i}} \cdot x' \in c(F) \text{ and } |x| \leq \text{big\_angle\_r}_F \\ &= \text{conj}_{c(F)}(\cos_{c(F)}(x + \hat{\mathbf{i}} \cdot 0)) && \text{if } x' = -\mathbf{0} \\ &= \cos_{c(F)}(0 + \hat{\mathbf{i}} \cdot \text{neg}_F(x')) && \text{if } x = -\mathbf{0} \text{ and } x' \neq -\mathbf{0} \\ &= +\infty + \hat{\mathbf{i}} \cdot (-\mathbf{0}) && \text{if } x' = +\infty \text{ and } x = 0 \\ &= +\infty + \hat{\mathbf{i}} \cdot 0 && \text{if } x' = -\infty \text{ and } x = 0 \\ &= \text{mul}_F(+\infty, \cos_F(x)) + \hat{\mathbf{i}} \cdot \text{mul}_F(-\infty, \sin_F(x)) && \text{if } x' = +\infty \text{ and } x \notin \{-\mathbf{0}, 0\} \\ &= \text{mul}_F(+\infty, \cos_F(x)) + \hat{\mathbf{i}} \cdot \text{mul}_F(+\infty, \sin_F(x)) && \text{if } x' = -\infty \text{ and } x \notin \{-\mathbf{0}, 0\} \\ &= \text{rad}_{c(F)}(x + \hat{\mathbf{i}} \cdot x') && \text{otherwise} \end{aligned}$$

### 5.3.2.4 Radian tangent

The  $\tan_{i(F)}$  operation:

$$\begin{aligned} \tan_{i(F)} : i(F) &\rightarrow i(F) \cup \{\mathbf{overflow}\} \\ \tan_{i(F)}(\hat{\mathbf{i}} \cdot x') &= \hat{\mathbf{i}} \cdot \tanh_F(x') \end{aligned}$$

The  $\tan_{c(F)}^*$  approximation helper function:

$$\tan_{c(F)}^* : \mathcal{C} \rightarrow \mathcal{C}$$

$\tan_{c(F)}^*(z)$  returns a close approximation to  $\tan(z)$  in  $\mathcal{C}$  with maximum error  $\text{max\_error\_tan}_{c(F)}$ .

Further requirements on the  $\tan_{c(F)}^*$  approximation helper function are:

$$\begin{aligned} \tan_{c(F)}^*(\text{conj}(z)) &= \text{conj}(\tan_{c(F)}^*(z)) && \text{if } z \in \mathcal{C} \\ \tan_{c(F)}^*(-z) &= -\tan_{c(F)}^*(z) && \text{if } z \in \mathcal{C} \end{aligned}$$

The relationship to the  $\tan_F^*$  and  $\tanh_F^*$  approximation helper functions in an associated library for real-valued operations shall be:

$$\begin{aligned} \tan_{c(F)}^*(x) &= \tan_F^*(x) && \text{if } x \in \mathcal{R} \\ \tan_{c(F)}^*(\hat{\mathbf{i}} \cdot x') &= \hat{\mathbf{i}} \cdot \tanh_F^*(x') && \text{if } x' \in \mathcal{R} \end{aligned}$$

The  $\tan_{c(F)}$  operation:

$$\begin{aligned} \tan_{c(F)} : c(F) &\rightarrow c(F) \cup \{\mathbf{underflow}, \mathbf{overflow}, \mathbf{absolute\_precision\_underflow}\} \\ \tan_{c(F)}(x + \hat{\mathbf{i}} \cdot x') &= \text{result}_{c(F)}^*(\tan_{c(F)}^*(x + \hat{\mathbf{i}} \cdot x'), \text{nearest}_F) && \text{if } x + \hat{\mathbf{i}} \cdot x' \in c(F) \text{ and } |x| \leq \text{big\_angle\_r}_F \\ &= \text{conj}_{c(F)}(\tan_{c(F)}(x + \hat{\mathbf{i}} \cdot 0)) && \text{if } x' = -\mathbf{0} \end{aligned}$$

$$\begin{aligned}
&= \mathit{neg}_F(\mathit{tan}_{c(F)}(0 + \hat{\mathbf{i}} \cdot \mathit{neg}_F(x'))) \\
&\quad \text{if } x = -\mathbf{0} \text{ and } x' \neq -\mathbf{0} \\
&= \mathit{mul}_F(0, \mathit{tan}_F(x)) + \hat{\mathbf{i}} \cdot 1 \\
&\quad \text{if } x' = +\infty \\
&= \mathit{mul}_F(0, \mathit{tan}_F(x)) + \hat{\mathbf{i}} \cdot (-1) \\
&\quad \text{if } x' = -\infty \\
&= \mathit{rad}_{c(F)}(x + \hat{\mathbf{i}} \cdot x') \quad \text{otherwise}
\end{aligned}$$

### 5.3.2.5 Radian cotangent

The  $\mathit{cot}_{i(F)}$  operation:

$$\begin{aligned}
\mathit{cot}_{i(F)} : i(F) &\rightarrow i(F) \cup \{\mathbf{underflow}, \mathbf{infinitary}\} \\
\mathit{cot}_{i(F)}(\hat{\mathbf{i}} \cdot x') &= \mathit{neg}_{i(F)}(\hat{\mathbf{i}} \cdot \mathit{coth}_F(x'))
\end{aligned}$$

The  $\mathit{cot}_{c(F)}^*$  approximation helper function:

$$\mathit{cot}_{c(F)}^* : \mathcal{C} \rightarrow \mathcal{C}$$

$\mathit{cot}_{c(F)}^*(z)$  returns a close approximation to  $\mathit{cot}(z)$  in  $\mathcal{C}$  with maximum error  $\mathit{max\_error\_tan}_{c(F)}$ .

Further requirements on the  $\mathit{cot}_{c(F)}^*$  approximation helper function are:

$$\begin{aligned}
\mathit{cot}_{c(F)}^*(\mathit{conj}(z)) &= \mathit{conj}(\mathit{cot}_{c(F)}^*(z)) && \text{if } z \in \mathcal{C} \\
\mathit{cot}_{c(F)}^*(-z) &= -\mathit{cot}_{c(F)}^*(z) && \text{if } z \in \mathcal{C}
\end{aligned}$$

The relationship to the  $\mathit{cot}_F^*$  and  $\mathit{coth}_F^*$  approximation helper functions in an associated library for real-valued operations shall be:

$$\begin{aligned}
\mathit{cot}_{c(F)}^*(x) &= \mathit{cot}_F^*(x) && \text{if } x \in \mathcal{R} \\
\mathit{cot}_{c(F)}^*(\hat{\mathbf{i}} \cdot x') &= -\hat{\mathbf{i}} \cdot \mathit{coth}_F^*(x') && \text{if } x' \in \mathcal{R}
\end{aligned}$$

The  $\mathit{cot}_{c(F)}$  operation:

$$\begin{aligned}
\mathit{cot}_{c(F)} : c(F) &\rightarrow c(F) \cup \{\mathbf{underflow}, \mathbf{overflow}, \mathbf{infinitary}, \mathbf{absolute\_precision\_underflow}\} \\
\mathit{cot}_{c(F)}(x + \hat{\mathbf{i}} \cdot x') &= \mathit{result}_{c(F)}^*(\mathit{cot}_{c(F)}^*(x + \hat{\mathbf{i}} \cdot x'), \mathit{nearest}_F) \\
&\quad \text{if } x + \hat{\mathbf{i}} \cdot x' \in c(F) \text{ and } |x| \leq \mathit{big\_angle\_r}_F \text{ and} \\
&\quad \quad (x \neq 0 \text{ or } x' \neq 0) \\
&= \mathit{conj}_{c(F)}(\mathit{tan}_{c(F)}(x + \hat{\mathbf{i}} \cdot 0)) \\
&\quad \text{if } x' = -\mathbf{0} \\
&= \mathit{neg}_F(\mathit{tan}_{c(F)}(0 + \hat{\mathbf{i}} \cdot \mathit{neg}_F(x'))) \\
&\quad \text{if } x = -\mathbf{0} \text{ and } x' \neq -\mathbf{0} \\
&= \mathit{mul}_F(0, \mathit{tan}_F(x)) + \hat{\mathbf{i}} \cdot (-1) \\
&\quad \text{if } x' = +\infty \\
&= \mathit{mul}_F(0, \mathit{tan}_F(x)) + \hat{\mathbf{i}} \cdot 1 \\
&\quad \text{if } x' = -\infty \\
&= \mathbf{infinitary}(+\infty + \hat{\mathbf{i}} \cdot (+\infty)) \\
&\quad \text{if } x = 0 \text{ and } x' = 0 \\
&= \mathit{rad}_{c(F)}(x + \hat{\mathbf{i}} \cdot x') \quad \text{otherwise}
\end{aligned}$$

### 5.3.2.6 Radian secant

The  $\mathit{sec}_{i(F)}$  operation:

$$\begin{aligned}
\mathit{sec}_{i(F)} : i(F) &\rightarrow c(F) \cup \{\mathbf{overflow}\} \\
\mathit{sec}_{i(F)}(\hat{\mathbf{i}} \cdot x') &= \mathit{neg}_{i(F)}(\hat{\mathbf{i}} \cdot \mathit{sech}_F(x'))
\end{aligned}$$

The  $sec_{c(F)}^*$  approximation helper function:

$$sec_{c(F)}^* : \mathcal{C} \rightarrow \mathcal{C}$$

$sec_{c(F)}^*(z)$  returns a close approximation to  $sec(z)$  in  $\mathcal{C}$  with maximum error  $max\_error\_tan_{c(F)}$ .

Further requirements on the  $sec_{c(F)}^*$  approximation helper function are:

$$\begin{aligned} sec_{c(F)}^*(conj(z)) &= conj(sec_{c(F)}^*(z)) && \text{if } z \in \mathcal{C} \\ sec_{c(F)}^*(-z) &= -sec_{c(F)}^*(z) && \text{if } z \in \mathcal{C} \end{aligned}$$

The relationship to the  $sec_F^*$  and  $sech_F^*$  approximation helper functions in an associated library for real-valued operations shall be:

$$\begin{aligned} sec_{c(F)}^*(x) &= sec_F^*(x) && \text{if } x \in \mathcal{R} \\ sec_{c(F)}^*(\hat{i} \cdot x') &= -\hat{i} \cdot sech_F(x') && \text{if } x' \in \mathcal{R} \end{aligned}$$

The  $sec_{c(F)}$  operation:

$$\begin{aligned} sec_{c(F)} : c(F) &\rightarrow c(F) \cup \{\mathbf{underflow}, \mathbf{overflow}, \mathbf{absolute\_precision\_underflow}\} \\ sec_{c(F)}(x + \hat{i} \cdot x') &= result_{c(F)}^*(sec_{c(F)}^*(x + \hat{i} \cdot x'), nearest_F) \\ &\quad \text{if } x + \hat{i} \cdot x' \in c(F) \text{ and } |x| \leq big\_angle\_r_F \\ &= conj_{c(F)}(sec_{c(F)}(x + \hat{i} \cdot 0)) \\ &\quad \text{if } x' = -\mathbf{0} \\ &= sec_{c(F)}(0 + \hat{i} \cdot neg_F(x')) \\ &\quad \text{if } x = -\mathbf{0} \text{ and } x' \neq -\mathbf{0} \\ &= mul_F(0, cos_F(x)) + \hat{i} \cdot mul_F(-\mathbf{0}, sin_F(x)) \\ &\quad \text{if } x' = +\infty \\ &= mul_F(0, cos_F(x)) + \hat{i} \cdot mul_F(0, sin_F(x)) \\ &\quad \text{if } x' = -\infty \\ &= rad_{c(F)}(x + \hat{i} \cdot x') \quad \text{otherwise} \end{aligned}$$

### 5.3.2.7 Radian cosecant

The  $csc_{i(F)}$  operation:

$$\begin{aligned} csc_{i(F)} : i(F) &\rightarrow c(F) \cup \{\mathbf{overflow}, \mathbf{infinitary}\} \\ csc_{i(F)}() &= neg_{i(F)}(\hat{i} \cdot csch_F(x')) \end{aligned}$$

The  $csc_{c(F)}^*$  approximation helper function:

$$csc_{c(F)}^* : \mathcal{C} \rightarrow \mathcal{C}$$

$csc_{c(F)}^*(z)$  returns a close approximation to  $csc(z)$  in  $\mathcal{C}$  with maximum error  $max\_error\_tan_{c(F)}$ .

Further requirements on the  $csc_{c(F)}^*$  approximation helper function are:

$$\begin{aligned} csc_{c(F)}^*(conj(z)) &= conj(csc_{c(F)}^*(z)) && \text{if } z \in \mathcal{C} \\ csc_{c(F)}^*(-z) &= -csc_{c(F)}^*(z) && \text{if } z \in \mathcal{C} \end{aligned}$$

The relationship to the  $csc_F^*$  and  $csch_F^*$  approximation helper functions in an associated library for real-valued operations shall be:

$$\begin{aligned} csc_{c(F)}^*(x) &= csc_F^*(x) && \text{if } x \in \mathcal{R} \\ csc_{c(F)}^*(\hat{i} \cdot x') &= -\hat{i} \cdot csch_F(x') && \text{if } x' \in \mathcal{R} \end{aligned}$$

The  $csc_{c(F)}$  operation:

$$csc_{c(F)} : c(F) \rightarrow c(F) \cup \{\mathbf{underflow}, \mathbf{overflow}, \mathbf{infinitary}, \mathbf{absolute\_precision\_underflow}\}$$

$$\begin{aligned}
csc_{c(F)}(x + \hat{\mathbf{i}} \cdot x') &= result_{c(F)}^*(csc_{c(F)}^*(x + \hat{\mathbf{i}} \cdot x'), nearest_F) \\
&\quad \text{if } x + \hat{\mathbf{i}} \cdot x' \in c(F) \text{ and } |x| \leq big\_angle\_r_F \text{ and} \\
&\quad \quad (x \neq 0 \text{ or } x' \neq 0) \\
&= conj_{c(F)}(csc_{c(F)}(x + \hat{\mathbf{i}} \cdot 0)) \\
&\quad \text{if } x' = -\mathbf{0} \\
&= neg_F(csc_{c(F)}(0 + \hat{\mathbf{i}} \cdot neg_F(x'))) \\
&\quad \text{if } x = -\mathbf{0} \text{ and } x' \neq -\mathbf{0} \\
&= mul_F(-\mathbf{0}, sin_F(x)) + \hat{\mathbf{i}} \cdot mul_F(-\mathbf{0}, cos_F(x)) \\
&\quad \text{if } x' = +\infty \\
&= mul_F(-\mathbf{0}, sin_F(x)) + \hat{\mathbf{i}} \cdot mul_F(0, cos_F(x)) \\
&\quad \text{if } x' = -\infty \\
&= \mathbf{infinitary}(+\infty + \hat{\mathbf{i}} \cdot (+\infty)) \\
&\quad \text{if } x = 0 \text{ and } x' = 0 \\
&= rad_{c(F)}(x + \hat{\mathbf{i}} \cdot x') \quad \text{otherwise}
\end{aligned}$$

### 5.3.2.8 Radian arc sine

The  $arcsin_{i(F)}$  operation:

$$\begin{aligned}
arcsin_{i(F)} : i(F) &\rightarrow i(F) \\
arcsin_{i(F)}(\hat{\mathbf{i}} \cdot x') &= \hat{\mathbf{i}} \cdot arcsinh_F(x')
\end{aligned}$$

The  $arcsin_{c(F)}^*$  approximation helper function:

$$arcsin_{c(F)}^* : \mathcal{C} \rightarrow \mathcal{C}$$

$arcsin_{c(F)}^*(z)$  returns a close approximation to  $\arcsin(z)$  in  $\mathcal{C}$  with maximum error  $max\_error\_sin_{c(F)}$ .

Further requirements on the  $arcsin_{c(F)}^*$  approximation helper function are:

$$\begin{aligned}
arcsin_{c(F)}^*(conj(z)) &= conj(arcsin_{c(F)}^*(z)) \quad \text{if } z \in \mathcal{C} \\
arcsin_{c(F)}^*(-z) &= -arcsin_{c(F)}^*(z) \quad \text{if } z \in \mathcal{C} \\
\Re(arcsin_{c(F)}^*(x)) &= -\pi/2 \quad \text{if } x \in \mathcal{R} \text{ and } x \leq -1 \\
\Re(arcsin_{c(F)}^*(x)) &= \pi/2 \quad \text{if } x \in \mathcal{R} \text{ and } x \geq 1
\end{aligned}$$

The relationship to the  $arcsin_F^*$  and  $arcsinh_F^*$  approximation helper functions in an associated library for real-valued operations shall be:

$$\begin{aligned}
arcsin_{c(F)}^*(x) &= arcsin_F^*(x) \quad \text{if } x \in \mathcal{R} \text{ and } |x| \leq 1 \\
arcsin_{c(F)}^*(\hat{\mathbf{i}} \cdot x') &= \hat{\mathbf{i}} \cdot arcsinh_F^*(x') \quad \text{if } x' \in \mathcal{R}
\end{aligned}$$

The  $arcsin_{c(F)}^\#$  range limitation helper function:

$$arcsin_{c(F)}^\#(z) = \max\{up_F(-\pi/2), \min\{\Re(arcsin_F^*(z)), down_F(\pi/2)\}\} + \hat{\mathbf{i}} \cdot \Im(arcsin_F^*(z))$$

The  $arcsin_{c(F)}$  operation:

$$\begin{aligned}
arcsin_{c(F)} : c(F) &\rightarrow c(F) \cup \{\mathbf{underflow}\} \\
arcsin_{c(F)}(x + \hat{\mathbf{i}} \cdot x') &= result_{c(F)}^*(arcsin_{c(F)}^\#(x + \hat{\mathbf{i}} \cdot x'), nearest_F) \\
&\quad \text{if } x + \hat{\mathbf{i}} \cdot x' \in c(F) \\
&= neg_{c(F)}(arcsin_{c(F)}(0 + \hat{\mathbf{i}} \cdot neg_F(x'))) \\
&\quad \text{if } x = -\mathbf{0} \\
&= conj_{c(F)}(arcsin_{c(F)}(x + \hat{\mathbf{i}} \cdot 0)) \\
&\quad \text{if } x' = -\mathbf{0} \text{ and } x \neq -\mathbf{0} \\
&= arc_F(x', x) + \hat{\mathbf{i}} \cdot (+\infty) \quad \text{if } x' = +\infty \text{ and } x \in \{-\infty, -\mathbf{0}, +\infty\}
\end{aligned}$$

$$\begin{aligned}
&= \text{arc}_F(\text{neg}_F(x'), x) + \hat{\mathbf{i}} \cdot (-\infty) && \text{if } x' = -\infty \text{ and } x \in \{-\infty, -\mathbf{0}, +\infty\} \\
&= \text{arc}_F(x', x) + \hat{\mathbf{i}} \cdot (+\infty) && \text{if } x \in \{-\infty, +\infty\} \text{ and } x' \in F \text{ and } x' \geq 0 \\
&= \text{arc}_F(\text{neg}_F(x'), x) + \hat{\mathbf{i}} \cdot (-\infty) && \text{if } x \in \{-\infty, +\infty\} \text{ and } x' \in F \text{ and } x' < 0 \\
&= \text{no\_result}_{c(F)}(x + \hat{\mathbf{i}} \cdot x') && \text{otherwise}
\end{aligned}$$

NOTE – The inverse of sin is multi-valued, the real part may have any integer multiple of  $2 \cdot \pi$  added to it, and the result is also in the solution set. The arcsin function (returning the principal value for the inverse) branch cuts at  $\{x \mid x \in \mathcal{R} \text{ and } |x| > 1\}$ . Thus  $\text{arcsin}_{c(F)}(x + \hat{\mathbf{i}} \cdot \mathbf{0}) \neq \text{arcsin}_{c(F)}(x + \hat{\mathbf{i}} \cdot (-\mathbf{0}))$  when  $|x| > 1$ .

### 5.3.2.9 Radian arc cosine

The  $\text{arccos}_{c(F)}^*$  approximation helper function:

$$\text{arccos}_{c(F)}^* : \mathcal{C} \rightarrow \mathcal{C}$$

$\text{arccos}_{c(F)}^*(z)$  returns a close approximation to  $\text{arccos}(z)$  in  $\mathcal{C}$  with maximum error  $\text{max\_error\_sin}_{c(F)}$ .

