INTERNATIONAL STANDARD

ISO/IEC 10967-2

First edition 2001-08-15

Information technology — Language independent arithmetic —

Part 2: Elementary numerical functions

Technologies de l'information — Arithmétique de langage indépendant —
Partie 2: Fonctions numériques élémentaires

L'ECHOCHIN.COM.



PDF disclaimer

This PDF file may contain embedded typefaces. In accordance with Adobe's licensing policy, this file may be printed or viewed but shall not be edited unless the typefaces which are embedded are licensed to and installed on the computer performing the editing. In downloading this file, parties accept therein the responsibility of not infringing Adobe's licensing policy. The ISO Central Secretariat accepts no liability in this area.

Adobe is a trademark of Adobe Systems Incorporated.

Details of the software products used to create this PDF file can be found in the General Info relative to the file; the PDF-creation parameters were optimized for printing. Every care has been taken to ensure that the file is suitable for use by ISO member bodies. In the unlikely event that a problem relating to it is found, please inform the Central Secretariat at the address given below.

ECNORM. Cick to View the full PDF of Econetic Appet 2:2001

© ISO/IEC 2001

All rights reserved. Unless otherwise specified, no part of this publication may be reproduced or utilized in any form or by any means, electronic or mechanical, including photocopying and microfilm, without permission in writing from either ISO at the address below or ISO's member body in the country of the requester.

ISO copyright office
Case postale 56 • CH-1211 Geneva 20
Tel. + 41 22 749 01 11
Fax + 41 22 749 09 47
E-mail copyright@iso.ch
Web www.iso.ch

Printed in Switzerland

Contents

			n	viii ix					
1	Sco. 1.1 1.2	Inclus	ions	1 1 2					
2	Cor	nformi	ty	2					
3	3 Normative references								
4	Symbols and definitions 4 4.1 Symbols								
	4.1	Symbo	ols	4					
		4.1.1	Sets and intervals	4					
		4.1.2	Operators and relations	4					
		4.1.3	Mathematical functions	5					
		4.1.4	Exceptional values	5					
		4.1.5	Datatypes	6					
	4.2	Defini	Datatypes	7					
5	\mathbf{Spe}	cificat	ions for integer and floating point operations	10					
	5.1	Basic	integer operations	10					
		5.1.1	The integer $result$ and $wrap$ helper functions	10					
		5.1.2	Integer maximum and minimum Integer diminish	11					
		5.1.3	Integer diminish	11					
		5.1.4	Integer power and arithmetic shift	12					
		5.1.5	Integer square root	12					
		5.1.6	Divisibility tests	12					
		5.1.7	Integer division (with floor, round, or ceiling) and remainder	13					
		5.1.8	Greatest common divisor and least common positive multiple	13					
		5.1.9	Support operations for extended integer range	14					
	5.2	Basic	floating point operations	15					
		5.2.1	The rounding and floating point $result$ helper functions	15					
		5.2.2	Floating point maximum and minimum	17					
		5.2.3	Floating point diminish	18					
		5.2.4	Proor, round, and ceiling	19					
		5.2.5	Remainder after division with round to integer	20					
		5.2.6	Square root and reciprocal square root	20					
		5.2.7	Multiplication to higher precision floating point datatype	20					
		5.2.8	Support operations for extended floating point precision	21					
	5.3	Eleme	entary transcendental floating point operations	22					
		5.3.1	Maximum error requirements	22					
		5.3.2	Sign requirements	23					
		5.3.3	Monotonicity requirements	23					
		5.3.4	The $result^*$ helper function	23					
		5.3.5	Hypotenuse	24					
		5.3.6	Operations for exponentiations and logarithms	24					