Further requirements on the  $\text{arccos}_{c(F)}^*$  approximation helper function are:

$$\begin{aligned}
\text{arccos}_{c(F)}^*(\text{conj}(z)) &= \text{conj}(\text{arccos}_{c(F)}^*(z)) && \text{if } z \in \mathcal{C} \\
\Re(\text{arccos}_{c(F)}^*(x)) &= \pi && \text{if } x \in \mathcal{R} \text{ and } x \leq -1
\end{aligned}$$

The relationship to the  $\text{arccos}_F^*$  and  $\text{arccosh}_F^*$  approximation helper functions in an associated library for real-valued operations shall be:

$$\text{arccos}_{c(F)}^*(x) = \text{arccos}_F^*(x) \quad \text{if } x \in \mathcal{R} \text{ and } |x| \leq 1$$

The  $\text{arccos}_{c(F)}^\#$  range limitation helper function:

$$\text{arccos}_{c(F)}^\#(x + \hat{\mathbf{i}} \cdot x') = \min\{\Re(\text{arccos}_F^*(x + \hat{\mathbf{i}} \cdot x')), \text{down}_F(\pi)\} + \hat{\mathbf{i}} \cdot \Im(\text{arccos}_F^*(x + \hat{\mathbf{i}} \cdot x'))$$

The  $\text{arccos}_{c(F)}$  operation:

$$\text{arccos}_{c(F)} : c(F) \rightarrow c(F)$$

$$\begin{aligned}
&\text{arccos}_{c(F)}(x + \hat{\mathbf{i}} \cdot x') \\
&= \text{result}_{c(F)}^*(\text{arccos}_{c(F)}^\#(x + \hat{\mathbf{i}} \cdot x'), \text{nearest}_F) && \text{if } x + \hat{\mathbf{i}} \cdot x' \in c(F) \\
&= \text{arccos}_{c(F)}(\mathbf{0} + \hat{\mathbf{i}} \cdot \text{neg}_F(x')) && \text{if } x = -\mathbf{0} \\
&= \text{conj}_{c(F)}(\text{arccos}_{c(F)}(x + \hat{\mathbf{i}} \cdot \mathbf{0})) && \text{if } x' = -\mathbf{0} \text{ and } x \neq -\mathbf{0} \\
&= \text{arc}_F(x, x') + \hat{\mathbf{i}} \cdot (-\infty) && \text{if } x' = +\infty \text{ and } x \in \{-\infty, -\mathbf{0}, +\infty\} \\
&= \text{arc}_F(x, \text{neg}_F(x')) + \hat{\mathbf{i}} \cdot (+\infty) && \text{if } x' = -\infty \text{ and } x \in \{-\infty, -\mathbf{0}, +\infty\} \\
&= \text{arc}_F(x, x') + \hat{\mathbf{i}} \cdot (-\infty) && \text{if } x \in \{-\infty, +\infty\} \text{ and } x' \in F \text{ and } x' \geq 0 \\
&= \text{arc}_F(x, \text{neg}_F(x')) + \hat{\mathbf{i}} \cdot (+\infty) && \text{if } x \in \{-\infty, +\infty\} \text{ and } x' \in F \text{ and } x' < 0 \\
&= \text{no\_result}_{c(F)}(x + \hat{\mathbf{i}} \cdot x') && \text{otherwise}
\end{aligned}$$

NOTE – The inverse of cos is multi-valued, the real part may have any integer multiple of  $2 \cdot \pi$  added to it, and the result is also in the solution set. The arccos function (returning the principal value for the inverse) branch cuts at  $\{x \mid x \in \mathcal{R} \text{ and } |x| > 1\}$ . Thus  $\text{arccos}_{c(F)}(x + \hat{\mathbf{i}} \cdot \mathbf{0}) \neq \text{arccos}_{c(F)}(x + \hat{\mathbf{i}} \cdot (-\mathbf{0}))$  when  $|x| > 1$ .



### 5.3.2.10 Radian arc tangent

The  $\arctan_{i(F)}$  operation:

$$\arctan_{i(F)} : i(F) \rightarrow i(F) \cup \{\mathbf{underflow}, \mathbf{infinitary}\}$$

$$\arctan_{i(F)}(\hat{\mathbf{i}} \cdot x') = \hat{\mathbf{i}} \cdot \arctanh_F(x')$$

The  $\arctan_{c(F)}^*$  approximation helper function:

$$\arctan_{c(F)}^* : \mathcal{C} \rightarrow \mathcal{C}$$

$\arctan_{c(F)}^*(z)$  returns a close approximation to  $\arctan(z)$  in  $\mathcal{C}$  with maximum error  $max\_error\_tan_{c(F)}$ .

Further requirements on the  $\arctan_{c(F)}^*$  approximation helper function are:

$$\arctan_{c(F)}^*(\mathit{conj}(z)) = \mathit{conj}(\arctan_{c(F)}^*(z)) \quad \text{if } z \in \mathcal{C}$$

$$\arctan_{c(F)}^*(-z) = -\arctan_{c(F)}^*(z)$$

$$\Re(\arctan_{c(F)}^*(\hat{\mathbf{i}} \cdot x')) = \pi/2 \quad \begin{array}{l} \text{if } z \in \mathcal{C} \\ \text{if } x' \in \mathcal{R} \text{ and } |x'| > 1 \end{array}$$

The relationship to the  $\arctan_F^*$  and  $\arctanh_F^*$  approximation helper functions in an associated library for real-valued operations shall be:

$$\arctan_{c(F)}^*(x) = \arctan_F^*(x) \quad \text{if } x \in \mathcal{R}$$

$$\arctan_{c(F)}^*(\hat{\mathbf{i}} \cdot x') = \hat{\mathbf{i}} \cdot \arctanh_F^*(x') \quad \text{if } x' \in \mathcal{R} \text{ and } |x'| < 1$$

The  $\arctan_{c(F)}^\#$  range limitation helper function:

$$\arctan_{c(F)}^\#(z) = \max\{up_F(-\pi/2), \min\{\Re(\arctan_F^*(z)), down_F(\pi/2)\}\} + \hat{\mathbf{i}} \cdot \Im(\arctan_F^*(z))$$

The  $\arctan_{c(F)}$  operation:

$$\arctan_{c(F)} : c(F) \rightarrow c(F) \cup \{\mathbf{underflow}, \mathbf{infinitary}\}$$

$$\arctan_{c(F)}(x + \hat{\mathbf{i}} \cdot x')$$

$$= \mathit{result}_{c(F)}^*(\arctan_{c(F)}^\#(x + \hat{\mathbf{i}} \cdot x'), \mathit{nearest}_F)$$

$$\quad \text{if } x + \hat{\mathbf{i}} \cdot x' \in c(F) \text{ and } ((x' \neq 1 \text{ and } x' \neq -1) \text{ or } x \neq 0)$$

$$= \mathit{neg}_{c(F)}(\arctan_{c(F)}(0 + \hat{\mathbf{i}} \cdot \mathit{neg}_F(x')))$$

$$\quad \text{if } x = -\mathbf{0}$$

$$= \mathit{conj}_{c(F)}(\arctan_{c(F)}(x + \hat{\mathbf{i}} \cdot \mathbf{0}))$$

$$\quad \text{if } x' = -\mathbf{0} \text{ and } x \neq -\mathbf{0}$$

$$= \mathit{signb}_F(x) \cdot \mathit{down}_F(\pi/2)$$

$$\quad \text{if } (x \in \{-\infty, +\infty\} \text{ and } x' \in F \cup \{-\infty, +\infty\}) \text{ or } \\ (x \in F \cup \{-\infty, +\infty\} \text{ and } x' \in \{-\infty, +\infty\})$$

$$= \mathbf{infinitary}(\mathit{down}_F(\pi/2) + \hat{\mathbf{i}} \cdot (+\infty))$$

$$\quad \text{if } x = 0 \text{ and } x' = 1$$

$$= \mathbf{infinitary}(\mathit{down}_F(\pi/2) + \hat{\mathbf{i}} \cdot (-\infty))$$

$$\quad \text{if } x = 0 \text{ and } x' = -1$$

$$= \mathit{no\_result}_{c(F)}(x + \hat{\mathbf{i}} \cdot x')$$

$$\quad \text{otherwise}$$

NOTE – The inverse of tan is multi-valued, the real part may have any integer multiple of  $2 \cdot \pi$  (even any integer multiple of  $\pi$ ) added to it, and the result is also in the solution set. The arctan function (returning the principal value for the inverse) branch cuts at  $\{\hat{\mathbf{i}} \cdot x' \mid x' \in F \text{ and } |x'| > 1\}$ . Thus  $\arctan_{c(F)}(0 + \hat{\mathbf{i}} \cdot x') \neq \arctan_{c(F)}(-\mathbf{0} + \hat{\mathbf{i}} \cdot x')$  when  $|x'| > 1$ .

### 5.3.2.11 Radian arc cotangent

The  $\text{arccot}_{i(F)}$  operation:

$$\text{arccot}_{i(F)} : i(F) \rightarrow i(F) \cup \{\mathbf{underflow}, \mathbf{infinitary}\}$$

$$\text{arccot}_{i(F)}(\hat{\mathbf{i}} \cdot x') = \hat{\mathbf{i}} \cdot \text{arccoth}_F(x')$$

The  $\text{arccot}_{c(F)}^*$  approximation helper function:

$$\text{arccot}_{c(F)}^* : \mathcal{C} \rightarrow \mathcal{C}$$

$\text{arccot}_{c(F)}^*(z)$  returns a close approximation to  $\text{arccot}(z)$  in  $\mathcal{C}$  with maximum error  $\text{max\_error\_tan}_{c(F)}$ .

Further requirements on the  $\text{arccot}_{c(F)}^*$  approximation helper function are:

$$\text{arccot}_{c(F)}^*(\text{conj}(z)) = \text{conj}(\text{arccot}_{c(F)}^*(z)) \quad \text{if } z \in \mathcal{C}$$

$$\text{arccot}_{c(F)}^*(-z) = -\text{arccot}_{c(F)}^*(z) \quad \text{if } z \in \mathcal{C}$$

$$\Re(\text{arccot}_{c(F)}^*(\hat{\mathbf{i}} \cdot x')) = \pi/2 \quad \text{if } x' \in \mathcal{R} \text{ and } |x'| < 1$$

The relationship to the  $\text{arccot}_F^*$  and  $\text{arccoth}_F^*$  approximation helper functions in an associated library for real-valued operations shall be:

$$\text{arccot}_{c(F)}^*(x) = \text{arccot}_F^*(x) \quad \text{if } x \in \mathcal{R}$$

$$\text{arccot}_{c(F)}^*(\hat{\mathbf{i}} \cdot x') = \hat{\mathbf{i}} \cdot \text{arccoth}_F^*(-x') \quad \text{if } x' \in \mathcal{R}$$

The  $\text{arccot}_{c(F)}^\#$  range limitation helper function:

$$\text{arccot}_{c(F)}^\#(z) = \max\{\text{up}_F(-\pi/2), \min\{\Re(\text{arccot}_F^*(z)), \text{down}_F(\pi/2)\}\} + \hat{\mathbf{i}} \cdot \Im(\text{arccot}_F^*(z))$$

The  $\text{arccot}_{c(F)}$  operation:

$$\text{arccot}_{c(F)} : c(F) \rightarrow c(F) \cup \{\mathbf{underflow}, \mathbf{infinitary}\}$$

$$\text{arccot}_{c(F)}(x + \hat{\mathbf{i}} \cdot x')$$

$$= \text{result}_{c(F)}^*(\text{arccot}_{c(F)}^\#(x + \hat{\mathbf{i}} \cdot x'), \text{nearest}_F)$$

$$\text{if } x + \hat{\mathbf{i}} \cdot x' \in c(F)$$

$$= \text{neg}_{c(F)}(\text{arccot}_{c(F)}(0 + \hat{\mathbf{i}} \cdot \text{neg}_F(x')))$$

$$\text{if } x = -\mathbf{0}$$

$$= \text{conj}_{c(F)}(\text{arccot}_{c(F)}(x + \hat{\mathbf{i}} \cdot \mathbf{0}))$$

$$\text{if } x' = -\mathbf{0} \text{ and } x \neq -\mathbf{0}$$

$$= \text{mul}_F(\text{signb}_F(x), \mathbf{0}) + \hat{\mathbf{i}} \cdot \text{mul}_F(\text{signb}_F(x'), \mathbf{0})$$

$$\text{if } (x \in \{-\infty, +\infty\} \text{ and } x' \in F \cup \{-\infty, +\infty\}) \text{ or}$$

$$(x \in F \cup \{-\infty, +\infty\} \text{ and } x' \in \{-\infty, +\infty\})$$

$$= \mathbf{infinitary}(\text{down}_F(\pi/2), x) + \hat{\mathbf{i}} \cdot (+\infty)$$

$$\text{if } x' = -1 \text{ and } x = 0$$

$$= \mathbf{infinitary}(\text{down}_F(\pi/2), x) + \hat{\mathbf{i}} \cdot (-\infty)$$

$$\text{if } x' = 1 \text{ and } x = 0$$

$$= \text{no\_result}_{c(F)}(x + \hat{\mathbf{i}} \cdot x')$$

otherwise

NOTE – The inverse of cot is multi-valued, the real part may have any integer multiple of  $2 \cdot \pi$  (even any integer multiple of  $\pi$ ) added to it, and the result is also in the solution set. The arccot function (returning the principal value for the inverse) branch cuts at  $\{\hat{\mathbf{i}} \cdot x' \mid x' \in \mathcal{R} \text{ and } |x'| < 1\}$ . Thus  $\text{arccot}_{c(F)}(0 + \hat{\mathbf{i}} \cdot x') \neq \text{arccot}_{c(F)}(-\mathbf{0} + \hat{\mathbf{i}} \cdot x')$  when  $|x'| < 1$  or  $x' = -\mathbf{0}$ .

### 5.3.2.12 Radian arc secant

The  $\text{arcsec}_{c(F)}^*$  approximation helper function:

$$\text{arcsec}_{c(F)}^* : \mathcal{C} \rightarrow \mathcal{C}$$

$\text{arcsec}_{c(F)}^*(z)$  returns a close approximation to  $\text{arcsec}(z)$  in  $\mathcal{C}$  with maximum error  $\text{max\_error\_tan}_{c(F)}$ .

Further requirements on the  $\text{arcsec}_{c(F)}^*$  approximation helper function are:

$$\text{arcsec}_{c(F)}^*(\text{conj}(z)) = \text{conj}(\text{arcsec}_{c(F)}^*(z)) \quad \text{if } z \in \mathcal{C}$$

The relationship to the  $\text{arcsec}_F^*$  and  $\text{arcsech}_F^*$  approximation helper functions in an associated library for real-valued operations shall be:

$$\text{arcsec}_{c(F)}^*(x) = \text{arcsec}_F^*(x) \quad \text{if } x \in \mathcal{R}$$

The  $\text{arcsec}_{c(F)}^\#$  range limitation helper function:

$$\begin{aligned} \text{arcsec}_{c(F)}^\#(z) &= \min\{\Re(\text{arcsec}_F^*(z)), \text{down}_F(\pi/2)\} + \hat{\mathbf{i}} \cdot \Im(\text{arcsec}_F^*(z)) && \text{if } \Re(z) \geq 1 \\ &= \min\{\Re(\text{arcsec}_F^*(z)), \text{down}_F(\pi)\} + \hat{\mathbf{i}} \cdot \Im(\text{arcsec}_F^*(z)) && \text{if } \Re(z) \leq -1 \\ &= \text{arcsec}_{c(F)}(z) && \text{otherwise} \end{aligned}$$

The  $\text{arcsec}_{c(F)}$  operation:

$$\begin{aligned} \text{arcsec}_{c(F)} : c(F) &\rightarrow c(F) \cup \{\mathbf{underflow}, \mathbf{infinitary}\} \\ \text{arcsec}_{c(F)}(x + \hat{\mathbf{i}} \cdot x') &= \text{result}_{c(F)}^*(\text{arcsec}_{c(F)}^\#(x + \hat{\mathbf{i}} \cdot x'), \text{nearest}_F) && \text{if } x + \hat{\mathbf{i}} \cdot x' \in c(F) \\ &= \text{arcsec}_{c(F)}(0 + \hat{\mathbf{i}} \cdot \text{neg}_F(x')) && \text{if } x = -\mathbf{0} \\ &= \text{conj}_{c(F)}(\text{arcsec}_{c(F)}(x + \hat{\mathbf{i}} \cdot 0)) && \text{if } x' = -\mathbf{0} \text{ and } x \neq -\mathbf{0} \\ &= \text{arccos}_{c(F)}(\text{div}_{c(F)}(1, x + \hat{\mathbf{i}} \cdot x')) && \text{if } x \in \{-\infty, +\infty\} \text{ and } x' \in F \cup \{-\infty, +\infty\} \\ &= \text{arccos}_{c(F)}(\text{div}_{c(F)}(1, x + \hat{\mathbf{i}} \cdot x')) && \text{if } x' \in \{-\infty, +\infty\} \text{ and } x \in F \cup \{-\infty, +\infty\} \\ &= \mathbf{infinitary}(? + \hat{\mathbf{i}} \cdot (+\infty)) && \text{if } x' = 0 \text{ and } x = 0 \\ &= \text{no\_result}_{c(F)}(x + \hat{\mathbf{i}} \cdot x') && \text{otherwise} \end{aligned}$$

NOTE – The inverse of sec is multi-valued, the real part may have any integer multiple of  $2 \cdot \pi$  added to it, and the result is also in the solution set. The arcsec function (returning the principal value for the inverse) branch cuts at  $\{x \mid x \in \mathcal{R} \text{ and } -1 < x < 1\}$ . Thus  $\text{arcsec}_{c(F)}(x + \hat{\mathbf{i}} \cdot 0) \neq \text{arcsec}_{c(F)}(x + \hat{\mathbf{i}} \cdot (-\mathbf{0}))$  when  $-1 < x < 1$  or  $x = -\mathbf{0}$ .

### 5.3.2.13 Radian arc cosecant

The  $\text{arccsc}_{i(F)}$  operation:

$$\begin{aligned} \text{arccsc}_{i(F)} : i(F) &\rightarrow i(F) \cup \{\mathbf{underflow}, \mathbf{infinitary}\} \\ \text{arccsc}_{i(F)}(\hat{\mathbf{i}} \cdot x') &= \hat{\mathbf{i}} \cdot \text{arccsch}_F(x') \end{aligned}$$

The  $\text{arccsc}_{c(F)}^*$  approximation helper function:

$$\text{arccsc}_{c(F)}^* : \mathcal{C} \rightarrow \mathcal{C}$$

$\text{arccsc}_{c(F)}^*(z)$  returns a close approximation to  $\text{arccsc}(z)$  in  $\mathcal{C}$  with maximum error  $\text{max\_error\_tan}_{c(F)}$ .

Further requirements on the  $\text{arccsc}_{c(F)}^*$  approximation helper function are:

$$\begin{aligned} \operatorname{arccsc}_{c(F)}^*(\operatorname{conj}(z)) &= \operatorname{conj}(\operatorname{arccsc}_{c(F)}^*(z)) && \text{if } z \in \mathcal{C} \\ \operatorname{arccsc}_{c(F)}^*(-z) &= -\operatorname{arccsc}_{c(F)}^*(z) && \text{if } z \in \mathcal{C} \end{aligned}$$

The relationship to the  $\operatorname{arccsc}_F^*$  and  $\operatorname{arccsch}_F^*$  approximation helper functions in an associated library for real-valued operations shall be:

$$\begin{aligned} \operatorname{arccsc}_{c(F)}^*(x) &= \operatorname{arccsc}_F^*(x) && \text{if } x \in \mathcal{R} \\ \operatorname{arccsc}_{c(F)}^*(\hat{\mathbf{i}} \cdot x') &= \hat{\mathbf{i}} \cdot \operatorname{arccsch}_F^*(-x') && \text{if } x' \in \mathcal{R} \end{aligned}$$

The  $\operatorname{arccsc}_{c(F)}^\#$  range limitation helper function:

$$\operatorname{arccsc}_{c(F)}^\#(z) = \max\{\operatorname{up}_F(-\pi/2), \min\{\Re(\operatorname{arccsc}_F^*(z)), \operatorname{down}_F(\pi/2)\}\} + \hat{\mathbf{i}} \cdot \Im(\operatorname{arccsc}_F^*(z))$$

The  $\operatorname{arccsc}_{c(F)}$  operation:

$$\begin{aligned} \operatorname{arccsc}_{c(F)} : c(F) &\rightarrow c(F) \cup \{\mathbf{underflow}, \mathbf{infinitary}\} \\ \operatorname{arccsc}_{c(F)}(x + \hat{\mathbf{i}} \cdot x') &= \operatorname{result}_{c(F)}^*(\operatorname{arccsc}_{c(F)}^\#(x + \hat{\mathbf{i}} \cdot x'), \operatorname{nearest}_F) && \text{if } x + \hat{\mathbf{i}} \cdot x' \in c(F) \\ &= \operatorname{neg}_{c(F)}(\operatorname{arccsc}_{c(F)}(0 + \hat{\mathbf{i}} \cdot \operatorname{neg}_F(x'))) && \text{if } x = -\mathbf{0} \\ &= \operatorname{conj}_{c(F)}(\operatorname{arccsc}_{c(F)}(x + \hat{\mathbf{i}} \cdot 0)) && \text{if } x' = -\mathbf{0} \text{ and } x \neq -\mathbf{0} \\ &= \operatorname{arcsin}_{c(F)}(\operatorname{div}_{c(F)}(1, x + \hat{\mathbf{i}} \cdot x')) && \text{if } x \in \{-\infty, +\infty\} \text{ and } x' \in F \cup \{-\infty, +\infty\} \\ &= \operatorname{arcsin}_{c(F)}(\operatorname{div}_{c(F)}(1, x + \hat{\mathbf{i}} \cdot x')) && \text{if } x' \in \{-\infty, +\infty\} \text{ and } x \in F \cup \{-\infty, +\infty\} \\ &= \mathbf{infinitary}(? + \hat{\mathbf{i}} \cdot (-\infty)) && \text{if } x = 0 \text{ and } x' = 0 \\ &= \operatorname{no\_result}_{c(F)}(x + \hat{\mathbf{i}} \cdot x') && \text{otherwise} \end{aligned}$$

NOTE – The inverse of csc is multi-valued, the real part may have any integer multiple of  $2 \cdot \pi$  added to it, and the result is also in the solution set. The arccsc function (returning the principal value for the inverse) branch cuts at  $\{x \mid x \in \mathcal{R} \text{ and } -1 < x < 1\}$ . Thus  $\operatorname{arccsc}_{c(F)}(x + \hat{\mathbf{i}} \cdot 0) \neq \operatorname{arccsc}_{c(F)}(x + \hat{\mathbf{i}} \cdot (-\mathbf{0}))$  when  $-1 < x < 1$  or  $x = -\mathbf{0}$ .

### 5.3.3 Operations for hyperbolic elementary functions

The complex trigonometric operations which are used in the specifications below are specified in clause 5.3.2. Note that the correspondences specified here are exact, not approximate.

#### 5.3.3.1 Hyperbolic normalisation

$$\operatorname{radh}_F : F \rightarrow F$$

$$\operatorname{radh}_F(x) = x$$

$$\operatorname{radh}_{i(F)} : i(F) \rightarrow i(F) \cup \{\mathbf{underflow}, \mathbf{absolute\_precision\_underflow}\}$$

$$\operatorname{radh}_{i(F)}(\hat{\mathbf{i}} \cdot x') = \hat{\mathbf{i}} \cdot \operatorname{rad}_F(x')$$

$$\operatorname{radh}_{c(F)} : c(F) \rightarrow c(F) \cup \{\mathbf{underflow}, \mathbf{absolute\_precision\_underflow}\}$$

$$\begin{aligned} \operatorname{radh}_{c(F)}(x + \hat{\mathbf{i}} \cdot x') &= \hat{\mathbf{i}} \cdot \operatorname{rad}_{c(F)}(x' + \hat{\mathbf{i}} \cdot \operatorname{neg}_F(x)) \end{aligned}$$

NOTE –  $rad_{c(F)}$  is defined in clause 5.3.2.1.

### 5.3.3.2 Hyperbolic sine

$$\sinh_{i(F)} : i(F) \rightarrow i(F) \cup \{\text{underflow, absolute\_precision\_underflow}\}$$

$$\sinh_{i(F)}(\hat{\mathbf{i}} \cdot x') = \hat{\mathbf{i}} \cdot (\sin_F(x'))$$

$$\sinh_{c(F)} : c(F) \rightarrow c(F) \cup \{\text{underflow, overflow, absolute\_precision\_underflow}\}$$

$$\begin{aligned} \sinh_{c(F)}(x + \hat{\mathbf{i}} \cdot x') \\ = \hat{\mathbf{i}} \cdot (\sin_{c(F)}(x' + \hat{\mathbf{i}} \cdot \text{neg}_F(x))) \end{aligned}$$

### 5.3.3.3 Hyperbolic cosine

$$\cosh_{i(F)} : i(F) \rightarrow F \cup \{\text{underflow, absolute\_precision\_underflow}\}$$

$$\cosh_{i(F)}(\hat{\mathbf{i}} \cdot x') = \cos_F(x')$$

$$\cosh_{c(F)} : c(F) \rightarrow c(F) \cup \{\text{underflow, overflow, absolute\_precision\_underflow}\}$$

$$\begin{aligned} \cosh_{c(F)}(x + \hat{\mathbf{i}} \cdot x') \\ = \cos_{c(F)}(x' + \hat{\mathbf{i}} \cdot \text{neg}_F(x)) \end{aligned}$$

### 5.3.3.4 Hyperbolic tangent

$$\tanh_{i(F)} : i(F) \rightarrow i(F) \cup \{\text{underflow, overflow, absolute\_precision\_underflow}\}$$

$$\tanh_{i(F)}(\hat{\mathbf{i}} \cdot x') = \hat{\mathbf{i}} \cdot (\tan_F(x'))$$

$$\tanh_{c(F)} : c(F) \rightarrow c(F) \cup \{\text{underflow, overflow, absolute\_precision\_underflow}\}$$

$$\begin{aligned} \tanh_{c(F)}(x + \hat{\mathbf{i}} \cdot x') \\ = \hat{\mathbf{i}} \cdot (\tan_{c(F)}(x' + \hat{\mathbf{i}} \cdot \text{neg}_F(x))) \end{aligned}$$

### 5.3.3.5 Hyperbolic cotangent

$$\coth_{i(F)} : i(F) \rightarrow i(F) \cup \{\text{underflow, overflow, infinitary, absolute\_precision\_underflow}\}$$

$$\coth_{i(F)}(\hat{\mathbf{i}} \cdot x') = \hat{\mathbf{i}} \cdot (\cot_F(\text{neg}_F(x')))$$

$$\coth_{c(F)} : c(F) \rightarrow c(F) \cup \{\text{underflow, overflow, infinitary, absolute\_precision\_underflow}\}$$

$$\begin{aligned} \coth_{c(F)}(x + \hat{\mathbf{i}} \cdot x') \\ = \hat{\mathbf{i}} \cdot (\cot_{c(F)}(\text{neg}_F(x') + \hat{\mathbf{i}} \cdot x)) \end{aligned}$$

### 5.3.3.6 Hyperbolic secant

$$\begin{aligned} sech_{i(F)} : i(F) &\rightarrow F \cup \{\text{underflow}, \text{overflow}, \text{absolute\_precision\_underflow}\} \\ sech_{i(F)}(\hat{\mathbf{i}} \cdot x') &= sec_F(neg_F(x')) \end{aligned}$$

$$\begin{aligned} sech_{c(F)} : c(F) &\rightarrow c(F) \cup \{\text{underflow}, \text{overflow}, \text{absolute\_precision\_underflow}\} \\ sech_{c(F)}(x + \hat{\mathbf{i}} \cdot x') &= sec_{c(F)}(neg_F(x') + \hat{\mathbf{i}} \cdot x) \end{aligned}$$

### 5.3.3.7 Hyperbolic cosecant

$$\begin{aligned} csch_{i(F)} : i(F) &\rightarrow i(F) \cup \{\text{underflow}, \text{overflow}, \text{infinitary}, \text{absolute\_precision\_underflow}\} \\ csch_{i(F)}(\hat{\mathbf{i}} \cdot x') &= \hat{\mathbf{i}} \cdot (csc_F(neg_F(x'))) \end{aligned}$$

$$\begin{aligned} csch_{c(F)} : c(F) &\rightarrow c(F) \cup \{\text{underflow}, \text{overflow}, \text{infinitary}, \text{absolute\_precision\_underflow}\} \\ csch_{c(F)}(x + \hat{\mathbf{i}} \cdot x') &= \hat{\mathbf{i}} \cdot (csc_{c(F)}(neg_F(x') + \hat{\mathbf{i}} \cdot x)) \end{aligned}$$

### 5.3.3.8 Inverse hyperbolic sine

$$\begin{aligned} arcsinh_{i(F)} : i(F) &\rightarrow i(F) \cup \{\text{underflow}\} \\ arcsinh_{i(F)}(\hat{\mathbf{i}} \cdot x') &= \hat{\mathbf{i}} \cdot (arcsin_F(x')) \end{aligned}$$

$$\begin{aligned} arcsinh_{c(F)} : c(F) &\rightarrow c(F) \cup \{\text{underflow}\} \\ arcsinh_{c(F)}(x + \hat{\mathbf{i}} \cdot x') &= \hat{\mathbf{i}} \cdot (arcsin_{c(F)}(x' + \hat{\mathbf{i}} \cdot neg_F(x))) \end{aligned}$$

NOTE – The inverse of sinh is multi-valued, the imaginary part may have any integer multiple of  $2 \cdot \pi$  added to it, and the result is also in the solution set. The arcsinh function (returning the principal value for the inverse) branch cuts at  $\{\hat{\mathbf{i}} \cdot x' \mid x' \in \mathcal{R} \text{ and } |x'| > 1\}$ . Thus  $arcsinh_{c(F)}(0 + \hat{\mathbf{i}} \cdot x') \neq arcsinh_{c(F)}(-0 + \hat{\mathbf{i}} \cdot x')$  when  $|x'| > 1$ .

### 5.3.3.9 Inverse hyperbolic cosine

$$\begin{aligned} arccosh_{c(F)} : c(F) &\rightarrow c(F) \\ arccosh_{c(F)}(x + \hat{\mathbf{i}} \cdot x') &= \hat{\mathbf{i}} \cdot (arccos_{c(F)}(x + \hat{\mathbf{i}} \cdot x')) \\ &\quad \text{if } (x' \in F \text{ and } x' \geq 0) \text{ or } x' = +\infty \\ &= neg_{c(F)}(\hat{\mathbf{i}} \cdot (arccos_{c(F)}(x + \hat{\mathbf{i}} \cdot x'))) \\ &\quad \text{if } (x' \in F \text{ and } x' < 0) \text{ or } x' \in \{-\infty, -0\} \end{aligned}$$

NOTE – The inverse of cosh is multi-valued, the imaginary part may have any integer multiple of  $2 \cdot \pi$  added to it, and the result is also in the solution set. The arccosh function (returning the principal value for the inverse) branch cuts at  $\{x \mid x \in \mathcal{R} \text{ and } x < 1\}$ . Thus  $arccosh_{c(F)}(x + \hat{\mathbf{i}} \cdot 0) \neq arccosh_{c(F)}(x + \hat{\mathbf{i}} \cdot (-0))$  when  $x < 1$  or  $x = -0$ .

### 5.3.3.10 Inverse hyperbolic tangent

$$\begin{aligned} \operatorname{arctanh}_{i(F)} : i(F) &\rightarrow i(F) \cup \{\mathbf{underflow}, \mathbf{infinitary}\} \\ \operatorname{arctanh}_{i(F)}(\hat{\mathbf{i}} \cdot x') &= \hat{\mathbf{i}} \cdot (\operatorname{arctan}_F(x')) \end{aligned}$$

$$\begin{aligned} \operatorname{arctanh}_{c(F)} : c(F) &\rightarrow c(F) \cup \{\mathbf{underflow}, \mathbf{infinitary}\} \\ \operatorname{arctanh}_{c(F)}(x + \hat{\mathbf{i}} \cdot x') &= \hat{\mathbf{i}} \cdot (\operatorname{arctan}_{c(F)}(x' + \hat{\mathbf{i}} \cdot \operatorname{neg}_F(x))) \end{aligned}$$

NOTE – The inverse of tanh is multi-valued, the imaginary part may have any integer multiple of  $2 \cdot \pi$  (even any integer multiple of  $\pi$ ) added to it, and the result is also in the solution set. The arctanh function (returning the principal value for the inverse) branch cuts at  $\{x \mid x \in \mathcal{R} \text{ and } |x| > 1\}$ . Thus  $\operatorname{arctanh}_{c(F)}(x + \hat{\mathbf{i}} \cdot 0) \neq \operatorname{arctanh}_{c(F)}(x + \hat{\mathbf{i}} \cdot (-0))$  when  $|x| > 1$ .

### 5.3.3.11 Inverse hyperbolic cotangent

$$\begin{aligned} \operatorname{arccoth}_{i(F)} : i(F) &\rightarrow i(F) \cup \{\mathbf{underflow}, \mathbf{infinitary}\} \\ \operatorname{arccoth}_{i(F)}(\hat{\mathbf{i}} \cdot x') &= \hat{\mathbf{i}} \cdot (\operatorname{arccot}_F(\operatorname{neg}_F(x'))) \end{aligned}$$

$$\begin{aligned} \operatorname{arccoth}_{c(F)} : c(F) &\rightarrow c(F) \cup \{\mathbf{underflow}, \mathbf{infinitary}\} \\ \operatorname{arccoth}_{c(F)}(x + \hat{\mathbf{i}} \cdot x') &= \hat{\mathbf{i}} \cdot (\operatorname{arccot}_{c(F)}(\operatorname{neg}_F(x') + \hat{\mathbf{i}} \cdot x)) \end{aligned}$$

NOTE – The inverse of coth is multi-valued, the imaginary part may have any integer multiple of  $2 \cdot \pi$  (even any integer multiple of  $\pi$ ) added to it, and the result is also in the solution set. The arccoth function (returning the principal value for the inverse) branch cuts at  $\{x \mid x \in \mathcal{R} \text{ and } |x| < 1\}$ . Thus  $\operatorname{arccoth}_{c(F)}(x + \hat{\mathbf{i}} \cdot 0) \neq \operatorname{arccoth}_{c(F)}(x + \hat{\mathbf{i}} \cdot (-0))$  when  $|x| < 1$  or  $x = -0$ .

### 5.3.3.12 Inverse hyperbolic secant

$$\begin{aligned} \operatorname{arcsech}_{c(F)} : c(F) &\rightarrow c(F) \cup \{\mathbf{underflow}, \mathbf{infinitary}\} \\ \operatorname{arcsech}_{c(F)}(x + \hat{\mathbf{i}} \cdot x') &= \hat{\mathbf{i}} \cdot (\operatorname{arcsec}_{c(F)}(x + \hat{\mathbf{i}} \cdot x')) \\ &\quad \text{if } (x' \in F \text{ and } x' \geq 0) \text{ or } x' = +\infty \\ &= \operatorname{neg}_{c(F)}(\hat{\mathbf{i}} \cdot (\operatorname{arcsec}_{c(F)}(x + \hat{\mathbf{i}} \cdot x'))) \\ &\quad \text{if } (x' \in F \text{ and } x' < 0) \text{ or } x' \in \{-\infty, -0\} \end{aligned}$$

NOTE – The inverse of sech is multi-valued, the imaginary part may have any integer multiple of  $2 \cdot \pi$  added to it, and the result is also in the solution set. The arcsech function (returning the principal value for the inverse) branch cuts at  $\{x \mid x \in \mathcal{R} \text{ and } x \leq 0 \text{ or } x > 1\}$ . Thus  $\operatorname{arcsech}_{c(F)}(x + \hat{\mathbf{i}} \cdot 0) \neq \operatorname{arcsech}_{c(F)}(x + \hat{\mathbf{i}} \cdot (-0))$  when  $x \leq 0$  or  $x = -0$  or  $x > 1$ .

### 5.3.3.13 Inverse hyperbolic cosecant

$$\begin{aligned} \operatorname{arccsch}_{i(F)} : i(F) &\rightarrow i(F) \cup \{\mathbf{underflow}, \mathbf{infinitary}\} \\ \operatorname{arccsch}_{i(F)}(\hat{\mathbf{i}} \cdot x') &= \hat{\mathbf{i}} \cdot (\operatorname{arccsc}_F(\operatorname{neg}_F(x'))) \end{aligned}$$

$$\begin{aligned} \operatorname{arccsch}_{c(F)} : c(F) &\rightarrow c(F) \cup \{\mathbf{underflow}, \mathbf{infinitary}\} \\ \operatorname{arccsch}_{c(F)}(x + \hat{\mathbf{i}} \cdot x') &= \hat{\mathbf{i}} \cdot (\operatorname{arccsc}_{c(F)}(\operatorname{neg}_F(x') + \hat{\mathbf{i}} \cdot x)) \end{aligned}$$

NOTE – The inverse of  $\operatorname{csch}$  is multi-valued, the imaginary part may have any integer multiple of  $2 \cdot \pi$  added to it, and the result is also in the solution set. The  $\operatorname{arccsch}$  function (returning the principal value for the inverse) branch cuts at  $\{\hat{\mathbf{i}} \cdot x' \mid x' \in \mathcal{R} \text{ and } |x'| < 1\}$ . Thus  $\operatorname{arccsch}_{c(F)}(0 + \hat{\mathbf{i}} \cdot x') \neq \operatorname{arccsch}_{c(F)}(-0 + \hat{\mathbf{i}} \cdot x')$  when  $-1 < x' < 1$  or  $x' = -0$ .

## 5.4 Operations for conversion between numeric datatypes

### 5.4.1 Integer to complex integer conversions

Let  $I$  be the non-special value set for an integer datatype.

$$\begin{aligned} \operatorname{convert}_{I \rightarrow c(I)} : I &\rightarrow c(I) \\ \operatorname{convert}_{I \rightarrow c(I)}(x) &= x + \hat{\mathbf{i}} \cdot 0 \end{aligned}$$

$$\begin{aligned} \operatorname{convert}_{i(I) \rightarrow c(I)} : i(I) &\rightarrow c(I) \\ \operatorname{convert}_{i(I) \rightarrow c(I)}(\hat{\mathbf{i}} \cdot x) &= 0 + \hat{\mathbf{i}} \cdot x \end{aligned}$$

### 5.4.2 Floating point to complex floating point conversions

Let  $F$  be the non-special value set for a floating point datatype.

The  $\operatorname{convert}_{F \rightarrow c(F)}$  operation:

$$\begin{aligned} \operatorname{convert}_{F \rightarrow c(F)} : F &\rightarrow c(F) \\ \operatorname{convert}_{F \rightarrow c(F)}(x) &= x + \hat{\mathbf{i}} \cdot 0 \end{aligned}$$

The  $\operatorname{convert}_{i(F) \rightarrow c(F)}$  operation:

$$\begin{aligned} \operatorname{convert}_{i(F) \rightarrow c(F)} : i(F) &\rightarrow c(F) \\ \operatorname{convert}_{i(F) \rightarrow c(F)}(\hat{\mathbf{i}} \cdot x) &= 0 + \hat{\mathbf{i}} \cdot x \end{aligned}$$

## 6 Notification

Notification is the process by which a user or program is informed that an arithmetic operation cannot return a suitable numeric result. Specifically, a notification shall occur when any arithmetic operation returns an exceptional value. Notification shall be performed according to the requirements of clause 6 of part 1.

An implementation shall not give notifications for operations conforming to this part, unless the specification requires that an exceptional value results for the given arguments.

The default method of notification should be recording of indicators.



## 6.1 Continuation values

If notifications are handled by a recording of indicators, in the event of notification the implementation shall provide a *continuation value* to be used in subsequent arithmetic operations. Continuation values may be in  $i(I)$ ,  $c(I)$ ,  $i(F)$  or  $c(F)$  (as appropriate), or be special values (where the real or imaginary component is  $-0$ ,  $-\infty$ ,  $+\infty$ , or a **qNaN**).

Floating point datatypes that satisfy the requirements of IEC 60559 have special values in addition to the values in  $F$ . These are:  $-0$ ,  $+\infty$ ,  $-\infty$ , *signaling NaNs* (**sNaN**), and *quiet NaNs* (**qNaN**). Such values may be components of complex floating point datatypes, and may be included in values passed as arguments to operations, and used as results or continuation values. Floating point types that do not fully conform to IEC 60559 can also have values corresponding to  $-0$ ,  $+\infty$ ,  $-\infty$ , or **NaN**.

## 7 Relationship with language standards

A computing system often provides some of the operations specified in this part within the context of a programming language. The requirements of the present standard shall be in addition to those imposed by the relevant programming language standards.

This part does not define the syntax of arithmetic expressions. However, programmers need to know how to reliably access the operations specified in this Part.

NOTE 1 – Providing the information required in this clause is properly the responsibility of programming language standards. An individual implementation would only need to provide details if it could not cite an appropriate clause of the language or binding standard.

An implementation shall document the notation that should be used to invoke an operation specified in this Part and made available. An implementation should document the notation that should be used to invoke an operation specified in this Part and that could be made available.

NOTE 2 – For example, the complex radian arc sine operation for an argument  $x$  ( $\text{arcsin}_{c(F)}(x)$ ) might be invoked as

<code>arcsin(x)</code>	in Ada [7]
<code>casin(x)</code>	in C [13]
<code>asin(x)</code>	in Fortran [18] and C++ [14]
<code>(asin x)</code>	in Common Lisp [38]

with a suitable expression of the argument ( $x$ ).

An implementation shall document the semantics of arithmetic expressions in terms of compositions of the operations specified in clause 5 of this Part and in clause 5 of Part 1.

Compilers often “optimize” code as part of compilation. Thus, an arithmetic expression might not be executed as written. An implementation shall document the possible transformations of arithmetic expressions (or groups of expressions) that it permits. Typical transformations include

- a) Insertion of operations, such as datatype conversions or changes in precision.
- b) Replacing operations (or entire subexpressions) with others, such as “`cos(-x)`”  $\rightarrow$  “`cos(x)`” (exactly the same result) or “`pi - arccos(x)`”  $\rightarrow$  “`arccos(-x)`” (more accurate result).
- c) Evaluating constant subexpressions.
- d) Eliminating unneeded subexpressions.

Only transformations which alter the semantics of an expression (the values produced, and the notifications generated) need be documented. Only the range of permitted transformations need be documented. It is not necessary to describe the specific choice of transformations that will be applied to a particular expression.

The textual scope of such transformations shall be documented, and any mechanisms that provide programmer control over this process should be documented as well.

## 8 Documentation requirements

In order to conform to this Part, an implementation shall include documentation providing the following information to programmers.

NOTE 1 – Much of the documentation required in this clause is properly the responsibility of programming language or binding standards. An individual implementation would only need to provide details if it could not cite an appropriate clause of the language or binding standard.

- a) A list of the provided operations that conform to this part.
- b) For each maximum error parameter, the value of that parameter or definition of that parameter function. Only maximum error parameters that are relevant to the provided operations need be given.
- c) The value of the parameters *big\_angle\_r<sub>F</sub>* and *big\_angle\_u<sub>F</sub>*. Only big angle parameters that are relevant to the provided operations need be given.
- d) For the *nearest<sub>F</sub>* function, the rule used for rounding halfway cases, unless *iec\_559<sub>F</sub>* is fixed to true.
- e) For each conforming operation, the continuation value provided for each notification condition. Specific continuation values that are required by this Part need not be documented. If the notification mechanism does not make use of continuation values (see clause 6), continuation values need not be documented.