5.3.6.1	Integer power of argument base	24
5.3.6.2	Natural exponentiation	25
5.3.6.3	Natural exponentiation, minus one	26
5.3.6.4	Exponentiation of 2	27
5.3.6.5	Exponentiation of 10	27
5.3.6.6		28
5.3.6.7	•	29
5.3.6.8		29
5.3.6.9	Natural logarithm of one plus the argument	30
5.3.6.10		30
5.3.6.11		31
5.3.6.12	Argument base logarithm	31
5.3.6.13	Argument base logarithm of one plus each argument	32
Introduc	tion to operations for trigonometric elementary functions of the second	32
	ons for radian trigonometric elementary functions	33
5.3.8.1	Radian angle normalisation	34
5.3.8.2	Radian sine	35
5.3.8.3	Padian agging	35
5.3.8.4	Radian tangent	36
5.3.8.5	Radian cotangent	36
5.3.8.6	Radian secant	37
5.3.8.7	Radian cosecant	37
5.3.8.8	Radian cosine with sine	38
5.3.8.9	Radian arc sine	38
5.3.8.10	Radian arc cosine	38
5.3.8.11	Radian arc tangent	39
5.3.8.12	Radian arc cotangent	40
5.3.8.13	Radian arc secant	41
		41
		42
		43
_		43
		44
		45
		45
/ 4	- 9	46
~ // .		47
		47
		48
		48
		48
		49
		50
5.3.9.13		51
5.3.9.14		51
5.3.9.15	9	52
		53
-	9	53
	5.3.6.2 5.3.6.3 5.3.6.4 5.3.6.5 5.3.6.6 5.3.6.7 5.3.6.8 5.3.6.9 5.3.6.10 5.3.6.11 5.3.6.12 5.3.6.13 Introduce Operation 5.3.8.1 5.3.8.2 5.3.8.3 5.3.8.4 5.3.8.5 5.3.8.6 5.3.8.7 5.3.8.8 5.3.8.9 5.3.8.10 5.3.8.11 5.3.8.12 5.3.8.13 5.3.8.14 5.3.8.15 Operation 5.3.8.14 5.3.8.15 5.3.8.15 5.3.8.14 5.3.8.15 5.3.8.14 5.3.8.15 5.3.8.14 5.3.8.15 5.3.8.14 5.3.8.15 5.3.8.16 5.3.9.1 5.3.9.1 5.3.9.1 5.3.9.1 5.3.9.1 5.3.9.10 5.3.9.11 5.3.9.11 5.3.9.12 5.3.9.13 5.3.9.14 5.3.9.15 Operation 5.3.9.11 5.3.9.15 Operation 5.3.9.11	5.3.6.2 Natural exponentiation 5.3.6.3 Natural exponentiation, minus one 5.3.6.4 Exponentiation of 2 5.3.6.5 Exponentiation of 10 5.3.6.6 Exponentiation of argument base 5.3.6.7 Exponentiation of one plus the argument base, minus one 5.3.6.8 Natural logarithm 5.3.6.10 2-logarithm 5.3.6.11 10-logarithm 5.3.6.13 Argument base logarithm 5.3.6.13 Argument base logarithm of one plus each argument Introduction to operations for trigonometric elementary functions Operations for radian trigonometric elementary functions 5.3.8.1 Radian angle normalisation 5.3.8.2 Radian sine 5.3.8.3 Radian cosine 5.3.8.4 Radian tangent 5.3.8.5 Radian cosenat 5.3.8.6 Radian cosecant 5.3.8.7 Radian cosecant 5.3.8.8 Radian arc sine 5.3.8.9 Radian arc sine 5.3.8.10 Radian arc sceant 5.3.8.11 Radian arc sceant 5.3.8.12 Radian arc sceant 5.3.8.13 Radian arc sceant 5.3.8.14 Radian arc sceant 5.3.8.15 Radian arc sceant 5.3.8.17 Radian arc sceant 5.3.8.18 Radian arc sceant 5.3.8.19 Radian arc sceant 5.3.8.10 Radian arc sceant 5.3.8.11 Radian arc sceant 5.3.8.12 Radian arc sceant 5.3.8.13 Radian arc sceant 5.3.8.14 Radian arc sceant 5.3.8.15 Radian arc sceant 5.3.8.16 Radian arc sceant 5.3.8.17 Radian arc sceant 5.3.8.18 Radian arc sceant 5.3.8.19 Argument angular-unit angle normalisation 5.3.9.1 Argument angular-unit sine 5.3.9.2 Argument angular-unit tosine 5.3.9.3 Argument angular-unit tosine 5.3.9.4 Argument angular-unit cosine 5.3.9.5 Argument angular-unit cosine with sine 5.3.9.9 Argument angular-unit cosine with sine 5.3.9.10 Argument angular-unit cosine with sine 5.3.9.10 Argument angular-unit cosine with sine 5.3.9.11 Argument angular-unit arc cosine 5.3.9.12 Argument angular-unit arc cosine 5.3.9.13 Argument angular-unit arc cosine 5.3.9.14 Argument angular-unit arc cosine 5.3.9.15 Argument angular-unit arc cosine 5.3.9.16 Argument angular-unit arc cosine 5.3.9.17 Argument angular-unit arc cosine 5.3.9.18 Argument angular-unit arc cosine 5.3.9.19 Argument angular-unit arc cosine 5.3.9.11 Argument angular-unit arc cosine 5.3.9.12 Argumen