NOTE 2 – Implementations that do not provide infinities or NaNs will have to document any continuation values used in place of such values.

- f) For each conforming operation, how the results depend on the rounding mode, if rounding modes are provided. Operations may be insensitive to the rounding mode, or sensitive to it, but even then need not heed the rounding mode.
- g) For each conforming operation, the notation to be used for invoking that operation.
- h) For each maximum error parameter, the notation to be used to access that parameter.
- i) The notation to be used to access the parameters *big\_angle\_r<sub>F</sub>* and *big\_angle\_u<sub>F</sub>*.
- j) For each of the provided operations where this Part specifies a relation to another operation specified in this Part, the binding for that other operation.
- k) For numerals conforming to this Part, which available string conversion operations, including reading from input, give exactly the same conversion results, even if the string syntaxes for ‘internal’ and ‘external’ numerals are different.

Since the integer and floating point datatypes used in conforming operations shall satisfy the requirements of Part 1, the following information shall also be provided by any conforming implementation.

- l) The means for selecting the modes of operation that ensure conformity.
- m) The translation of arithmetic expressions into combinations of the operations provided by any part of ISO/IEC 10967, including any use made of higher precision. (See clause 7 of Part 1.)

- n) The methods used for notification, and the information made available about the notification. (See clause 6 of Part 1.)
- o) The means for selecting among the notification methods, and the notification method used in the absence of a user selection. (See clause 6.3 of Part 1.)
- p) When “recording of indicators” is the method of notification, the datatype used to represent *Ind* (see clause 6.1.2 of Part 1), the method for denoting the values of *Ind*, and the notation for invoking each of the “indicator” operations. *E* is the set of notification indicators. The association of values in *Ind* with subsets of *E* must be clear. In interpreting clause 6.1.2 of Part 1, the set of indicators *E* shall be interpreted as including all exceptional values listed in the signatures of conforming operations. In particular, *E* may need to contain **infinitary** and **absolute\_precision\_underflow**.



## Annex A (normative)

### Partial conformity

If an implementation of an operation fulfills all relevant requirements according to the main normative text in this Part, except the ones relaxed in this Annex, the implementation of that operation is said to *partially conform* to this Part.

Conformity to this Part shall not be claimed for operations that only fulfill Partial conformity.

Partial conformity shall not be claimed for operations that relax other requirements than those relaxed in this Annex.

#### A.1 Maximum error relaxation

This part has the following maximum error requirements for conformity.

$$\begin{aligned} \mathit{max\_error\_mul}_{c(F)} &\in [0.5, ??] \\ \mathit{max\_error\_div}_{c(F)} &\in [0.5, ??] \\ \mathit{max\_error\_exp}_{c(F)} &\in [0.5, ??] \\ \mathit{max\_error\_power}_{c(F)} &\in [0.5, ??] \\ \mathit{max\_error\_sin}_{c(F)} &\in [0.5, ??] \\ \mathit{max\_error\_tan}_{c(F)} &\in [0.5, ??] \end{aligned}$$

In a partially conforming implementation the maximum error parameters may be greater than what is specified by this part. The maximum error parameter values given by an implementation shall still adequately reflect the accuracy of the relevant operations, if a claim of partial conformity is made.

A partially conforming implementation shall document which maximum error parameters have greater values than specified by this part, and their values.

#### A.2 Extra accuracy requirements relaxation

This Part has a number of extra accuracy requirements. These are detailed in the paragraphs beginning “Further requirements on the  $op_F^*$  approximation helper function are:”.

In a partially conforming implementation these further requirements need not be fulfilled. The values returned must still be within the maximum error bounds that are given by the maximum error parameters, if a claim of partial conformity is made.

A partially conforming implementation shall document which extra accuracy requirements are not fulfilled by the implementation.

#### A.3 Partial conformity to part 1 or to part 2

...



## Annex B (informative)

### Rationale

This annex explains and clarifies some of the ideas behind *Information technology – Language independent arithmetic – Part 2: Elementary numerical functions* (LIA-2).

#### B.1 Scope

##### B.1.1 Inclusions

LIA-2 is intended to define the meaning of some operations on integer and floating point types as specified in LIA-1 (ISO/IEC 10967-1), in addition to the operations specified in LIA-1. LIA-2 does not specify operations for any additional arithmetic datatypes, though fixed point datatypes are used in some of the specifications for conversion operations.

The specifications for the operations covered by LIA-2 are given in sufficient detail to

- a) support detailed and accurate numerical analysis of arithmetic algorithms,
- b) enable a precise determination of conformity or non-conformity, and
- c) prevent exceptions (like overflow) from going undetected.

LIA-2 does in no way prevent language standards or implementations including further arithmetic operations, other variations of included arithmetic operations, or the inclusion of further arithmetic datatypes, like rational number or fixed point datatypes. Some of these may become the topic of standardisation in other parts of LIA.

##### B.1.2 Exclusions

LIA-2 is not concerned with techniques for the *implementation* of numerical functions. Even when an LIA-2 specification is made in terms of other LIA-1 or LIA-2 operations, that does *not* imply a requirement that an implementation implements the operation in terms of those other operations. It is sufficient that the result (returned value or returned continuation value, and exception behaviour) is *as if* it was implemented in terms of those other operations.

LIA-2 does not provide specifications for operations which involve no arithmetic processing, like assignment and parameter passing, though any implicit conversions done in association with such operations are in scope. The implicit conversions should be made available to the programmer as explicit conversions.

LIA-2 does not cover operations for the support of domains such as linear algebra, statistics, and symbolic processing. Such domains deserve *separate* standardisation, if standardised.

LIA-2 only covers operations that involve integer or floating point datatypes, as specified in LIA-1, and in some cases also a Boolean datatype, but then only as result. The operations covered by LIA-2 are often to some extent covered by programming language standards, like the operations `sin`, `cos`, `tan`, `arctan`, and so on.

## B.2 Conformity

Conformity to this standard is dependent on the existence of language binding standards. Each programming language committee (or other organisation responsible for a programming language or other specification to which LIA-1 and LIA-2 may apply) is encouraged to produce a binding covering at least those operations already required by the programming language (or similar) and also specified in LIA-2.

The term “programming language” is here used in a generalised sense to include other computing entities such as calculators, spread sheets, page description languages, web-script languages, and database query languages to the extent that they provide the operations covered by LIA-2.

Suggestions for bindings are provided in Annex C. Annex C has partial binding examples for a number of existing programming languages and LIA-2. In addition to the bindings for the operations in LIA-2, it is also necessary to provide bindings for the maximum error parameters and big angle parameters specified by LIA-2. Annex C contains suggestions for these bindings. To conform to this standard, in the absence of a binding standard, an implementation should create a binding, following the suggestions in Annex C.

## B.3 Normative references

The referenced IEC 60559 standard is identical to the IEEE 754 standard and the former IEC 559 standard.

## B.4 Symbols and definitions

### B.4.1 Symbols

#### B.4.1.1 Sets and intervals

The interval notation is in common use. It has been chosen over the other commonly used interval notation because the chosen notation has no risk of confusion with the pair notation.

#### B.4.1.2 Operators and relations

Note that all operators, relations, and other mathematical notation used in LIA-2 is used in their conventional exact mathematical sense. They are not used to stand for operations specified by IEC 60559, LIA-1, LIA-2, or, with the exception of programme excerpts which are clearly marked, any programming language. E.g.  $x/u$  stands for the mathematically exact result of dividing  $x$  by  $u$ , whether that value is representable in any floating point datatype or not, and  $x/u \neq \text{div}_F(x, u)$  is often the case. Likewise,  $=$  is the mathematical equality, not the  $eq_F$  operation:  $0 \neq -\mathbf{0}$ , while  $eq_F(0, -\mathbf{0}) = \text{true}$ .

#### B.4.1.3 Mathematical functions

The elementary functions named  $\sin$ ,  $\cos$ , etc., used in LIA-2 are the exact mathematical functions, not any approximation. The approximations to these mathematical functions are introduced in clauses 5.3 and 5.4 and are written in a way clearly distinct from the mathematical functions. E.g.,  $\sin_F^*$ ,  $\cos_F^*$ , etc., which are unspecified mathematical functions approximating the targeted exact mathematical functions to a specified degree;  $\sin_F$ ,  $\cos_F$ , etc., which are the operations specified by LIA-2 based on the respective approximating function; **sin**, **cos**, etc., which are programming language names bound to LIA-2 operations. **sin** is thus very different from  $\sin$ .



#### B.4.1.4 Datatypes and exceptional values

The sequence types  $[I]$  and  $[F]$  appear as input datatypes to a few operations:  $max\_seq_I$ ,  $min\_seq_I$ ,  $gcd\_seq_I$ ,  $lcm\_seq_I$ ,  $max\_seq_F$ ,  $min\_seq_F$ ,  $mmax\_seq_F$ , and  $mmin\_seq_F$ .

In effect, a sequence is a finite linearly ordered collection of elements which can be indexed from 1 to the length of the sequence. Equality of two or more elements with different indices is possible. Sequences are used in LIA-2 as an abstraction of arrays, lists, other kinds of one-dimensional sequenced collections, and even variable length argument lists. As used in LIA-2 the order of the elements and number of occurrences of each element, as long as it is more than one, does not matter, so multi-sets (bags) and sets also qualify.

LIA-2 uses a modified set of exceptional values compared to LIA-1. Instead of LIA-1's **undefined**, LIA-2 uses **invalid** and **infinitary**. IEC 60559 distinguishes between **invalid** and **divide\_by\_zero** (the latter is called **infinitary** by LIA-2). The distinction is valid and should be recognised, since **infinitary** indicates that an infinite but *exact* result is (or can be, if it were available) returned, while **invalid** indicates that a result in the target datatype (extended with infinities) cannot, or should not, be returned with adequate accuracy.

LIA-1 distinguished between **integer\_overflow** and **floating\_overflow**. This distinction is moot, since no distinction was made between **integer\_undefined** and **floating\_undefined**. In addition, continuing this distinction would force LIA to start distinguishing not only **integer\_overflow** and **floating\_overflow**, but also **fixed\_overflow**, **complex\_floating\_overflow**, **complex\_integer\_overflow**, etc. Further, there is no general consensus that maintaining this distinction is useful, and many programming languages do not require a distinction. A binding standard can still maintain this distinction, if desired.

LIA allows for three methods for handing notifications: recording of indicators, change of control flow (returnable or not), and termination of program. The preferred method is recording of indicators. This allows the computation to continue using the continuation values. For **underflow** and **infinitary** notifications this course of action is strongly preferred, provided that a suitable continuation value can be represented in the result datatype.

Not all occurrences of the same exceptional value need be handled the same. There may be explicit mode changes in how notifications are handled, and there may be implicit changes. E.g., **invalid** without a specified (by LIA-2 or binding) continuation value to cause change of control flow (like an Ada [7] exception), while **invalid** with a specified continuation value use recording of indicators. This should be specified by bindings or by implementations.

The operations may return any of the exceptional values **overflow**, **underflow**, **invalid**, **infinitary**, or **absolute\_precision\_underflow**. This does *not* imply that the implemented operations are to actually *return* any of these values. When these values are returned according to the LIA specification, that means that the implementation is to perform a notification handling for that exceptional value. If the notification handling is by recording of indicators, then what is actually returned by the implemented operation is the continuation value.

#### B.4.2 Definitions of terms

Note the LIA distinction between exceptional values, exceptions, and exception handling (handling of notification by non-returnable change of control flow; as in e.g. Ada). LIA exceptional values are not the same as Ada **exceptions**, nor are they the same as IEC 60559 special values.

Note also the LIA distinction between denormal and subnormal. Subnormal include zero values, while denormal does not.

## B.5 Specifications for the complex datatypes and operations

## Annex C (informative)

### Example bindings for specific languages

This annex describes how a computing system can simultaneously conform to a language standard (or publicly available specification) and to LIA-3. It contains suggestions for binding the “abstract” operations specified in LIA-3 to concrete language syntax. The format used for these example bindings in this annex is a short form version, suitable for the purposes of this annex. An actual binding is under no obligation to follow this format. An actual binding should, however, as in the bindings examples, give the LIA-3 operation name, or parameter name, bound to an identifier by the binding.

Portability of programs can be improved if two conforming LIA-3 systems using the same language agree in the manner with which they adhere to LIA-3. For instance, LIA-3 requires that the parameter *big\_angle\_r<sub>F</sub>* be provided (if any conforming radian trigonometric operations are provided), but if one system provides it by means of the identifier **BigAngle** and another by the identifier **MaxAngle**, portability is impaired. Clearly, it would be best if such names were defined in the relevant language standards or binding standards, but in the meantime, suggestions are given here to aid portability.

The following clauses are suggestions rather than requirements because the areas covered are the responsibility of the various language standards committees. Until binding standards are in place, implementors can promote “de facto” portability by following these suggestions on their own.

The languages covered in this annex are

- Ada
- C
- C++
- Fortran
- Haskell
- Java
- Common Lisp
- ISLisp
- Modula-2
- PL/I
- SML

This list is not exhaustive. Other languages and other computing devices (like ‘scientific’ calculators, ‘web script’ languages, and database ‘query languages’) are suitable for conformity to LIA-2.

In this annex, the parameters, operations, and exception behaviour of each language are examined to see how closely they fit the requirements of LIA-2. Where parameters, constants, or operations are not provided by the language, names and syntax are suggested. (Already provided syntax is marked with a ★.)

This annex describes only the language-level support for LIA-2. An implementation that wishes to conform must ensure that the underlying hardware and software is also configured to conform to LIA-2 requirements.

A complete binding for LIA-2 will include, or refer to, a binding for LIA-1. In turn, a complete binding for the LIA-1 may include, or refer to, a binding for IEC 60559. Such a joint LIA-2/LIA-

1/IEC 60559 binding should be developed as a single binding standard. To avoid conflict with ongoing development, only the LIA-2 specific portions of such a binding are exemplified in this annex.

Most language standards permit an implementation to provide, by some means, the parameters and operations required by LIA-2 that are not already part of the language. The method for accessing these additional parameters and operations depends on the implementation and language, and is not specified in LIA-2 nor exemplified in this annex. It could include external subroutine libraries; new intrinsic functions supported by the compiler; constants and functions provided as global “macros”; and so on. The actual method of access through libraries, macros, etc. should of course be given in a real binding.

Most language standards do not constrain the accuracy of elementary numerical functions, or specify the subsequent behaviour after an arithmetic notification occurs.

In the event that there is a conflict between the requirements of the language standard and the requirements of LIA-2, the language binding standard should clearly identify the conflict and state its resolution of the conflict.

## C.1 Ada

The programming language Ada is defined by ISO/IEC 8652:1995, *Information Technology – Programming Languages – Ada* [7].

An implementation should follow all the requirements of LIA-3 unless otherwise specified by this language binding.

The operations or parameters marked “†” are not part of the language and should be provided by an implementation that wishes to conform to LIA-3 for that operation or parameter. For each of the marked items a suggested identifier is provided.

The Ada datatype `Boolean` corresponds to the LIA datatype **Boolean**.

Every implementation of Ada has at least one integer datatype, and at least one floating point datatype. The notations *INT* and *FLT* are used to stand for the names of one of these datatypes (respectively) in what follows.

Ada has an overloading system, so that the same name can be used for different types of arguments to the operations.

The LIA-3 complex integer operations are listed below, along with the syntax used to invoke them:

$itimes_{i(I)}(x)$	<code>..x</code>	†
$itimes_{c(I)}(x)$	<code>..x</code>	†
$re_I(x)$	<code>re(x)</code>	†
$re_{i(I)}(x)$	<code>re(x)</code>	†
$re_{c(I)}(x)$	<code>re(x)</code>	†
$im_I(x)$	<code>im(x)</code>	†
$im_{i(I)}(x)$	<code>im(x)</code>	†
$im_{c(I)}(x)$	<code>im(x)</code>	†
$plusitimes_I(x, y)$	<code>..x y</code>	†
$neg_{i(I)}(x)$	<code>-x</code>	†
$neg_{c(I)}(x)$	<code>-x</code>	†
$conj_I(x)$	<code>conj(x)</code>	†
$conj_{i(I)}(x)$	<code>conj(x)</code>	†
$conj_{c(I)}(x)$	<code>conj(x)</code>	†
$add_{i(I)}(x, y)$	<code>x + y</code>	†

$add_{c(I)}(x, y)$	$x + y$	†
$sub_{i(I)}(x, y)$	$x - y$	†
$sub_{c(I)}(x, y)$	$x - y$	†
$mul_{i(I)}(x, y)$	$x * y$	†
$mul_{c(I)}(x, y)$	$x * y$	†
$.eq_{i(I)}(x, y)$	$x = y$	†
$eq_{c(I)}(x, y)$	$x = y$	†
$.neq_{i(I)}(x, y)$	$x /= y$	†
$neq_{c(I)}(x, y)$	$x /= y$	†
$.abs_{i(I)}(x)$	$abs(x)$	†
$.sign_{i(I)}(x)$	$sign(x)$	†

where  $x$  and  $y$  are expressions of type *CINT*.

The LIA-3 basic complex floating point operations are listed below, along with the syntax used to invoke them:

$itimes_F(x)$	$i * x$ or $j * x$	†
$re_F(x)$	$Re(x)$	†
$re_{i(F)}(x)$	$Re(x)$	†
$re_{c(F)}(x)$	$Re(x)$	*
$im_F(x)$	$Im(x)$	†
$im_{i(F)}(x)$	$Im(x)$	†
$im_{c(F)}(x)$	$Im(x)$	*
$plusitimes_F(x, y)$	$Compose\_From\_Cartesian(x, y)$	*
$plusitimes_F(x, y)$	$x + i * y$ or $x + j * y$	*
$..._F(x, y)$	$Compose\_From\_Polar(x, y)$	*
$..._F(u, x, y)$	$Compose\_From\_Polar(x, y, u)$	*
$neg_{i(F)}(x)$	$- x$	*
$neg_{c(F)}(x)$	$- x$	*
$conj_F(x)$	$Conjugate(x)$	†
$conj_{i(F)}(x)$	$Conjugate(x)$	*
$conj_{c(F)}(x)$	$Conjugate(x)$	*
$.add_{F,i(F)}(x, y)$	$x + y$	*
$.add_{F,c(F)}(x, y)$	$x + y$	*
$.add_{i(F),F}(x, y)$	$x + y$	*
$.add_{c(F),F}(x, y)$	$x + y$	*
$.add_{i(F),c(F)}(x, y)$	$x + y$	*
$.add_{c(F),i(F)}(x, y)$	$x + y$	*
$add_{i(F)}(x, y)$	$x + y$	*
$add_{c(F)}(x, y)$	$x + y$	*
$.sub_{F,i(F)}(x, y)$	$x - y$	*
$.sub_{F,c(F)}(x, y)$	$x - y$	*
$.sub_{i(F),F}(x, y)$	$x - y$	*
$.sub_{c(F),F}(x, y)$	$x - y$	*
$.sub_{i(F),c(F)}(x, y)$	$x - y$	*
$.sub_{c(F),i(F)}(x, y)$	$x - y$	*
$sub_{i(F)}(x, y)$	$x - y$	*
$sub_{c(F)}(x, y)$	$x - y$	*
$.mul_{F,i(F)}(x, y)$	$x * y$	*
$.mul_{F,c(F)}(x, y)$	$x * y$	*
$.mul_{i(F),F}(x, y)$	$x * y$	*
$.mul_{c(F),F}(x, y)$	$x * y$	*

$.mul_{i(F),c(F)}(x, y)$	$x * y$	*
$.mul_{c(F),i(F)}(x, y)$	$x * y$	*
$mul_{i(F)}(x, y)$	$x * y$	*
$.div_{F,i(F)}(x, y)$	$x / y$	*
$.div_{F,c(F)}(x, y)$	$x / y$	*
$.div_{i(F),F}(x, y)$	$x / y$	*
$.div_{c(F),F}(x, y)$	$x / y$	*
$.div_{i(F),c(F)}(x, y)$	$x / y$	*
$.div_{c(F),i(F)}(x, y)$	$x / y$	*
$div_{i(F)}(x, y)$	$x / y$	*
$eq_{i(F)}(x, y)$	$x = y$	*
$eq_{c(F)}(x, y)$	$x = y$	*
$neq_{i(F)}(x, y)$	$x \neq y$	*
$neq_{c(F)}(x, y)$	$x \neq y$	*
...lt, gt, leq, geq on imaginary...		*
$abs_{i(F)}(x)$	$abs(x)$	*
$abs_{c(F)}(x)$	$abs(x)$ or $Modulus(x)$	*
$phase_F(x)$	$Argument(x)$	†
$phase_{i(F)}(x)$	$Argument(x)$	*
$phase_{c(F)}(x)$	$Argument(x)$	*
$.phaseu_{c(F)}(u, x)$	$Argument(x, u)$	*
$sign_{i(F)}(x)$	$Sign(x)$	†
$sign_{c(F)}(x)$	$Sign(x)$	†

where  $x$ ,  $y$ , and  $z$  are expressions of type *CFLT*.

The parameters for LIA-3 operations approximating real complex valued functions can be accessed by the following syntax:

$max\_err\_mul_{c(F)}$	$Err\_Mul(x)$	†
$max\_err\_div_{c(F)}$	$Err\_Div(x)$	†
$max\_err\_exp_{c(F)}$	$Err\_Exp(x)$	†
$max\_err\_power_{c(F)}$	$Err\_Power(x)$	†
$max\_err\_sin_{c(F)}$	$Err\_Sin(x)$	†
$max\_err\_tan_{c(F)}$	$Err\_Tan(x)$	†

where  $x$  is an expression of type *CFLT*. Several of the parameter functions are constant for each type (and library), the argument is then used only to differentiate among the floating point types.

The LIA-3 elementary floating point operations are listed below, along with the syntax used to invoke them:

$mul_{c(F)}(x, y)$	$x * y$	*
$div_{c(F)}(x, y)$	$x / y$	*
$.exp_{i(F)}(x)$	$Exp(x)$	*
$exp_{c(F)}(x)$	$Exp(x)$	*
$.power_{i(F),I}(b, a)$	$b ** a$	*
$.power_{c(F),I}(x, a)$	$x ** a$	*
$.power_{c(F),F}(b, y)$	$b ** y$	*
$.power_{F,c(F)}(b, x)$	$b ** x$	*
$power_{c(F)}(x, y)$	$x ** y$	*
$sqrt_{c(F)}(x)$	$Sqrt(x)$	*

$\ln_{c(F)}(x)$	$\text{Log}(x)$	★
$\text{.logbase}_{c(F)}(b, x)$	$\text{Log}(x, b)$ (note parameter order)	†
$\text{radh}_F(x)$	$\text{RadH}(x)$	†
$\text{radh}_{i(F)}(x)$	$\text{RadH}(x)$	†
$\text{radh}_{c(F)}(x)$	$\text{RadH}(x)$	†
$\text{sinh}_{i(F)}(x)$	$\text{SinH}(x)$	†
$\text{sinh}_{c(F)}(x)$	$\text{SinH}(x)$	★
$\text{cosh}_{i(F)}(x)$	$\text{CosH}(x)$	†
$\text{cosh}_{c(F)}(x)$	$\text{CosH}(x)$	★
$\text{tanh}_{i(F)}(x)$	$\text{TanH}(x)$	†
$\text{tanh}_{c(F)}(x)$	$\text{TanH}(x)$	★
$\text{coth}_{i(F)}(x)$	$\text{CotH}(x)$	†
$\text{coth}_{c(F)}(x)$	$\text{CotH}(x)$	★
$\text{sech}_{i(F)}(x)$	$\text{SecH}(x)$	†
$\text{sech}_{c(F)}(x)$	$\text{SecH}(x)$	†
$\text{csch}_{i(F)}(x)$	$\text{Csch}(x)$	†
$\text{csch}_{c(F)}(x)$	$\text{Csch}(x)$	†
$\text{arcsinh}_{i(F)}(x)$	$\text{ArcSinH}(x)$	†
$\text{arcsinh}_{c(F)}(x)$	$\text{ArcSinH}(x)$	★
$\text{arccosh}_{c(F)}(x)$	$\text{ArcCosH}(x)$	★
$\text{arctanh}_{i(F)}(x)$	$\text{ArcTanH}(x)$	†
$\text{arctanh}_{c(F)}(x)$	$\text{ArcTanH}(x)$	★
$\text{arccoth}_{i(F)}(x)$	$\text{ArcCotH}(x)$	†
$\text{arccoth}_{c(F)}(x)$	$\text{ArcCotH}(x)$	★
$\text{arcsech}_{c(F)}(x)$	$\text{ArcSecH}(x)$	†
$\text{arccsch}_{i(F)}(x)$	$\text{ArcCsch}(x)$	†
$\text{arccsch}_{c(F)}(x)$	$\text{ArcCsch}(x)$	†
$\text{rad}_{i(F)}(x)$	$\text{Rad}(x)$	†
$\text{rad}_{c(F)}(x)$	$\text{Rad}(x)$	†
$\text{sin}_{i(F)}(x)$	$\text{Sin}(x)$	†
$\text{sin}_{c(F)}(x)$	$\text{Sin}(x)$	★
$\text{cos}_{i(F)}(x)$	$\text{Cos}(x)$	†
$\text{cos}_{c(F)}(x)$	$\text{Cos}(x)$	★
$\text{tan}_{i(F)}(x)$	$\text{Tan}(x)$	†
$\text{tan}_{c(F)}(x)$	$\text{Tan}(x)$	★
$\text{cot}_{i(F)}(x)$	$\text{Cot}(x)$	†
$\text{cot}_{c(F)}(x)$	$\text{Cot}(x)$	★
$\text{sec}_{i(F)}(x)$	$\text{Sec}(x)$	†
$\text{sec}_{c(F)}(x)$	$\text{Sec}(x)$	†
$\text{csc}_{i(F)}(x)$	$\text{Csc}(x)$	†
$\text{csc}_{c(F)}(x)$	$\text{Csc}(x)$	†
$\text{arcsin}_{i(F)}(x)$	$\text{ArcSin}(x)$	†
$\text{arcsin}_{c(F)}(x)$	$\text{ArcSin}(x)$	★
$\text{arccos}_{c(F)}(x)$	$\text{ArcCos}(x)$	†
$\text{arctan}_{i(F)}(x)$	$\text{ArcTan}(x)$	★

$arctan_{c(F)}(x)$	<code>ArcTan(x)</code>	†
$arccot_{i(F)}(x)$	<code>ArcCot(x)</code>	*
$arccot_{c(F)}(x)$	<code>ArcCot(x)</code>	†
$arcsec_{c(F)}(x)$	<code>ArcSec(x)</code>	†
$arccsc_{i(F)}(x)$	<code>ArcCsc(x)</code>	†
$arccsc_{c(F)}(x)$	<code>ArcCsc(x)</code>	†

where  $x$  and  $y$  are expressions of type *CFLT*.

Arithmetic value conversions in Ada are always explicit.

$convert_{I \rightarrow c(I)}(x)$	<code>Compose_From_Cartesian(x)</code>	†
$convert_{i(I) \rightarrow c(I)}(x)$	<code>Compose_From_Cartesian(x)</code>	†
$convert_{F \rightarrow c(F)}(x)$	<code>Compose_From_Cartesian(x)</code>	*
$convert_{F \rightarrow c(F)}(x)$	$x - i * 0$ or $x - j * 0$	*
$convert_{i(F) \rightarrow c(F)}(x)$	<code>Compose_From_Cartesian(x)</code>	*
$convert_{i(F) \rightarrow c(F)}(x)$	$-0 + x$	*
complex IO...		

Numerals...:

$imaginary\_unit_F$	<code>i</code> or <code>j</code>	*
---------------------	----------------------------------	---

## C.2 C

The programming language C is defined by ISO/IEC 9899:1999, *Information technology – Programming languages – C* [13].

An implementation should follow all the requirements of LIA-3 unless otherwise specified by this language binding.

The operations or parameters marked “†” are not part of the language and should be provided by an implementation that wishes to conform to the LIA-3 for that operation. For each of the marked items a suggested identifier is provided.

The LIA datatype **Boolean** is implemented by the C datatype `int` ( $1 = \text{true}$  and  $0 = \text{false}$ ), or the new `_Bool` datatype.

Every implementation of C has integral datatypes `int`, `long int`, `unsigned int`, and `unsigned long int`. *INT* is used below to designate one of the integer datatypes.

C99 has three complex floating point datatypes: `_Complex float`, `_Complex double`, and `_Complex long double`. *CFLT* is used below to designate one of the complex floating point datatypes. These datatypes are, however, not required for all C99 implementations.

The LIA-3 complex integer operations are listed below, along with the syntax used to invoke them:

$itimes_{i(I)}(x)$	<code>..x</code>	†
$itimes_{c(I)}(x)$	<code>..x</code>	†
$re_I(x)$	<code>re(x)</code>	†
$re_{i(I)}(x)$	<code>re(x)</code>	†
$re_{c(I)}(x)$	<code>re(x)</code>	†
$im_I(x)$	<code>im(x)</code>	†
$im_{i(I)}(x)$	<code>im(x)</code>	†
$im_{c(I)}(x)$	<code>im(x)</code>	†
$plusitimes_I(x, y)$	<code>..x y</code>	†
$neg_{i(I)}(x)$	<code>-x</code>	†
$neg_{c(I)}(x)$	<code>-x</code>	†



$conj_I(x)$	<code>conj(x)</code>	†
$conj_{i(I)}(x)$	<code>conj(x)</code>	†
$conj_{c(I)}(x)$	<code>conj(x)</code>	†
$add_{i(I)}(x, y)$	<code>x + y</code>	†
$add_{c(I)}(x, y)$	<code>x + y</code>	†
$sub_{i(I)}(x, y)$	<code>x - y</code>	†
$sub_{c(I)}(x, y)$	<code>x - y</code>	†
$mul_{i(I)}(x, y)$	<code>x * y</code>	†
$mul_{c(I)}(x, y)$	<code>x * y</code>	†
$.eq_{i(I)}(x, y)$	<code>x = y</code>	†
$eq_{c(I)}(x, y)$	<code>x = y</code>	†
$.neq_{i(I)}(x, y)$	<code>x /= y</code>	†
$neq_{c(I)}(x, y)$	<code>x /= y</code>	†
$.abs_{i(I)}(x)$	<code>abs(x)</code>	†
$.sign_{i(I)}(x)$	<code>sign(x)</code>	†

where  $x$  and  $y$  are expressions of type *CINT*. The LIA-3 non-transcendental complex floating point operations are listed below, along with the syntax used to invoke them:

$itimes_F(x)$	<code>i * x</code> or <code>j * x</code>	†
$re_F(x)$	<code>Re(x)</code>	†
$re_{i(F)}(x)$	<code>Re(x)</code>	†
$re_{c(F)}(x)$	<code>Re(x)</code>	*
$im_F(x)$	<code>Im(x)</code>	†
$im_{i(F)}(x)$	<code>Im(x)</code>	†
$im_{c(F)}(x)$	<code>Im(x)</code>	*
$plusitimes_F(x, y)$	<code>Compose_From_Cartesian(x, y)</code>	*
$plusitimes_F(x, y)$	<code>x + i * y</code> or <code>x + j * y</code>	*
$..._F(x, y)$	<code>Compose_From_Polar(x, y)</code>	*
$..._F(u, x, y)$	<code>Compose_From_Polar(x, y, u)</code>	*
$neg_{i(F)}(x)$	<code>- x</code>	*
$neg_{c(F)}(x)$	<code>- x</code>	*
$conj_F(x)$	<code>Conjugate(x)</code>	†
$conj_{i(F)}(x)$	<code>Conjugate(x)</code>	*
$conj_{c(F)}(x)$	<code>Conjugate(x)</code>	*
$.add_{F,i(F)}(x, y)$	<code>x + y</code>	*
$.add_{F,c(F)}(x, y)$	<code>x + y</code>	*
$.add_{i(F),F}(x, y)$	<code>x + y</code>	*
$.add_{c(F),F}(x, y)$	<code>x + y</code>	*
$.add_{i(F),c(F)}(x, y)$	<code>x + y</code>	*
$.add_{c(F),i(F)}(x, y)$	<code>x + y</code>	*
$add_{i(F)}(x, y)$	<code>x + y</code>	*
$add_{c(F)}(x, y)$	<code>x + y</code>	*
$.sub_{F,i(F)}(x, y)$	<code>x - y</code>	*
$.sub_{F,c(F)}(x, y)$	<code>x - y</code>	*
$.sub_{i(F),F}(x, y)$	<code>x - y</code>	*
$.sub_{c(F),F}(x, y)$	<code>x - y</code>	*
$.sub_{i(F),c(F)}(x, y)$	<code>x - y</code>	*
$.sub_{c(F),i(F)}(x, y)$	<code>x - y</code>	*
$sub_{i(F)}(x, y)$	<code>x - y</code>	*
$sub_{c(F)}(x, y)$	<code>x - y</code>	*
$.mul_{F,i(F)}(x, y)$	<code>x * y</code>	*

$.mul_{F,c(F)}(x, y)$	$x * y$	*
$.mul_{i(F),F}(x, y)$	$x * y$	*
$.mul_{c(F),F}(x, y)$	$x * y$	*
$.mul_{i(F),c(F)}(x, y)$	$x * y$	*
$.mul_{c(F),i(F)}(x, y)$	$x * y$	*
$mul_{i(F)}(x, y)$	$x * y$	*
$.div_{F,i(F)}(x, y)$	$x / y$	*
$.div_{F,c(F)}(x, y)$	$x / y$	*
$.div_{i(F),F}(x, y)$	$x / y$	*
$.div_{c(F),F}(x, y)$	$x / y$	*
$.div_{i(F),c(F)}(x, y)$	$x / y$	*
$.div_{c(F),i(F)}(x, y)$	$x / y$	*
$div_{i(F)}(x, y)$	$x / y$	*
$eq_{i(F)}(x, y)$	$x = y$	*
$eq_{c(F)}(x, y)$	$x = y$	*
$neq_{i(F)}(x, y)$	$x \neq y$	*
$neq_{c(F)}(x, y)$	$x \neq y$	*
...lt, gt, leq, geq on imaginary...		*
$abs_{i(F)}(x)$	$abs(x)$	*
$abs_{c(F)}(x)$	$abs(x)$ or $Modulus(x)$	*
$phase_F(x)$	$Argument(x)$	†
$phase_{i(F)}(x)$	$Argument(x)$	*
$phase_{c(F)}(x)$	$Argument(x)$	*
$.phaseu_{c(F)}(u, x)$	$Argument(x, u)$	*
$sign_{i(F)}(x)$	$Sign(x)$	†
$sign_{c(F)}(x)$	$Sign(x)$	†

where  $x$ ,  $y$  and  $z$  are expressions of the same complex floating point type.

The LIA-3 parameters for operations approximating complex real valued transcendental functions can be accessed by the following syntax:

$max\_err\_mul_{c(F)}$	$err\_cmult$	†
$max\_err\_div_{c(F)}$	$err\_cdivt$	†
$max\_err\_exp_{c(F)}$	$err\_cexpt$	†
$max\_err\_power_{c(F)}$	$err\_cpowert$	†
$max\_err\_sin_{c(F)}$	$err\_csint$	†
$max\_err\_tan_{c(F)}$	$err\_ctant$	†

where  $t$  ...

The LIA-3 elementary complex floating point operations are listed below, along with the syntax used to invoke them:

$mul_{c(F)}(x, y)$	$x * y$	*
$div_{c(F)}(x, y)$	$x / y$	*
$power_{c(F)I}(b, z)$	$cpowerit(b, z)$	†
$exp_{c(F)}(x)$	$cexpt(x)$	*
$power_{c(F)}(b, y)$	$cpowert(b, y)$	†
$pow_{c(F)}(b, y)$	$cpowt(b, y)$	* Not LIA-3!
$ln_{c(F)}(x)$	$clogt(x)$	*
$.logbase_{c(F)}(b, x)$	$clogbaset(b, x)$	†

$\sinh_{i(F)}(x)$	$\text{sinht}(x)$	*
$\sinh_{c(F)}(x)$	$\text{sinht}(x)$	*
$\cosh_{i(F)}(x)$	$\text{cosht}(x)$	*
$\cosh_{c(F)}(x)$	$\text{cosht}(x)$	*
$\tanh_{i(F)}(x)$	$\text{tanht}(x)$	*
$\tanh_{c(F)}(x)$	$\text{tanht}(x)$	*
$\coth_{i(F)}(x)$	$\text{cotht}(x)$	†
$\coth_{c(F)}(x)$	$\text{cotht}(x)$	†
$\text{sech}_{i(F)}(x)$	$\text{secht}(x)$	†
$\text{sech}_{c(F)}(x)$	$\text{secht}(x)$	†
$\text{csch}_{i(F)}(x)$	$\text{cscht}(x)$	†
$\text{csch}_{c(F)}(x)$	$\text{cscht}(x)$	†
$\text{radh}_F(x)$	$\text{radianht}(x)$	†
$\text{radh}_{i(F)}(x)$	$\text{radianht}(x)$	†
$\text{radh}_{c(F)}(x)$	$\text{radianht}(x)$	†
$\text{arcsinh}_{i(F)}(x)$	$\text{asinht}(x)$	*
$\text{arcsinh}_{c(F)}(x)$	$\text{asinht}(x)$	*
$\text{arccosh}_{c(F)}(x)$	$\text{acosht}(x)$	*
$\text{arctanh}_{i(F)}(x)$	$\text{atanht}(x)$	*
$\text{arctanh}_{c(F)}(x)$	$\text{atanht}(x)$	*
$\text{arccoth}_{i(F)}(x)$	$\text{acotht}(x)$	†
$\text{arccoth}_{c(F)}(x)$	$\text{acotht}(x)$	†
$\text{arcsech}_{c(F)}(x)$	$\text{asecht}(x)$	†
$\text{arccsch}_{i(F)}(x)$	$\text{acscht}(x)$	†
$\text{arccsch}_{c(F)}(x)$	$\text{acscht}(x)$	†
$\text{rad}_{i(F)}(x)$	$\text{radian}(x)$	†
$\text{rad}_{c(F)}(x)$	$\text{radian}(x)$	†
$\sin_{i(F)}(x)$	$\text{sint}(x)$	*
$\sin_{c(F)}(x)$	$\text{sint}(x)$	*
$\cos_{i(F)}(x)$	$\text{cost}(x)$	*
$\cos_{c(F)}(x)$	$\text{cost}(x)$	*
$\tan_{i(F)}(x)$	$\text{tant}(x)$	*
$\tan_{c(F)}(x)$	$\text{tant}(x)$	*
$\cot_{i(F)}(x)$	$\text{cott}(x)$	†
$\cot_{c(F)}(x)$	$\text{cott}(x)$	†
$\sec_{i(F)}(x)$	$\text{sect}(x)$	†
$\sec_{c(F)}(x)$	$\text{sect}(x)$	†
$\csc_{i(F)}(x)$	$\text{csct}(x)$	†
$\csc_{c(F)}(x)$	$\text{csct}(x)$	†
$\text{arcsin}_{i(F)}(x)$	$\text{asint}(x)$	*
$\text{arcsin}_{c(F)}(x)$	$\text{asint}(x)$	*
$\text{arccos}_{c(F)}(x)$	$\text{acost}(x)$	*
$\text{arctan}_{i(F)}(x)$	$\text{atant}(x)$	*
$\text{arctan}_{c(F)}(x)$	$\text{atant}(x)$	*
$\text{arccot}_{i(F)}(x)$	$\text{acott}(x)$	†
$\text{arccot}_{c(F)}(x)$	$\text{acott}(x)$	†

$arcsec_{c(F)}(x)$	<code>asect(x)</code>	†
$arccsc_{i(F)}(x)$	<code>acsct(x)</code>	†
$arccsc_{c(F)}(x)$	<code>acsct(x)</code>	†

where  $b$ ,  $x$ , and  $y$  are expressions of the same complex floating point type.  $t$  is a string, part of the operation name, and is “f” for `_Complex float`, the empty string for `_Complex double`, and “l” for `_Complex long float`.

Arithmetic value conversions in C can be explicit or implicit. The explicit arithmetic value conversions are usually expressed as ‘casts’, except when converting to/from strings. The rules for when implicit conversions are applied is not repeated here, but work as if a cast had been applied.

$convert_{I \rightarrow c(I)}(x)$	<code>..(x)</code>	†
$convert_{i(I) \rightarrow c(I)}(x)$	<code>..(x)</code>	†
$convert_{F \rightarrow c(F)}(x)$	<code>x - I * 0</code> or <code>x - _IMAGINARY_I * 0</code>	*
$convert_{i(F) \rightarrow c(F)}(x)$	<code>-0 + x</code>	*
complex IO...		

where  $x$  is an expression of type *INT*,  $y$  is an expression of type *FLT*, and  $z$  is an expression of type *FXD*, where *FXD* is a fixed point type. *INT2* is the integer datatype that corresponds to *I'*. A ? above indicates that the parameter is optional.  $e$  is greater than 0.

Numerals...:

$imaginary\_unit_F$	<code>I</code> or <code>_IMAGINARY_I</code>	*
---------------------	---	---

C99 has two ways of handling arithmetic errors. One, for backwards compatibility, is by assigning to `errno`. The other is by recording of indicators, the method preferred by LIA-3, which can be used for floating point errors. For C99, the **absolute\_precision\_underflow** notification is ignored.

### C.3 C++

The programming language C++ is defined by ISO/IEC 14882:1998, *Programming languages - C++* [14].

An implementation should follow all the requirements of LIA-2 unless otherwise specified by this language binding.

The operations or parameters marked “†” are not part of the language and should be provided by an implementation that wishes to conform to the LIA-2 for that operation. For each of the marked items a suggested identifier is provided.

This example binding recommends that all identifiers suggested here be defined in the namespace `std::math`.

The LIA-1 datatype **Boolean** is implemented by the C++ datatype `bool`.

Every implementation of C++ has integral datatypes `int`, `long int`, `unsigned int`, and `unsigned long int`. *INT* is used below to designate one of the integer datatypes.

C++ has three floating point datatypes: `float`, `double`, and `long double`. *FLT* is used below to designate one of the floating point datatypes.

The LIA-2 integer operations are listed below, along with the syntax used to invoke them:

where  $x$  and  $y$  are expressions of the same integer type and where  $xs$  is an expression of type valarray of an integer type.

The LIA-2 non-transcendental floating point operations are listed below, along with the syntax used to invoke them:

where  $x$ ,  $y$  and  $z$  are expressions of the same floating point type, and where  $xs$  is an expression of type valarray of a floating point type.