			5.3.10.2 Converting argument angular-unit angle to radian angle	54
			5.3.10.3 Converting argument angular-unit angle to (another) argument	
			angular-unit angle	55
		5.3.11	Operations for hyperbolic elementary functions	56
			5.3.11.1 Hyperbolic sine	56
			5.3.11.2 Hyperbolic cosine	56
			5.3.11.3 Hyperbolic tangent	57
			5.3.11.4 Hyperbolic cotangent	58
			5.3.11.5 Hyperbolic secant	58
			5.3.11.6 Hyperbolic cosecant	59
			5.3.11.7 Inverse hyperbolic sine	59
			5.3.11.8 Inverse hyperbolic cosine	60
			5.3.11.9 Inverse hyperbolic tangent	60
			5.3.11.10 Inverse hyperbolic cotangent	60
			5.3.11.11 Inverse hyperbolic secant	61
			5.3.11.12 Inverse hyperbolic cosecant	61
	5.4	Opera	tions for conversion between numeric datatypes	62
		5.4.1	Integer to integer conversions	63
		5.4.2	Integer to integer conversions	63
		5.4.3	Integer to floating point conversions	64
		5.4.4	Floating point to floating point conversions	64
		5.4.5	Floating point to fixed point conversions.	65
		5.4.6	Fixed point to floating point conversions	66
	5.5		rals as operations in a programming language	67
	0.0	5.5.1	Numerals for integer datatypes	67
		5.5.2	Numerals for floating point datatypes	68
		0.0.2	realistate for housing point data.	00
6	Not	ificatio	on a la l	68
	6.1	Contin	nuation values	69
			on nuation values	
7	\mathbf{Rel}	ationsl	hip with language standards	69
0	ъ	,	alick and the second se	=0
8	Doo	cument	tation requirements	70
۸ -	nnov	A (no	ormative) Partial conformity	73
-1 .	A.1	`	num error relaxation	73
	A.2		accuracy requirements relaxation	74
	A.3		onships to other operations relaxation	74
	A.4		close-to-axis angular normalisation relaxation	$\frac{74}{74}$
	A.5		requirements relaxation	75
	Α.υ	atti	requirements relaxation	10
A	nnex	B (in	formative) Rationale	77
	B.1		,	77
	2.1	B.1.1	Inclusions	77
		B.1.2	Exclusions	78
	B.2		rmity	78
	-	B.2.1	Validation	79
	В.3		ative references	79
	B.4		ols and definitions	79
		~		

	B.4.1	Symbols		9
		B.4.1.1	Sets and intervals	9
		B.4.1.2	Operators and relations	80
		B.4.1.3	-	80
		B.4.1.4	Exceptional values	80
		B.4.1.5	Datatypes	31
	B.4.2	Definitio	ons of terms	31
B.5	Specifi	cations for	or the numerical functions	31
	-			32
		B.5.1.1		32
		B.5.1.2	Integer maximum and minimum	32
		B.5.1.3	Integer diminish	32
		B.5.1.4	Integer power and arithmetic shift	3
		B.5.1.5	Integer square root	3
		B.5.1.6		3
		B.5.1.7	Integer division (with floor, round, or ceiling) and remainder 8	3
		B.5.1.8		34
		B.5.1.9	Support operations for extended integer range	34
	B.5.2	Basic flo		34
		B.5.2.1		86
		B.5.2.2	Floating point maximum and minimum 8	86
		B.5.2.3	Floating point diminish	86
		B.5.2.4	, , , , , , , , , , , , , , , , , , , ,	86
		B.5.2.5	Remainder after division and round to integer	37
		B.5.2.6	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	37
		B.5.2.7		88
		B.5.2.8		88
	B.5.3			39
		B.5.3.1	1 · · ·	39
		B.5.3.2	g 1 VO	0
		B.5.3.3		0
		B.5.3.4	*	0
		B.5.3.5		1
		B.5.3.6	1	1
				93
		B.5.3.8	1 0)4
		B.5.3.9	1 0	96
			1	7
	~ C)			8
	B.5.4	-	<i>u</i> 1	8
	B.5.5		1 0 0 0	9
		B.5.5.1	9 01	9
D.C	NT	B.5.5.2		9
B.6	Notific			
D 7	B.6.1		ation values	
B.7		_	th language standards	
B.8	Docun	nentation	requirements	11