The parameters for operations approximating real valued transcendental functions can be accessed by the following syntax:

$max\_err\_mul_{c(F)}$	<code>err_mul&lt;CFLT&gt;()</code>	†
$max\_err\_div_{c(F)}$	<code>err_div&lt;CFLT&gt;()</code>	†
$max\_err\_exp_{c(F)}$	<code>err_exp&lt;CFLT&gt;()</code>	†
$max\_err\_power_{c(F)}$	<code>err_power&lt;CFLT&gt;()</code>	†
$max\_err\_sin_{c(F)}$	<code>err_sin&lt;CFLT&gt;()</code>	†
$max\_err\_tan_{c(F)}$	<code>err_tan&lt;CFLT&gt;()</code>	†

where  $u$  is an expression of a floating point type. Several of the parameter functions are constant for each type (and library).

The LIA-2 elementary floating point operations are listed below, along with the syntax (type generic macros) used to invoke them:

$exp_{c(F)}(x)$	<code>exp(x)</code>	★
$power_{c(F)}(b, y)$	<code>power(b, y)</code>	†
$pow_{c(F)}(b, y)$	<code>pow(b, y)</code>	★ Not LIA-2! (See C.)
$ln_{c(F)}(x)$	<code>log(x)</code>	★
$logbase_{c(F)}(b, x)$	<code>logbase(b, x)</code>	†
$radh_{c(F)}(x)$	<code>radh(x)</code>	†
$sinh_{c(F)}(x)$	<code>sinh(x)</code>	★
$cosh_{c(F)}(x)$	<code>cosh(x)</code>	★
$tanh_{c(F)}(x)$	<code>tanh(x)</code>	★
$coth_{c(F)}(x)$	<code>coth(x)</code>	†
$sech_{c(F)}(x)$	<code>sech(x)</code>	†
$csch_{c(F)}(x)$	<code>csch(x)</code>	†
$arcsinh_{c(F)}(x)$	<code>asinh(x)</code>	★
$arccosh_{c(F)}(x)$	<code>acosh(x)</code>	★
$arctanh_{c(F)}(x)$	<code>atanh(x)</code>	★
$arcoth_{c(F)}(x)$	<code>acoth(x)</code>	†
$arcsech_{c(F)}(x)$	<code>asech(x)</code>	†
$arccsch_{c(F)}(x)$	<code>acsch(x)</code>	†
$rad_{c(F)}(x)$	<code>rad(x)</code>	†

$\sin_{c(F)}(x)$	<code>sin(x)</code>	*
$\cos_{c(F)}(x)$	<code>cos(x)</code>	*
$\tan_{c(F)}(x)$	<code>tan(x)</code>	*
$\cot_{c(F)}(x)$	<code>cot(x)</code>	†
$\sec_{c(F)}(x)$	<code>sec(x)</code>	†
$\csc_{c(F)}(x)$	<code>csc(x)</code>	†
$\arcsin_{c(F)}(x)$	<code>asin(x)</code>	*
$\arccos_{c(F)}(x)$	<code>acos(x)</code>	*
$\arctan_{c(F)}(x)$	<code>atan(x)</code>	*
$\text{arccot}_{c(F)}(x)$	<code>acot(x)</code>	†
$\text{arctg}_{c(F)}(x)$	<code>actg(x)</code>	†
$\text{arcsec}_{c(F)}(x)$	<code>asec(x)</code>	†
$\text{arccsc}_{c(F)}(x)$	<code>acsc(x)</code>	†

where  $b$ ,  $x$ ,  $y$ ,  $u$ , and  $v$  are expressions of type *CFLT*.

Arithmetic value conversions in C++ are can be explicit or implicit. The rules for when implicit conversions are applied is not repeated here. The explicit arithmetic value conversions are usually expressed as ‘casts’, except when converting to/from strings.

where  $x$  is an expression of type *INT*,  $y$  is an expression of type *FLT*, and  $z$  is an expression of type *FXD*, where *FXD* is a fixed point type. *INT?* is the integer datatype that corresponds to  $I'$ . A ? above indicates that the parameter is optional.  $e$  is greater than 0.

C++ provides non-negative numerals for all its integer and floating point types in base 10. Numerals for different integer types are distinguished by suffixes. Numerals for different floating point types are distinguished by suffix: `f` for `float`, no suffix for `double`, `l` for `long double`. Numerals for floating point types must have a ‘.’ in them. The details are not repeated in this example binding, see ISO/IEC 14882, clause 2.9.1 Integer literals, and clause 2.9.4 Floating literals.

## C.4 Fortran

The programming language Fortran is defined by ISO/IEC 1539-1:1997, *Information technology – Programming languages – Fortran – Part 1: Base language* [18].

An implementation should follow all the requirements of LIA-2 unless otherwise specified by this language binding.

The operations or parameters marked “†” are not part of the language and should be provided by an implementation that wishes to conform to the LIA-2 for that operation. For each of the marked items a suggested identifier is provided.

The Fortran datatype `LOGICAL` corresponds to the LIA datatype `Boolean`.

Every implementation of Fortran has one integer datatype, denoted as `INTEGER`, and two floating point data type denoted as `REAL` (single precision) and `DOUBLE PRECISION`.

An implementation is permitted to offer additional `INTEGER` types with a different range and additional `REAL` types with different precision or range, parameterised with the `KIND` parameter.

The LIA-2 integer operations are listed below, along with the syntax used to invoke them:

$add_{c(I)}(x, y)$	$x + y$	†
$sub_{c(I)}(x, y)$	$x - y$	†
$mul_{c(I)}(x, y)$	$x * y$	†

where  $x$  and  $y$  are expressions of type `CINTEGER` and where  $xs$  is an expression of type array of `CINTEGER`.

The additional non-transcendental floating point operations are listed below, along with the syntax used to invoke them:

$add_{c(F)}(x, y)$	$x + y$	*
$sub_{c(F)}(x, y)$	$x - y$	*
$mul_{c(F)}(x, y)$	$x * y$	*

where  $x$ ,  $y$  and  $z$  are expressions of type `FLT`, and where  $xs$  is an expression of type array of `FLT`.

The LIA-3 parameters for operations approximating real valued transcendental functions can be accessed by the following syntax:

$max\_err\_mul_F$	<code>ERR_MUL(x)</code>	†
$max\_err\_div_F$	<code>ERR_DIV(x)</code>	†
$max\_err\_exp_F$	<code>ERR_EXP(x)</code>	†
$max\_err\_power_F$	<code>ERR_POWER(x)</code>	†
$max\_err\_sin_F$	<code>ERR_SIN(x)</code>	†
$max\_err\_tan_F$	<code>ERR_TAN(x)</code>	†

where  $b$ ,  $x$  and  $u$  are expressions of type `CFLT`. Several of the parameter functions are constant for each type (and library), the argument is then used only to differentiate among the floating point types.

The LIA-3 elementary floating point operations are listed below, along with the syntax used to invoke them:

$div_{c(F)}(x, y)$	$x / y$	*
$sqrt_{c(F)}(x)$	<code>SQRT(x)</code>	*
$exp_{c(F)}(x)$	<code>EXP(x)</code>	*
$power_{c(F)}(b, y)$	$b ** y$	*
$ln_{c(F)}(x)$	<code>LOG(x)</code>	*
$logbase_{c(F)}(b, x)$	<code>LOGBASE(b, x)</code>	†
$sinh_{c(F)}(x)$	<code>SINH(x)</code>	*
$cosh_{c(F)}(x)$	<code>COSH(x)</code>	*
$tanh_{c(F)}(x)$	<code>TANH(x)</code>	*
$coth_{c(F)}(x)$	<code>COTH(x)</code>	†
$sech_{c(F)}(x)$	<code>SECH(x)</code>	†
$csch_{c(F)}(x)$	<code>CSCH(x)</code>	†
$radh_{c(F)}(x)$	<code>RADH(x)</code>	†

$\operatorname{arcsinh}_{c(F)}(x)$	$\operatorname{ASINH}(x)$	†
$\operatorname{arccosh}_{c(F)}(x)$	$\operatorname{ACOSH}(x)$	†
$\operatorname{arctanh}_{c(F)}(x)$	$\operatorname{ATANH}(x)$	†
$\operatorname{arccoth}_{c(F)}(x)$	$\operatorname{ACOTH}(x)$	†
$\operatorname{arcsech}_{c(F)}(x)$	$\operatorname{ASECH}(x)$	†
$\operatorname{arccsch}_{c(F)}(x)$	$\operatorname{ACSCH}(x)$	†
$\operatorname{rad}_{c(F)}(x)$	$\operatorname{RAD}(x)$	†
$\operatorname{sin}_{c(F)}(x)$	$\operatorname{SIN}(x)$	*
$\operatorname{cos}_{c(F)}(x)$	$\operatorname{COS}(x)$	*
$\operatorname{tan}_{c(F)}(x)$	$\operatorname{TAN}(x)$	*
$\operatorname{cot}_{c(F)}(x)$	$\operatorname{COT}(x)$	†
$\operatorname{sec}_{c(F)}(x)$	$\operatorname{SEC}(x)$	†
$\operatorname{csc}_{c(F)}(x)$	$\operatorname{CSC}(x)$	†
$\operatorname{arcsin}_{c(F)}(x)$	$\operatorname{ASIN}(x)$	*
$\operatorname{arccos}_{c(F)}(x)$	$\operatorname{ACOS}(x)$	*
$\operatorname{arctan}_{c(F)}(x)$	$\operatorname{ATAN}(x)$	*
$\operatorname{arccot}_{c(F)}(x)$	$\operatorname{ACOT}(x)$	†
$\operatorname{arcctg}_{c(F)}(x)$	$\operatorname{ACTG}(x)$	†
$\operatorname{arcsec}_{c(F)}(x)$	$\operatorname{ASEC}(x)$	†
$\operatorname{arccsc}_{c(F)}(x)$	$\operatorname{ACSC}(x)$	†

where  $b$ ,  $x$ ,  $y$ ,  $u$ , and  $v$  are expressions of type *CFLT*.

Arithmetic value conversions in Fortran are always explicit, and the conversion function is named like the target type, except when converting to/from strings.

where  $x$  is an expression of type *INT*,  $y$  is an expression of type *FLT*, and  $z$  is an expression of type *FXD*, where *FXD* is a fixed point type. *INT2* is the integer datatype that corresponds to  $I'$ .

## C.5 Haskell

The programming language Haskell is defined by *Report on the programming language Haskell 98* [56], together with *Standard libraries for the Haskell 98 programming language* [57].

An implementation should follow all the requirements of LIA-2 unless otherwise specified by this language binding.

The operations or parameters marked “†” are not part of the language and should be provided by an implementation that wishes to conform to the LIA-3 for that operation. For each of the marked items a suggested identifier is provided.

The Haskell datatype `Bool` corresponds to the LIA datatype **Boolean**.

Every implementation of Haskell has at least two integer datatypes **Integer**, which is unlimited, and **Int**, and at least two floating point datatypes, **Float**, and **Double**. The notation *INT* is used to stand for the name of one of the integer datatypes, and *FLT* is used to stand for the name of one of the floating point datatypes in what follows.

The LIA-2 integer operations are listed below, along with the syntax used to invoke them:



$add_{c(I)}(x, y)$	$x + y$ or $(+) x y$	†
$sub_{c(I)}(x, y)$	$x - y$ or $(-) x y$	†
$mul_{c(I)}(x, y)$	$x * y$ or $(*) x y$	†

where  $x$  and  $y$  are expressions of type *complex INT* [not in Haskell98].

The LIA-3 non-transcendental floating point operations are listed below, along with the syntax used to invoke them:

$add_{c(F)}(x, y)$	$x + y$ or $(+) x y$	*
$sub_{c(F)}(x, y)$	$x - y$ or $(-) x y$	*

where  $x$ ,  $y$  and  $z$  are expressions of type *complex FLT*.

The LIA-3 parameters for operations approximating real valued transcendental functions can be accessed by the following syntax:

$max\_err\_mul_F$	<code>err_mul x</code>	†
$max\_err\_div_F$	<code>err_div x</code>	†
$max\_err\_exp_F$	<code>err_exp x</code>	†
$max\_err\_power_F$	<code>err_power x</code>	†
$max\_err\_sin_F$	<code>err_sin x</code>	†
$max\_err\_tan_F$	<code>err_tan x</code>	†

where  $x$  is an expression of type *complex FLT*. Several of the parameter functions are constant for each type (and library), the argument is then used only to differentiate among the floating point types.

The LIA-3 elementary floating point operations are listed below, along with the syntax used to invoke them:

$mul_{c(F)}(x, y)$	$x * y$ or $(*) x y$	*
$sqrt_{c(F)}(x)$	<code>sqrt x</code>	*
$exp_{c(F)}(x)$	<code>exp x</code>	*
$power_{c(F)}(b, y)$	$b ** y$ or $(**) b y$	*
$ln_{c(F)}(x)$	<code>log x</code>	*
$logbase_{c(F)}(b, x)$	<code>logBase b x</code> or $b$ ‘logBase’ $x$	*
$radh_{c(F)}(x)$	<code>hypradians x</code>	†
$sinh_{c(F)}(x)$	<code>sinh x</code>	*
$cosh_{c(F)}(x)$	<code>cosh x</code>	*
$tanh_{c(F)}(x)$	<code>tanh x</code>	*
$coth_{c(F)}(x)$	<code>coth x</code>	†
$sech_{c(F)}(x)$	<code>sech x</code>	†
$csch_{c(F)}(x)$	<code>csch x</code>	†
$arcsinh_{c(F)}(x)$	<code>asinh x</code>	*
$arccosh_{c(F)}(x)$	<code>acosh x</code>	*

$\operatorname{arctanh}_{c(F)}(x)$	<code>atanh <math>x</math></code>	★
$\operatorname{arccoth}_{c(F)}(x)$	<code>acoth <math>x</math></code>	†
$\operatorname{arcsech}_{c(F)}(x)$	<code>asech <math>x</math></code>	†
$\operatorname{arccsch}_{c(F)}(x)$	<code>acsch <math>x</math></code>	†
$\operatorname{rad}_{c(F)}(x)$	<code>radians <math>x</math></code>	†
$\operatorname{sin}_{c(F)}(x)$	<code>sin <math>x</math></code>	★
$\operatorname{cos}_{c(F)}(x)$	<code>cos <math>x</math></code>	★
$\operatorname{tan}_{c(F)}(x)$	<code>tan <math>x</math></code>	★
$\operatorname{cot}_{c(F)}(x)$	<code>cot <math>x</math></code>	†
$\operatorname{sec}_{c(F)}(x)$	<code>sec <math>x</math></code>	†
$\operatorname{csc}_{c(F)}(x)$	<code>csc <math>x</math></code>	†
$\operatorname{arcsin}_{c(F)}(x)$	<code>asin <math>x</math></code>	★
$\operatorname{arccos}_{c(F)}(x)$	<code>acos <math>x</math></code>	★
$\operatorname{arctan}_{c(F)}(x)$	<code>atan <math>x</math></code>	★
$\operatorname{arccot}_{c(F)}(x)$	<code>acot <math>x</math></code>	†
$\operatorname{arcctg}_{c(F)}(x)$	<code>actg <math>x</math></code>	†
$\operatorname{arcsec}_{c(F)}(x)$	<code>asec <math>x</math></code>	†
$\operatorname{arccsc}_{c(F)}(x)$	<code>acsc <math>x</math></code>	†

where  $b$ ,  $x$ ,  $y$ ,  $u$ , and  $v$  are expressions of type *complex FLT*.

Arithmetic value conversions in Haskell are always explicit. They are done with the overloaded `fromIntegral` and `fromFractional` operations.

where  $x$  is an expression of type *INT*,  $y$  is an expression of type *FLT*, and  $z$  is an expression of type *FXD*, where *FXD* is a fixed point type.

## C.6 Java

The programming language Java is defined by *The Java Language Specification* [55].

An implementation should follow all the requirements of LIA-3 unless otherwise specified by this language binding.

The operations or parameters marked “†” are not part of the language and should be provided by an implementation that wishes to conform to the LIA-3 for that operation. For each of the marked items a suggested identifier is provided. The LIA-3 operations that are provided in Java 2 (marked “★” below) are in the final class `java.lang.Math` [not in Java 2...].

The Java datatype `boolean` corresponds to the LIA datatype **Boolean**.

Every implementation of Java has the integral datatypes `int`, and `long`.

Java has two floating point datatypes, `float` and `double`, which must conform to IEC 60559.

The LIA-3 integer operations are listed below, along with the syntax used to invoke them:

$\operatorname{add}_{c(I)}(x, y)$	<code>add(<math>x</math>, <math>y</math>)</code>	★
$\operatorname{sub}_{c(I)}(x, y)$	<code>sub(<math>x</math>, <math>y</math>)</code>	★
$\operatorname{mul}_{c(I)}(x, y)$	<code>mul(<math>x</math>, <math>y</math>)</code>	★

where  $x$  and  $y$  are expressions of type *CINT*.

The LIA-3 non-transcendental floating point operations are listed below, along with the syntax used to invoke them:

$add_{c(F)}(x, y)$	<code>add(x, y)</code>	★
$sub_{c(F)}(x, y)$	<code>sub(x, y)</code>	★

where  $x$ ,  $y$  and  $z$  are expressions of type *FLT*, and where  $xs$  is an expression of type **array of FLT**.

The LIA-2 parameters for operations approximating real valued transcendental functions can be accessed by the following syntax:

$max\_err\_mul_{c(F)}$	<code>err_mul(x)</code>	†
$max\_err\_div_{c(F)}$	<code>err_div(x)</code>	†
$max\_err\_exp_{c(F)}$	<code>err_exp(x)</code>	†
$max\_err\_power_{c(F)}(b, x)$	<code>err_power(b, x)</code>	†
$max\_err\_sin_{c(F)}$	<code>err_sin(x)</code>	†
$max\_err\_tan_{c(F)}$	<code>err_tan(x)</code>	†

where  $b$ ,  $x$  and  $u$  are expressions of type *FLT*. Several of the parameter functions are constant for each type (and library), the argument is then used only to differentiate among the floating point types.

The LIA-2 elementary floating point operations are listed below, along with the syntax used to invoke them. These are defined only for **double** not for **float**.

$mul_{c(F)}(x, y)$	<code>mul(x, y)</code>	★
$div_{c(F)}(x, y)$	<code>div(x, y)</code>	★
$exp_{c(F)}(x)$	<code>exp(x)</code>	★
$power_{c(F)}(b, y)$	<code>power(b, y)</code>	†
$pow_{c(F)}(b, y)$	<code>pow(b, y)</code>	★ Not LIA-2!
$ln_{c(F)}(x)$	<code>log(x)</code>	★
$logbase_{c(F)}(b, x)$	<code>log(b, x)</code>	†
$radh_{c(F)}(x)$	<code>hypradian(x)</code>	†
$sinh_{c(F)}(x)$	<code>sinh(x)</code>	†
$cosh_{c(F)}(x)$	<code>cosh(x)</code>	†
$tanh_{c(F)}(x)$	<code>tanh(x)</code>	†
$coth_{c(F)}(x)$	<code>coth(x)</code>	†
$sech_{c(F)}(x)$	<code>sech(x)</code>	†
$csch_{c(F)}(x)$	<code>csch(x)</code>	†
$arcsinh_{c(F)}(x)$	<code>asinh(x)</code>	†
$arccosh_{c(F)}(x)$	<code>acosh(x)</code>	†
$arctanh_{c(F)}(x)$	<code>atanh(x)</code>	†
$arcoth_{c(F)}(x)$	<code>acoth(x)</code>	†
$arcsech_{c(F)}(x)$	<code>asech(x)</code>	†

$arccsch_{c(F)}(x)$	<code>acsch(x)</code>	†
$rad_{c(F)}(x)$	<code>radian(x)</code>	†
$sin_{c(F)}(x)$	<code>sin(x)</code>	★
$cos_{c(F)}(x)$	<code>cos(x)</code>	★
$tan_{c(F)}(x)$	<code>tan(x)</code>	★
$cot_{c(F)}(x)$	<code>cot(x)</code>	†
$sec_{c(F)}(x)$	<code>sec(x)</code>	†
$csc_{c(F)}(x)$	<code>csc(x)</code>	†
$arcsin_{c(F)}(x)$	<code>asin(x)</code>	★
$arccos_{c(F)}(x)$	<code>acos(x)</code>	★
$arctan_{c(F)}(x)$	<code>atan(x)</code>	★
$arccot_{c(F)}(x)$	<code>acot(x)</code>	†
$arcctg_{c(F)}(x)$	<code>actg(x)</code>	†
$arcsec_{c(F)}(x)$	<code>asec(x)</code>	†
$arccsc_{c(F)}(x)$	<code>acsc(x)</code>	†

where  $b$ ,  $x$ ,  $y$ ,  $u$ , and  $v$  are expressions of type *CFLT*.

Arithmetic value conversions in Java can be explicit or implicit. The rules for when implicit conversions are applied is not repeated here. The explicit arithmetic value conversions are usually expressed as ‘casts’, except when converting to/from strings.

where  $x$  is an expression of type *INT*,  $y$  is an expression of type *FLT*, and  $z$  is an expression of type *FXD*, where *FXD* is a fixed point type. *INT2* is the integer datatype that corresponds to  $I'$ . A ? above indicates that the parameter is optional.  $e$  is greater than 0.

## C.7 Common Lisp

The programming language Common Lisp is defined by ANSI X3.226-1994, *Information Technology – Programming Language – Common Lisp* [38].