Annex C (informative) Example bindings for specific languages	103
C.1 Ada	104
C.2 BASIC	
C.3 C	
$C.4$ $C++\dots$	
C.5 Fortran	
C.6 Haskell	
C.7 Java	
C.8 Common Lisp	
C.9 ISLisp	
C.10 Modula-2	
C.11 Pascal and Extended Pascal	162
C.12 PL/I	167
0.13 SML	10
C.12 PL/I C.13 SML Annex D (informative) Bibliography Annex E (informative) Possible changes to part 1	173
Annex E (informative) Possible changes to part 1	177
and the second s	
No.	
a i CK	
· C.	
\mathcal{N} .	
COX CONTRACTOR OF THE CONTRACT	
D. W.	
10 ,	
ECNORM. COM.	

Foreword

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

International Standards are drafted in accordance with the rules given 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. 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 10967-2 was prepared by Joint Technical Committee ISO/IEC JTC 1, 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 integer and floating point arithmetic and complex elementary numerical functions

Additional parts will specify other arithmetic datatypes or arithmetic operations.

Annex A forms a normative part of this part of ISO/IEC 10967. Annexes B to E are for information only.

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.

Programmers writing programs that perform a significant amount of numeric processing have often not been certain how a program will perform when run under a given language processor. Programming language standards have traditionally been somewhat weak in the area of numeric processing, seldom providing an adequate specification of the properties of arithmetic datatypes, particularly floating point numbers. Often they do not even require much in the way of documentation of the actual arithmetic operations by a conforming language processor.

It is the intent of this part to help to redress these shortcomings, by setting out precise definitions of elementary numerical functions, and requirements for documentation.

It is not claimed that this part will ensure complete certainty of arithmetic behaviour in all circumstances; the complexity of numeric software and the difficulties of analysing and proving algorithms are too great for that to be attempted. Rather, this international Standard will provide a firmer basis than hitherto for attempting such analysis.

The aims for this part, part 2 of ISO/IEC 10967, are extensions of the aims for part 1: to ensure adequate accuracy for numerical computation, predictability, notification on the production of exceptional results, and compatibility with programming language standards.

The content

The content of this part is based on part 1, and extends part 1's specifications to specifications for operations approximating real elementary functions, operations often required (usually without a detailed specification) by the standards for programming languages widely used for scientific software. This part also provides specifications for conversions between the "internal" values of an arithmetic datatype, and a very close approximation in, e.g., the decimal radix. It does not cover the further transformation to decimal string format, which is usually provided by language standards. This part also includes specifications for a number of other functions deemed useful, even though they may not be stipulated by programming language standards.

The numerical functions covered by this part are computer approximations to mathematical functions of one or more real 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 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.

ECHORAN.COM. Chakto view the full Path of South Cha 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.

Х

Information technology — Language independent arithmetic —

Part 2: Elementary numerical functions

1 Scope

This part of ISO/IEC 10967 defines the properties of numerical approximations for many of the real 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 are included reference in this part.

This part provides specifications for numerical functions for which all operand values are of integer or 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 covers most numerical functions required by the ISO/IEC standards for Ada [11], Basic [16], C [17], C++ [18], Fortran [22], ISLisp [24], Pascal [27], and PL/I [29]. In particular, specifications are provided for:

- a) Some additional integer operations.
- b) Some additional non-transcendental floating point operations, including maximum and minimum operations.
- c) Exponentiations, logarithms, and hyperbolics.
- d) Trigonometrics, both in radians and for argument-given angular unit with degrees as a special case.

This part also provides specifications for:

1. Scope 1

- e) Conversions between integer and floating point datatypes (possibly with different radices) conforming to the requirements of part 1, and the conversion operations used, for example, in text input and output of integer and floating point numbers.
- f) The results produced by an included floating point operation when one or more argument values are IEC 60559 special values.
- g) Program-visible parameters that characterise certain aspects of the operations.

1.2 Exclusions

This part provides no specifications for

- a) Numerical functions whose operands are of more than one datatype (with one exception). This part neither requires nor excludes the presence of such "mixed operand" operations.
- b) An interval datatype, or the operations on such data. This part neither requires nor excludes such data or operations.
- c) A fixed point datatype, or the operations on such data. This part neither requires nor excludes such data or operations.
- d) A rational datatype, or the operations on such data. This part neither requires nor excludes such data or operations.
- e) Complex, matrix, statistical, or symbolic operations. This part neither requires nor excludes such data or operations.
- f) The properties of arithmetic datatypes that are not related to the numerical process, such as the representation of values on physical media.
- g) The properties of integer and floating point datatypes that properly belong in programming 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 in the programming language,
 - 4) the presence or absence of automatic datatype coercions,
 - 5) the rules for assignment, parameter passing, and returning value,
 - 6) the consequences of applying an operation to values of improper datatype, or to uninitialised data.

Furthermore, this part does not provide specifications for how the operations should be implemented or 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 programming language standards and in binding standards.

2 Conformity

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

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 specification of the binding standard takes precedence.

When a binding standard covers only a subset of the operations specified in this part, an implementation remains free to conform to this part with respect to other 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 operations that together satisfy all the requirements of clauses 5 through 8 that are relevant to those operations. The implementation shall then document the binding.

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

An implementation is free to provide 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 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. However, a mode of operation that conforms to this part should be the default mode of operation.

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 1 as well, which in turn may include (explicitly or by reference) a binding for IEC 60559 as well.
- 3 This part does not require a particular set of operations to be provided. It is not possible to conform to this part without specifying to which datatypes and set of operations (and modes of operation) conformity is claimed.

3 Normative references

The following normative documents contain provisions which, through reference in this text, constitute provisions of this part of ISO/IEC 10967. For dated references, subsequent amendments to, or revisions of, any of these publications do not apply. However, parties to agreements based on this part of ISO/IEC 10967 are encouraged to investigate the possibility of applying the most recent editions of the normative documents indicated below. For undated references, the latest edition of the normative document referred to applies. Members of ISO and IEC maintain registers of currently valid International Standards.

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

ISO/IEC 10967-1:1994, Information technology – Language independent arithmetic – Part 1: Integer and floating point arithmetic.

NOTE - See also annex E.

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{R} denotes the set of classical real numbers, and \mathcal{C} denotes the set of complex numbers over \mathcal{R} . Note that $\mathcal{Z} \subset \mathcal{R} \subset \mathcal{C}$.

The conventional notation for set definition and manipulation is used.

In this part, the following notation for intervals is used

- [x, z] designates the interval $\{y \in \mathcal{R} \mid x \leqslant y \leqslant z\},\$
- [x, z] designates the interval $\{y \in \mathcal{R} \mid x < y \leqslant 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. In particular this part uses

- \Rightarrow and \Leftrightarrow for logical implication and equivalence
- +, -, /, |x|, |x|, [x], and round(x) on reals
- · for multiplication on reals
- $\langle , \leq , \geq , \text{ and } \rangle$ between reals
- = and \neq between real as well as special values

max on non-empty upwardly closed sets of reals min on non-empty downwardly closed sets of reals

 \cup , \cap , \times , \in , \notin , \subset , \subseteq , $\not\subseteq$, and = with sets

- \times for the Cartesian product of sets
- \rightarrow for a mapping between sets

| for the divides relation between integers

For $x \in \mathcal{R}$, the notation |x| designates the largest integer not greater than x:

$$\lfloor x \rfloor \in \mathcal{Z}$$
 and $x - 1 < \lfloor x \rfloor \leqslant x$

the notation [x] designates the smallest integer not less than x:

$$\lceil x \rceil \in \mathcal{Z}$$
 and $x \leqslant \lceil x \rceil < x + 1$

and the notation round(x) designates the integer closest to x:

```
\operatorname{round}(x) \in \mathcal{Z} and x - 0.5 \leqslant \operatorname{round}(x) \leqslant x + 0.5
```

where in case x is exactly half-way between two integers, the even integer is the result.

The divides relation () on integers tests whether an integer i divides an integer j exactly:

```
i|j \Leftrightarrow (i \neq 0 \text{ and } i \cdot n = j \text{ for some } n \in \mathcal{Z})
```

NOTE -i|j is true exactly when j/i is defined and $j/i \in \mathcal{Z}$).