An implementation should follow all the requirements of LIA-3 unless otherwise specified by this language binding.

The operations or parameters marked “†” are not part of the language and should be provided by an implementation that wishes to conform to the LIA-2 for that operation. For each of the marked items a suggested identifier is provided.

Common Lisp does not have a single datatype that corresponds to the LIA-1 datatype **Boolean**. Rather, NIL corresponds to **false** and T corresponds to **true**.

Every implementation of Common Lisp has one unbounded integer datatype. Any mathematical integer is assumed to have a representation as a Common Lisp data object, subject only to total memory limitations.

Common Lisp has four floating point types: **short-float**, **single-float**, **double-float**, and **long-float**. Not all of these floating point types must be distinct.

The additional integer operations are listed below, along with the syntax used to invoke them:

$add_{c(I)}(x, y)$	<code>(+ x y)</code>	★
$sub_{c(I)}(x, y)$	<code>(- x y)</code>	★
$mul_{c(I)}(x, y)$	<code>(* x y)</code>	★

where  $x$  and  $y$  are expressions of type *CINT*.

The LIA-3 non-transcendental floating point operations are listed below, along with the syntax used to invoke them:

$add_{c(F)}(x, y)$	<code>(+ x y)</code>	★
$sub_{c(F)}(x, y)$	<code>(- x y)</code>	★

where  $x$ ,  $y$  and  $z$  are data objects of the same floating point type, and where  $xs$  is a data objects that is a list of data objects of (the same, in this binding) floating point type. Note that Common Lisp allows mixed number types in many of its operations. This example binding does not explain that in detail.

The LIA-2 parameters for operations approximating real valued transcendental functions can be accessed by the following syntax:

$max\_err\_mul_F$	<code>(err-mul x)</code>	†
$max\_err\_div_F$	<code>(err-div x)</code>	†
$max\_err\_exp_F$	<code>(err-exp x)</code>	†
$max\_err\_power_F$	<code>(err-power x)</code>	†
$max\_err\_sin_F$	<code>(err-sin x)</code>	†
$max\_err\_tan_F$	<code>(err-tan x)</code>	†

where  $b$ ,  $x$  and  $u$  are expressions of type *CFLT*. Several of the parameter functions are constant for each type (and library), the argument is then used only to differentiate among the floating point types.

The LIA-3 elementary floating point operations are listed below, along with the syntax used to invoke them:

$mul_{c(F)}(x, y)$	<code>(* x y)</code>	★
$div_{c(F)}(x, y)$	<code>(/ x y)</code>	★
$sqrt_{c(F)}(x)$	<code>(sqrt x)</code>	★
$exp_{c(F)}(x)$	<code>(exp x)</code>	★
$power_{c(F)}(b, y)$	<code>(expt b y)</code> (deviation: <code>(expt 0.0 0.0)</code> is 1)	★
$ln_{c(F)}(x)$	<code>(log x)</code>	★
$logbase_{c(F)}(b, x)$	<code>(log x b)</code> (note parameter order)	★
$radh_{c(F)}(x)$	<code>(hypradians x)</code>	†
$sinh_{c(F)}(x)$	<code>(sinh x)</code>	★
$cosh_{c(F)}(x)$	<code>(cosh x)</code>	★
$tanh_{c(F)}(x)$	<code>(tanh x)</code>	★
$coth_F(x)$	<code>(coth x)</code>	†

$sech_F(x)$	(sech $x$ )	†
$csch_{c(F)}(x)$	(csch $x$ )	†
$arcsinh_{c(F)}(x)$	(asinh $x$ )	*
$arccosh_{c(F)}(x)$	(acosh $x$ )	*
$arctanh_{c(F)}(x)$	(atanh $x$ )	*
$arccoth_{c(F)}(x)$	(acoth $x$ )	†
$arcsech_{c(F)}(x)$	(asech $x$ )	†
$arccsch_{c(F)}(x)$	(acsch $x$ )	†
$rad_{c(F)}(x)$	(radians $x$ )	†
$sin_{c(F)}(x)$	(sin $x$ )	*
$cos_{c(F)}(x)$	(cos $x$ )	*
$tan_{c(F)}(x)$	(tan $x$ )	*
$cot_{c(F)}(x)$	(cot $x$ )	†
$sec_{c(F)}(x)$	(sec $x$ )	†
$csc_{c(F)}(x)$	(csc $x$ )	†
$arcsin_{c(F)}(x)$	(asin $x$ )	*
$arccos_{c(F)}(x)$	(acos $x$ )	*
$arctan_{c(F)}(x)$	(atan $x$ )	*
$arccot_{c(F)}(x)$	(acot $x$ )	†
$arcctg_{c(F)}(x)$	(actg $x$ )	†
$arcsec_{c(F)}(x)$	(asec $x$ )	†
$arccsc_{c(F)}(x)$	(acsc $x$ )	†

where  $b$ ,  $x$ ,  $y$ ,  $u$ , and  $v$  are expressions of type *CFLT*.

Arithmetic value conversions in Common Lisp are can be explicit or implicit. The rules for when implicit conversions are done is implementation defined.

where  $x$  is an expression of type *INT*,  $y$  is an expression of type *FLT*, and  $z$  is an expression of type *FXD*, where *FXD* is a fixed point type. Conversion from string to numeric value is in Common Lisp done via a general read procedure, which reads Common Lisp ‘S-expressions’.

## C.8 ISLisp

The programming language ISLisp is defined by ISO/IEC 13816:1997, *Information technology – Programming languages, their environments and system software interfaces – Programming language ISLISP* [20].

An implementation should follow all the requirements of LIA-2 unless otherwise specified by this language binding.

The operations or parameters marked “†” are not part of the language and should be provided by an implementation that wishes to conform to the LIA-2 for that operation. For each of the marked items a suggested identifier is provided.

ISLisp does not have a single datatype that corresponds to the LIA datatype **Boolean**. Rather, NIL corresponds to **false** and T corresponds to **true**.

Every implementation of ISLisp has one unbounded integer datatype. Any mathematical integer is assumed to have a representation as a ISLisp data object, subject only to total memory limitations.

ISLisp has one floating point type required to conform to IEC 60559.

The additional integer operations are listed below, along with the syntax used to invoke them:

$min_I(x, y)$	(min $x$ $y$ )	*
$max_I(x, y)$	(max $x$ $y$ )	*
$min\_seq_I(xs)$	(min . $xs$ ) or (min $x_1$ $x_2$ ... $x_n$ )	*
$max\_seq_I(xs)$	(max . $xs$ ) or (max $x_1$ $x_2$ ... $x_n$ )	*
$dim_I(x, y)$	(dim $x$ $y$ )	†
$sqrt_I(x)$	(isqrt $x$ )	*
$power_I(x, y)$	(expt $x$ $y$ ) (deviation: (expt 0 0) is 1)	*
$shift2_I(x, y)$	(shift2 $x$ $y$ )	†
$shift10_I(x, y)$	(shift10 $x$ $y$ )	†
$divides_I(x, y)$	(dividesp $x$ $y$ )	†
$even_I(x)$	(evenp $x$ )	†
$odd_I(x)$	(oddp $x$ )	†
$div_I(x, y)$	(div $x$ $y$ )	*
$moda_I(x, y)$	(mod $x$ $y$ )	*
$group_I(x, y)$	(group $x$ $y$ )	†
$pad_I(x, y)$	(pad $x$ $y$ )	†
$quot_I(x, y)$	(quot $x$ $y$ )	†
$remr_I(x, y)$	(remainder $x$ $y$ )	†
$gcd_I(x, y)$	(gcd $x$ $y$ ) (deviation: (gcd 0 0) is 0)	*
$lcm_I(x, y)$	(lcm $x$ $y$ )	*
$gcd\_seq_I(xs)$	(gcds $xs$ )	†
$lcm\_seq_I(xs)$	(lcms $xs$ )	†
$add\_wrap_I(x, y)$	(add_wrap $x$ $y$ )	†
$add\_ov_I(x, y)$	(add_over $x$ $y$ )	†
$sub\_wrap_I(x, y)$	(sub_wrap $x$ $y$ )	†
$sub\_ov_I(x, y)$	(sub_over $x$ $y$ )	†
$mul\_wrap_I(x, y)$	(mul_wrap $x$ $y$ )	†
$mul\_ov_I(x, y)$	(mul_over $x$ $y$ )	†

where  $x$  and  $y$  are expressions of type *INT* and where  $xs$  is an expression of type list of *INT*.

The LIA-2 non-transcendental floating point operations are listed below, along with the syntax used to invoke them:

$min_F(x, y)$	(min $x$ $y$ )	*
$max_F(x, y)$	(max $x$ $y$ )	*
$mmin_F(x, y)$	(mmin $x$ $y$ )	†
$mmax_F(x, y)$	(mmax $x$ $y$ )	†
$min\_seq_F(xs)$	(min . $xs$ ) or (min $x_1$ $x_2$ ... $x_n$ )	*
$max\_seq_F(xs)$	(max . $xs$ ) or (max $x_1$ $x_2$ ... $x_n$ )	*
$mmin\_seq_F(xs)$	(mmin . $xs$ ) or (mmin $x_1$ $x_2$ ... $x_n$ )	†
$mmax\_seq_F(xs)$	(mmax . $xs$ ) or (mmax $x_1$ $x_2$ ... $x_n$ )	†

$\mathit{floor}_F(x)$	(float (floor $x$ ))	★
$\mathit{rounding}_F(x)$	(float (round $x$ ))	★
$\mathit{ceiling}_F(x)$	(float (ceiling $x$ ))	★
$\mathit{dim}_F(x, y)$	(dim $x y$ )	†
$\mathit{dprod}_{F \rightarrow F'}(x, y)$	(prod $x y$ )	†
$\mathit{remr}_F(x, y)$	(remainder $x y$ )	†
$\mathit{sqr}_F(x)$	(sqrt $x$ )	★
$\mathit{rsqr}_F(x)$	(rsqrt $x$ )	†
$\mathit{add\_lo}_F(x, y)$	(add_low $x y$ )	†
$\mathit{sub\_lo}_F(x, y)$	(sub_low $x y$ )	†
$\mathit{mul\_lo}_F(x, y)$	(mul_low $x y$ )	†
$\mathit{div\_rest}_F(x, y)$	(div_rest $x y$ )	†
$\mathit{sqr\_rest}_F(x)$	(sqr_rest $x$ )	†

where  $x$ ,  $y$  and  $z$  are data objects of the same floating point type, and where  $xs$  is an data objects that are lists of data objects of the same floating point type.

The LIA-2 parameters for operations approximating real valued transcendental functions can be accessed by the following syntax:

$\mathit{max\_err\_hypot}_F$	(err-hypotenuse $x$ )	†
$\mathit{max\_err\_exp}_F$	(err-exp $x$ )	†
$\mathit{max\_err\_power}_F$	(err-power $x$ )	†
$\mathit{max\_err\_sinh}_F$	(err-sinh $x$ )	†
$\mathit{max\_err\_tanh}_F$	(err-tanh $x$ )	†
$\mathit{big\_radian\_angle}_F$	(big-radian-angle $x$ )	†
$\mathit{max\_err\_sin}_F$	(err-sin $x$ )	†
$\mathit{max\_err\_tan}_F$	(err-tan $x$ )	†
$\mathit{min\_angular\_unit}_F$	(minimum-angular-unit $x$ )	†
$\mathit{big\_angle}_F$	(big-angle $x$ )	†
$\mathit{max\_err\_sinu}_F(u)$	(err-sin-cycle $u$ )	†
$\mathit{max\_err\_tanu}_F(u)$	(err-tan-cycle $u$ )	†
$\mathit{max\_err\_convert}_F$	err-convert-to-string	†
$\mathit{max\_err\_convert}_D$	err-convert-to-string	†

where  $b$ ,  $x$  and  $u$  are expressions of type *FLT*. Several of the parameter functions are constant for each type (and library), the argument is then used only to differentiate among the floating point types.

The LIA-2 elementary floating point operations are listed below, along with the syntax used to invoke them:

$\mathit{hypot}_F(x, y)$	(hypotenuse $x y$ )	†
$\mathit{power}_{FI}(b, z)$	(expt $b z$ )	★



$exp_F(x)$	(exp $x$ )	★
$expm1_F(x)$	(expm1 $x$ )	†
$exp2_F(x)$	(exp2 $x$ )	†
$exp10_F(x)$	(exp10 $x$ )	†
$power_F(b, y)$	(expt $b y$ )	★
$power1pm1_F(b, y)$	(expm1 $b y$ )	†
$ln_F(x)$	(log $x$ )	★
$ln1p_F(x)$	(log1p $x$ )	†
$log2_F(x)$	(log2 $x$ )	†
$log10_F(x)$	(log10 $x$ )	†
$logbase_F(b, x)$	(logbase $b x$ )	†
$logbase1p1p_F(b, x)$	(logbase1p $b x$ )	†
$sinh_F(x)$	(sinh $x$ )	★
$cosh_F(x)$	(cosh $x$ )	★
$tanh_F(x)$	(tanh $x$ )	★
$coth_F(x)$	(coth $x$ )	†
$sech_F(x)$	(sech $x$ )	†
$csch_F(x)$	(csch $x$ )	†
$arcsinh_F(x)$	(asinh $x$ )	†
$arccosh_F(x)$	(acosh $x$ )	†
$arctanh_F(x)$	(atanh $x$ )	★
$arccoth_F(x)$	(acoth $x$ )	†
$arcsech_F(x)$	(asech $x$ )	†
$arccsch_F(x)$	(acsch $x$ )	†
$axis\_rad_F(x)$	(axis_rad $x$ )	†
$rad_F(x)$	(radians $x$ )	†
$sin_F(x)$	(sin $x$ )	★
$cos_F(x)$	(cos $x$ )	★
$tan_F(x)$	(tan $x$ )	★
$cot_F(x)$	(cot $x$ )	†
$sec_F(x)$	(sec $x$ )	†
$csc_F(x)$	(csc $x$ )	†
$arcsin_F(x)$	(asin $x$ )	★
$arccos_F(x)$	(acos $x$ )	★
$arctan_F(x)$	(atan $x$ )	★
$arccot_F(x)$	(acot $x$ )	†
$arcctg_F(x)$	(actg $x$ )	†
$arcsec_F(x)$	(asec $x$ )	†
$arccsc_F(x)$	(acsc $x$ )	†
$arc_F(x, y)$	(atan2 $y x$ )	★
$axis\_cycle_F(u, x)$	(axis_cycle $u x$ )	†
$cycle_F(u, x)$	(cycle $u x$ )	†
$sinu_F(u, x)$	(sinU $u x$ )	†

$\text{cosu}_F(u, x)$	<code>(cosU u x)</code>	†
$\text{tanu}_F(u, x)$	<code>(tanU u x)</code>	†
$\text{cotu}_F(u, x)$	<code>(cotU u x)</code>	†
$\text{secu}_F(u, x)$	<code>(secU u x)</code>	†
$\text{cscu}_F(u, x)$	<code>(cscU u x)</code>	†
$\text{arcsinu}_F(u, x)$	<code>(asinU u x)</code>	†
$\text{arccosu}_F(u, x)$	<code>(acosU u x)</code>	†
$\text{arctanu}_F(u, x)$	<code>(atanU u x)</code>	†
$\text{arccotu}_F(u, x)$	<code>(acotU u x)</code>	†
$\text{arctgu}_F(u, x)$	<code>(actgU u x)</code>	†
$\text{arcsecu}_F(u, x)$	<code>(asecU u x)</code>	†
$\text{arccscu}_F(u, x)$	<code>(acscU u x)</code>	†
$\text{arcu}_F(u, x, y)$	<code>(atan2U u y x)</code>	†
$\text{rad\_to\_cycle}_F(x, u)$	<code>(rad_to_cycle x u)</code>	†
$\text{cycle\_to\_rad}_F(u, x)$	<code>(cycle_to_rad u x)</code>	†
$\text{cycle\_to\_cycle}_F(u, x, v)$	<code>(cycle_to_cycle u x v)</code>	†

where  $b$ ,  $x$ ,  $y$ ,  $u$ , and  $v$  are expressions of type *FLT*, and  $z$  is an expressions of type *INT*.

Arithmetic value conversions in ISLisp are can be explicit or implicit. The rules for when implicit conversions are done is implementation defined.

$\text{convert}_{I \rightarrow I''}(x)$	<code>(format g "~B" x)</code>	★(binary)
$\text{convert}_{I \rightarrow I''}(x)$	<code>(format g "~O" x)</code>	★(octal)
$\text{convert}_{I \rightarrow I''}(x)$	<code>(format g "~D" x)</code>	★(decimal)
$\text{convert}_{I \rightarrow I''}(x)$	<code>(format g "~X" x)</code>	★(hexadecimal)
$\text{convert}_{I \rightarrow I''}(x)$	<code>(format g "~rR" x)</code>	★(radix $r$ )
$\text{convert}_{I \rightarrow I''}(x)$	<code>(format-integer g x r)</code>	★(radix $r$ )
$\text{rounding}_{F \rightarrow I}(y)$	<code>(round y)</code>	★
$\text{floor}_{F \rightarrow I}(y)$	<code>(floor y)</code>	★
$\text{ceiling}_{F \rightarrow I}(y)$	<code>(ceiling y)</code>	★
$\text{convert}_{I \rightarrow F}(x)$	<code>(float x kind)</code>	★
$\text{convert}_{F \rightarrow F'}(y)$	<code>(float y kind)</code>	★
$\text{convert}_{F \rightarrow F''}(y)$	<code>(format g "~G" y)</code>	★
$\text{convert}_{F \rightarrow F''}(y)$	<code>(format-float g y)</code>	★

where  $x$  is an expression of type *INT*,  $y$  is an expression of type *FLT*, and  $z$  is an expression of type *FXD*, where *FXD* is a fixed point type. Conversion from string to numeric value is in ISLisp done via a general read procedure, which reads ISLisp ‘S-expressions’.

ISLisp provides non-negative numerals for its integer and floating point types in base is 10.

ISLisp does not specify numerals for infinities and NaNs. Suggestion:

$+\infty$	<code>infinity</code>	†
<code>qNaN</code>	<code>nan</code>	†
<code>sNaN</code>	<code>signan</code>	†

as well as string formats for reading and writing these values as character strings.

ISLisp has a notion of ‘error’ that implies a catchable, possibly returnable, change of control flow. ISLisp uses its exception mechanism as its default means of notification. ISLisp ignores **underflow** notifications. On **underflow** the continuation value (specified in LIA-2) is used directly without recording the **underflow** itself. ISLisp uses the error **domain-error** for **invalid** and some **infinitary** notifications, the error **arithmetic-error** for **overflow** notifications, and the error **division-by-zero** for other **infinitary** notifications.

An implementation that wishes to follow LIA-2 should provide recording of indicators as an alternative means of handling numeric notifications. Recording of indicators is the LIA-2 preferred means of handling numeric notifications.

## C.9 Modula-2

The programming language Modula-2 is defined by ISO/IEC 10514-1:1996, *Information technology – Programming languages – Part 1: Modula-2, Base Language* [21]. An implementation should follow all the requirements of LIA-3 unless otherwise specified by this language binding.

The operations or parameters marked “†” are not part of the language and should be provided by an implementation that wishes to conform to the LIA-3 for that operation. For each of the marked items a suggested identifier is provided.

The Modula-2 datatype **Boolean** corresponds to the LIA datatype **Boolean**.

The additional integer operations are listed below, along with the syntax used to invoke them:

$add_{c(I)}(x, y)$	<b>add</b> ( $x, y$ )	†
--------------------	-----------------------	---

where  $x$  and  $y$  are expressions of type *INT* and where  $xs$  is an expression of type **array** [] of *INT*.

The additional non-transcendental floating point operations are listed below, along with the syntax used to invoke them:

$add_{c(F)}(x, y)$	<b>add</b> ( $x, y$ )	†
$sub_{c(F)}(x, y)$	<b>sub</b> ( $x, y$ )	†

where  $x$ ,  $y$  and  $z$  are expressions of type *FLT*, and where  $xs$  is an expression of type **array** [] of *FLT*.

The LIA-3 parameters for operations approximating real valued transcendental functions can be accessed by the following syntax:

$max\_err\_mul_{c(F)}$	<b>err_mul</b> ( $x$ )	†
$max\_err\_div_{c(F)}$	<b>err_div</b> ( $x$ )	†
$max\_err\_exp_{c(F)}$	<b>err_exp</b> ( $x$ )	†
$max\_err\_power_{c(F)}$	<b>err_power</b> ( $x$ )	†
$max\_err\_sin_{c(F)}$	<b>err_sin</b> ( $x$ )	†
$max\_err\_tan_{c(F)}$	<b>err_tan</b> ( $x$ )	†

where  $x$  and  $u$  are expressions of type *FLT*. Several of the parameter functions are constant for each type (and library), the argument is then used only to differentiate among the floating point types.

The LIA-2 elementary floating point operations are listed below, along with the syntax used to invoke them:

$mul_{c(F)}(x, y)$	<code>mul(x, y)</code>	†
$div_{c(F)}(x, y)$	<code>div(x, y)</code>	†
$sqrt_{c(F)}(x)$	<code>sqrt(x)</code>	*
$exp_{c(F)}(x)$	<code>exp(x)</code>	*
$power_{c(F)}(b, y)$	<code>power(b, y)</code>	*
$ln_{c(F)}(x)$	<code>ln(x)</code>	*
$logbase_{c(F)}(b, x)$	<code>log(x, b)</code>	†
$radh_{c(F)}(x)$	<code>hypradian(x)</code>	†
$sinh_{c(F)}(x)$	<code>sinh(x)</code>	†
$cosh_{c(F)}(x)$	<code>cosh(x)</code>	†
$tanh_{c(F)}(x)$	<code>tanh(x)</code>	†
$coth_{c(F)}(x)$	<code>coth(x)</code>	†
$sech_{c(F)}(x)$	<code>sech(x)</code>	†
$csch_{c(F)}(x)$	<code>csch(x)</code>	†
$arcsinh_{c(F)}(x)$	<code>arcsinh(x)</code>	†
$arccosh_{c(F)}(x)$	<code>arccosh(x)</code>	†
$arctanh_{c(F)}(x)$	<code>arctanh(x)</code>	†
$arccoth_{c(F)}(x)$	<code>arccoth(x)</code>	†
$arcsech_{c(F)}(x)$	<code>arcsech(x)</code>	†
$arccsch_{c(F)}(x)$	<code>arccsch(x)</code>	†
$rad_{c(F)}(x)$	<code>radian(x)</code>	†
$sin_{c(F)}(x)$	<code>sin(x)</code>	*
$cos_{c(F)}(x)$	<code>cos(x)</code>	*
$tan_{c(F)}(x)$	<code>tan(x)</code>	*
$cot_{c(F)}(x)$	<code>cot(x)</code>	†
$sec_{c(F)}(x)$	<code>sec(x)</code>	†
$csc_{c(F)}(x)$	<code>csc(x)</code>	†
$arcsin_{c(F)}(x)$	<code>arcsin(x)</code>	*
$arccos_{c(F)}(x)$	<code>arccos(x)</code>	*
$arctan_{c(F)}(x)$	<code>arctan(x)</code>	*
$arccot_{c(F)}(x)$	<code>arccot(x)</code>	†
$arcctg_{c(F)}(x)$	<code>arcctg(x)</code>	†
$arcsec_{c(F)}(x)$	<code>arcsec(x)</code>	†
$arccsc_{c(F)}(x)$	<code>arccsc(x)</code>	†

where  $b$ ,  $x$ ,  $y$ ,  $u$ , and  $v$  are expressions of type *CFLT*.

Arithmetic value conversions in Modula-2 are can be explicit or implicit. The rules for when implicit conversions are applied is not repeated here. The explicit arithmetic value conversions

are usually expressed as ‘casts’, except when converting to/from strings.

where  $x$  is an expression of type *INT*,  $y$  is an expression of type *FLT*, and  $z$  is an expression of type *FXD*, where *FXD* is a fixed point type. *INT2* is the integer datatype that corresponds to  $I'$ . A ? above indicates that the parameter is optional.  $e$  is greater than 0.

## C.10 PL/I

The programming language PL/I is defined by ANSI X3.53-1976 (R1998), *Programming languages – PL/I* [39], and endorsed by ISO 6160:1979, *Programming languages – PL/I* [25]. The programming language General Purpose PL/I is defined by ISO/IEC 6522:1992, *Information technology – Programming languages – PL/I general-purpose subset* [26], also: ANSI X3.74-1987 (R1998).

An implementation should follow all the requirements of LIA-3 unless otherwise specified by this language binding.

The operations or parameters marked “†” are not part of the language and should be provided by an implementation that wishes to conform to the LIA-3 for that operation. For each of the marked items a suggested identifier is provided.

The LIA datatype **Boolean** is implemented in the PL/I datatype **BIT(1)** (1 = **true** and 0 = **false**).

An implementation of PL/I provides at least one integer data type, and at least one floating point data type. The attribute **FIXED**( $n,0$ ) selects a signed integer datatype with at least  $n$  (decimal or binary) digits of storage. The attribute **FLOAT**( $k$ ) selects a floating point datatype with at least  $n$  (decimal or binary) digits of precision.

The LIA-3 integer operations are listed below, along with the syntax used to invoke them:

$add_{c(I)}(x, y)$	$x + y$	†
$sub_{c(I)}(x, y)$	$x - y$	†
$mul_{c(I)}(x, y)$	$x * y$	†

where  $x$  and  $y$  are expressions of type *CINT*.

The LIA-3 non-transcendental floating point operations are listed below, along with the syntax used to invoke them:

$add_{c(F)}(x, y)$	$x + y$	*
$sub_{c(F)}(x, y)$	$x - y$	*

where  $x$ ,  $y$  and  $z$  are expressions of type *FLT*, and where  $xs$  is an expression of type **array** of *FLT*.

The parameters for operations approximating real valued transcendental functions can be accessed by the following syntax:

$max\_err\_mul_{c(F)}$	<b>err_mul</b> ( $x$ )	†
$max\_err\_div_{c(F)}$	<b>err_div</b> ( $x$ )	†
$max\_err\_exp_{c(F)}$	<b>err_exp</b> ( $x$ )	†
$max\_err\_power_{c(F)}$	<b>err_power</b> ( $x$ )	†

$max\_err\_sin_{c(F)}$	<code>err_sin(x)</code>	†
$max\_err\_tan_{c(F)}$	<code>err_tan(x)</code>	†

where  $x$  and  $u$  are expressions of type *CFLT*. Several of the parameter functions are constant for each type (and library), the argument is then used only to differentiate among the floating point types.

The LIA-3 elementary floating point operations are listed below, along with the syntax used to invoke them:

$mul_{c(F)}(x, y)$	<code>x * y</code>	*
$div_{c(F)}(x, y)$	<code>x / y</code>	*
$sqrt_{c(F)}(x)$	<code>sqrt(x)</code>	*
$exp_{c(F)}(x)$	<code>exp(x)</code>	*
$power_{c(F)}(b, y)$	<code>power(b, y)</code>	†
$ln_{c(F)}(x)$	<code>log(x)</code>	*
$logbase_{c(F)}(b, x)$	<code>log(b, x)</code>	†
$radh_{c(F)}(x)$	<code>radh(x)</code>	†
$sinh_{c(F)}(x)$	<code>sinh(x)</code>	*
$cosh_{c(F)}(x)$	<code>cosh(x)</code>	*
$tanh_{c(F)}(x)$	<code>tanh(x)</code>	*
$coth_{c(F)}(x)$	<code>coth(x)</code>	†
$sech_{c(F)}(x)$	<code>sech(x)</code>	†
$csch_{c(F)}(x)$	<code>csch(x)</code>	†
$arcsinh_{c(F)}(x)$	<code>arcsinh(x)</code>	*
$arccosh_{c(F)}(x)$	<code>arccosh(x)</code>	*
$arctanh_{c(F)}(x)$	<code>arctanh(x)</code>	*
$arcoth_{c(F)}(x)$	<code>arcoth(x)</code>	†
$arcsech_{c(F)}(x)$	<code>arcsech(x)</code>	†
$arccsch_{c(F)}(x)$	<code>arccsch(x)</code>	†
$rad_{c(F)}(x)$	<code>rad(x)</code>	†
$sin_{c(F)}(x)$	<code>sin(x)</code>	*
$cos_{c(F)}(x)$	<code>cos(x)</code>	*
$tan_{c(F)}(x)$	<code>tan(x)</code>	*
$cot_{c(F)}(x)$	<code>cot(x)</code>	*
$sec_{c(F)}(x)$	<code>sec(x)</code>	†
$csc_{c(F)}(x)$	<code>csc(x)</code>	†
$arcsin_{c(F)}(x)$	<code>arcsin(x)</code>	*
$arccos_{c(F)}(x)$	<code>arccos(x)</code>	*
$arctan_{c(F)}(x)$	<code>arctan(x)</code>	*
$arccot_{c(F)}(x)$	<code>arccot(x)</code>	†
$arcctg_{c(F)}(x)$	<code>arcctg(x)</code>	†
$arcsec_{c(F)}(x)$	<code>arcsec(x)</code>	†
$arccsc_{c(F)}(x)$	<code>arccsc(x)</code>	†

where  $b$ ,  $x$ ,  $y$ ,  $u$ , and  $v$  are expressions of type *CFLT*.

Arithmetic value conversions in PL/I are can be explicit or implicit. The rules for when implicit conversions are applied is not repeated here. The explicit arithmetic value conversions are usually expressed as ‘casts’, except when converting to/from strings.

where  $x$  is an expression of type  $INT$ ,  $y$  is an expression of type  $FLT$ , and  $z$  is an expression of type  $FXD$ , where  $FXD$  is a fixed point type.  $INT2$  is the integer datatype that corresponds to  $I'$ . A ? above indicates that the parameter is optional.  $a$  is greater than 0.

## C.11 SML

The programming language SML is defined by *The Definition of Standard ML (Revised)* [58].

An implementation should follow all the requirements of LIA-2 unless otherwise specified by this language binding.

The operations or parameters marked “†” are not part of the language and should be provided by an implementation that wishes to conform to the LIA-3 for that operation. For each of the marked items a suggested identifier is provided.

The SML datatype **Boolean** corresponds to the LIA datatype **Boolean**.

Every implementation of SML has at least one integer datatype, **int**, and at least one floating point datatype, **real**. The notation  $INT$  is used to stand for the name of one of the integer datatypes, and  $FLT$  is used to stand for the name of one of the floating point datatypes in what follows.

The LIA-2 integer operations are listed below, along with the syntax used to invoke them:

$add_{c(I)}(x, y)$	$x + y$ or <b>op +</b> $(x, y)$	★
$sub_{c(I)}(x, y)$	$x - y$ or <b>op -</b> $(x, y)$	★
$mul_{c(I)}(x, y)$	$x * y$ or <b>op *</b> $(x, y)$	★

where  $x$  and  $y$  are expressions of type  $CINT$ .

The additional non-transcendental floating point operations are listed below, along with the syntax used to invoke them:

$add_{c(F)}(x, y)$	$x + y$ or <b>op +</b> $(x, y)$	★
$sub_{c(F)}(x, y)$	$x - y$ or <b>op -</b> $(x, y)$	★

where  $x$ ,  $y$  and  $z$  are expressions of type  $CFLT$ .

The binding for the floor, round, and ceiling operations here take advantage of the unlimited **Integer** type in SML, and that **Integer** is the default integer type.

The parameters for operations approximating real valued transcendental functions can be accessed by the following syntax:

$max\_err\_mul_F$	<b>err_mul</b> $x$	†
$max\_err\_div_F$	<b>err_div</b> $x$	†
$max\_err\_exp_F$	<b>err_exp</b> $x$	†
$max\_err\_power_F$	<b>err_power</b> $x$	†
$max\_err\_sin_F$	<b>err_sin</b> $x$	†
$max\_err\_tan_F$	<b>err_tan</b> $x$	†

where  $x$  and  $u$  are expressions of type *FLT*. Several of the parameter functions are constant for each type (and library), the argument is then used only to differentiate among the floating point types.

The LIA-3 elementary floating point operations are listed below, along with the syntax used to invoke them:

$mul_{c(F)}(x, y)$	$x * y$ or $op * (x, y)$	*
$div_{c(F)}(x, y)$	$x / y$ or $op / (x, y)$	*
$sqrt_{c(F)}(x)$	$sqrt\ x$	*
$exp_{c(F)}(x)$	$exp\ x$	*
$power_{c(F)}(b, y)$	$b ** y$	†
$pow_{c(F)}(b, y)$	$b\ pow\ y$ or $op\ pow\ (x, y)$	* Not LIA-2! (See C.)
$ln_{c(F)}(x)$	$ln\ x$	*
$logbase_{c(F)}(b, x)$	$log\_base\ (b, x)$	†
$radh_{c(F)}(x)$	$hypradians\ x$	†
$sinh_{c(F)}(x)$	$sinh\ x$	*
$cosh_{c(F)}(x)$	$cosh\ x$	*
$tanh_{c(F)}(x)$	$tanh\ x$	*
$coth_{c(F)}(x)$	$coth\ x$	†
$sech_{c(F)}(x)$	$sech\ x$	†
$csch_{c(F)}(x)$	$csch\ x$	†
$arcsinh_{c(F)}(x)$	$arcsinh\ x$	†
$arccosh_{c(F)}(x)$	$arccosh\ x$	†
$arctanh_{c(F)}(x)$	$arctanh\ x$	†
$arccoth_{c(F)}(x)$	$arccoth\ x$	†
$arcsech_{c(F)}(x)$	$arcsech\ x$	†
$arccsch_{c(F)}(x)$	$arccsch\ x$	†
$rad_{c(F)}(x)$	$radians\ x$	†
$sin_{c(F)}(x)$	$sin\ x$	*
$cos_{c(F)}(x)$	$cos\ x$	*
$tan_{c(F)}(x)$	$tan\ x$	*
$cot_{c(F)}(x)$	$cot\ x$	†
$sec_{c(F)}(x)$	$sec\ x$	†
$csc_{c(F)}(x)$	$csc\ x$	†
$arcsin_{c(F)}(x)$	$arcsin\ x$	*
$arccos_{c(F)}(x)$	$arccos\ x$	*
$arctan_{c(F)}(x)$	$arctan\ x$	*
$arccot_{c(F)}(x)$	$arccot\ x$	†
$arcctg_{c(F)}(x)$	$arcctg\ x$	†
$arcsec_{c(F)}(x)$	$arcsec\ x$	†
$arccsc_{c(F)}(x)$	$arccsc\ x$	†



where  $b$ ,  $x$ ,  $y$ ,  $u$ , and  $v$  are expressions of type  $FLT$ , and  $z$  is an expressions of type  $INT$

Type conversions in SML are always explicit.

where  $x$  is an expression of type  $INT$ ,  $y$  is an expression of type  $FLT$ , and  $z$  is an expression of type  $FXD$ , where  $FXD$  is a fixed point type.  $INT2$  is the integer datatype that corresponds to  $I'$ .



## Annex D (informative)

### Bibliography

This annex gives references to publications relevant to LIA-3.

#### International standards documents

- [1] ISO/IEC JTC1 Directives, Part 3: *Drafting and presentation of International Standards*, 1989.
- [2] ISO 6093:1985, *Information processing – Representation of numerical values in character strings for information interchange*.
- [3] ISO/IEC TR 10176:1998, *Information technology – Guidelines for the preparation of programming language standards*.
- [4] ISO/IEC TR 10182:1993, *Information technology – Programming languages, their environments and system software interfaces – Guidelines for language bindings*.
- [5] ISO/IEC 13886:1996, *Information technology – Language-Independent Procedure Calling, (LIPC)*.
- [6] ISO/IEC 11404:1996, *Information technology – Programming languages, their environments and system software interfaces – Language-independent datatypes, (LID)*.
- [7] ISO/IEC 8652:1995, *Information technology – Programming languages – Ada*.
- [8] ISO/IEC 13813:1998, *Information technology – Programming languages – Generic packages of real and complex type declarations and basic operations for Ada (including vector and matrix types)*.
- [9] ISO/IEC 13814:1998, *Information technology – Programming languages – Generic package of complex elementary functions for Ada*.
- [10] ISO 8485:1989, *Programming languages – APL*.
- [11] ISO/IEC DIS 13751, *Information technology – Programming languages, their environments and system software interfaces – Programming language APL, extended, 1999*.
- [12] ISO/IEC 10279:1991, *Information technology – Programming languages – Full BASIC*. (Essentially an endorsement of ANSI X3.113-1987 (R1998) [36].)
- [13] ISO/IEC 9899:1999, *Programming languages – C*.
- [14] ISO/IEC 14882:1998, *Programming languages – C++*.
- [15] ISO 1989:1985, *Programming languages – COBOL*. (Endorsement of ANSI X3.23-1985 (R1991) [37].) Currently under revision (1998).
- [16] ISO/IEC 16262:1998, *Information technology – ECMAScript language specification*.
- [17] ISO/IEC 15145:1997, *Information technology – Programming languages – FORTH*. (Also: ANSI X3.215-1994.)

- [18] ISO/IEC 1539-1:1997, *Information technology – Programming languages – Fortran - Part 1: Base language.*
- [19] ISO/IEC TR 15580:1998, *Information technology – Programming languages – Fortran – Floating-point exception handling.*
- [20] ISO/IEC 13816:1997, *Information technology – Programming languages, their environments and system software interfaces – Programming language ISLISP.*
- [21] ISO/IEC 10514-1:1996, *Information technology – Programming languages – Part 1: Modula-2, Base Language.*
- [22] ISO/IEC 10514-2:1998, *Information technology – Programming languages – Part 2: Generics Modula-2.*
- [23] ISO 7185:1990, *Information technology – Programming languages – Pascal.*
- [24] ISO/IEC 10206:1991, *Information technology – Programming languages – Extended Pascal.*
- [25] ISO 6160:1979, *Programming languages – PL/I.* (Endorsement of ANSI X3.53-1976 (R1998) [39].)
- [26] ISO/IEC 6522:1992, *Information technology – Programming languages – PL/I general-purpose subset.* (Also: ANSI X3.74-1987 (R1998).)
- [27] ISO/IEC 13211-1:1995, *Information technology – Programming languages – Prolog – Part 1: General core.*
- [28] ISO/IEC 9075:1992, *Information technology – Database languages – SQL.*
- [29] ISO/IEC 8824-1:1995, *Information technology – Abstract Syntax Notation One (ASN.1) – Part 1: Specification of basic notation.*
- [30] ISO/IEC 9001:1994, *Quality systems – Model for quality assurance in design, development, production, installation and servicing.*
- [31] ISO/IEC 9126:1991, *Information technology – Software product evaluation – Quality characteristics and guidelines for their use.*
- [32] ISO/IEC 12119:1994, *Information technology – Software packages – Quality requirements and testing.*
- [33] ISO/IEC 14598-1:1999, *Information technology – Software product evaluation – Part 1: General overview.*

### National standards documents

- [34] ANSI/IEEE Standard 754-1985, *IEEE Standard for Binary Floating-Point Arithmetic*.
- [35] ANSI/IEEE Standard 854-1987, *IEEE Standard for Radix-Independent Floating-Point Arithmetic*.
- [36] ANSI X3.113-1987 (R1998), *Information technology – Programming languages – Full BASIC*.
- [37] ANSI X3.23-1985 (R1991), *Programming languages – COBOL*.
- [38] ANSI X3.226-1994, *Information Technology – Programming Language – Common Lisp*.
- [39] ANSI X3.53-1976 (R1998), *Programming languages – PL/I*.
- [40] ANSI/IEEE 1178-1990, *IEEE Standard for the Scheme Programming Language*.
- [41] ANSI/NCITS 319-1998, *Information Technology – Programming Languages – Smalltalk*.

### Books, articles, and other documents

- [42] J. S. Squire (ed), *Ada Letters*, vol. XI, No. 7, ACM Press (1991).
- [43] M. Abramowitz and I. Stegun (eds), *Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables*, Tenth Printing, 1972, Superintendent of Documents, U.S. Government Printing Office, Washington, D.C. 20402.
- [44] J. Du Croz and M. Pont, *The Development of a Floating-Point Validation Package*, *NAG Newsletter*, No. 3, 1984.
- [45] J. W. Demmel and X. Li, *Faster Numerical Algorithms via Exception Handling*, 11th International Symposium on Computer Arithmetic, Winsor, Ontario, June 29 - July 2, 1993.
- [46] D. Goldberg, *What Every Computer Scientist Should Know about Floating-Point Arithmetic*. ACM Computing Surveys, Vol. 23, No. 1, March 1991.
- [47] J. R. Hauser, *Handling Floating-Point Exceptions in Numeric Programs*. ACM Transactions on Programming Languages and Systems, Vol. 18, No. 2, March 1986, Pages 139-174.
- [48] W. Kahan, *Branch Cuts for Complex Elementary Functions, or Much Ado about Nothing's Sign Bit*, Chapter 7 in *The State of the Art in Numerical Analysis* ed. by M Powell and A Iserles (1987) Oxford.
- [49] W. Kahan, *Lecture Notes on the Status of IEEE Standard 754 for Binary Floating-Point Arithmetic*, Panel Discussion of *Floating-Point Past, Present and Future*, May 23, 1995, in a series of San Francisco Bay Area Computer Historical Perspectives, sponsored by SUN Microsystems Inc.
- [50] U. Kulisch and W. L. Miranker, *Computer Arithmetic in Theory and Practice*, Academic Press, 1981.
- [51] U. Kulisch and W. L. Miranker (eds), *A New Approach to Scientific Computation*, Academic Press, 1983.
- [52] D. C. Sorenson and P. T. P. Tang, *On the Orthogonality of Eigenvectors Computed by Divide-and-Conquer Techniques*, *SIAM Journal of Numerical Analysis*, Vol. 28, No. 6, p. 1760, algorithm 5.3.

- [53] *Floating-Point C Extensions* in Technical Report Numerical C Extensions Committee X3J11, April 1995, SC22/WG14 N403, X3J11/95-004.
- [54] David M. Gay, *Correctly Rounded Binary-Decimal and Decimal-Binary Conversions*, AT&T Bell Laboratories, Numerical Analysis Manuscript 90-10, November 1990.
- [55] James Gosling, Bill Joy, Guy Steele, *The Java Language Specification*.
- [56] Simon Peyton Jones et al., *Report on the programming language Haskell 98*, February 1999.
- [57] Simon Peyton Jones et al., *Standard libraries for the Haskell 98 programming language*, February 1999.
- [58] Robin Milner, Mads Tofte, Robert Harper, and David MacQueen, *The Definition of Standard ML (Revised)*, The MIT Press, 1997, ISBN: 0-262-63181-4.