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 [47] (e is the Napierian base)

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

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 $x \in [0, \infty[$. The only one over which there is some difference of conventions is the arccot function. Conventions there vary for negative arguments; either a negative value (giving a sign symmetric function), or a positive return value (giving a function that is continuous over zero). In this part, arccot refers to the sign symmetric inverse function (with a branch cut at 0), and arccotc refers to the continuous inverse function.

```
\operatorname{arccosh}(x) \geqslant 0, \ \operatorname{arcsech}(x) \geqslant 0, \\ \operatorname{arcsin}(x) \in [-\pi/2, \pi/2], \ \operatorname{arccos}(x) \in [0, \pi], \ \operatorname{arctan}(x) \in ]-\pi/2, \pi/2[, \\ \operatorname{arccot}(x) \in ]-\pi/2, \pi/2], \ \operatorname{arccotc}(x) \geqslant [0, \pi], \ \operatorname{arccsc}(x) \in [0, \pi], \ \operatorname{arccsc}(x) \in [-\pi/2, \pi/2].
```

NOTE -e = 2.71828... e is not in any floating point datatype conforming to part 1, unless added as a special value, which is usually not done.

4.1.4 Exceptional values

The exceptional value underflow is used in this part as it is in part 1.

Three new exceptional values, **overflow**, **invalid**, and **infinitary**, are introduced in this part replacing three other exceptional values used in part 1. **invalid** and **infinitary** are in this part used instead of the **undefined** of part 1. **overflow** is used instead of the **integer_overflow** and **floating_overflow** of part 1. Bindings may still distinguish between **integer_overflow** and **floating_overflow**.

One new exceptional value, **absolute_precision_underflow**, is introduced in this part with no correspondence in part 1. The exceptional value **absolute_precision_underflow** is used when the given floating point angle value argument is so big that even a highly accurate result from a trigonometric operation is questionable, due to the fact that the density of floating point values has decreased significantly at these big angle values.

For the exceptional values, a continuation value may be given in parenthesis after the exceptional value.

4.1.5 Datatypes

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

NOTE 1 – Mathematical relations are true or false (or undefined, if an operand is undefined). In contrast, **true** and **false** are values in **Boolean**.

For pairs, define:

$$fst((x,y)) = x$$
$$snd((x,y)) = y$$

Square brackets are used to write finite sequences of values. [] is the sequence containing no values. [s], is the sequence of one value, s. $[s_1, s_2]$, is the sequence of two values, s_1 and then s_2 , etc. The colon operator is used to prepend a value to a sequence: $x : [x_1, ..., x_n] = [x, x_1, ..., x_n]$. [S], where S is a set, denotes the set of finite sequences, where each value in a sequence is in S.

NOTE 2 – It is always clear from context, in the text of this part, if [X] is a sequence of one element, or the set of sequences with values from X. It is also clear from context if $[x_1, x_2]$ is a sequence of two values or an interval.

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

NOTES

- 3 This implies that $fminN_F < 0.5 \cdot epsilon_F/r_F$ in this part, rather than just $fminN_F \leq epsilon_F$.
- 4 These extra requirements, which do not limit the use of any existing floating point datatype, are made 1) so that angles in radians are not too degenerate within the first two cycles, plus and minus, when represented in F, and 2) in order to be able to avoid the underflow notification in specifications for the $expm1_F$ and $ln1p_F$ operations.
- 5 F should also be such that $p_F \ge 2 + \lceil \log_{r_F}(1000) \rceil$, to allow for a not too coarse angle resolution anywhere in the interval $[-big_angle_r_F, big_angle_r_F]$. See clause 5.3.8.

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

$$-0, +\infty, -\infty$$
 qNaN, and sNaN.

These values are not part of the set I or 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 6 – This part uses the above five special values for compatibility with IEC 60559. In particular, the symbol $-\mathbf{0}$ (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, $+\infty$, $-\infty$, 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 $+\infty$, $-\infty$, or qNaN, the actual result or continuation value is binding or implementation defined.

The following symbols used in this part are defined in part 1:

Integer parameters:

 $bounded_I$, $maxint_I$, and $minint_I$.

Integer helper function:

 $wrap_I$.

Integer operations:

 neg_I , add_I , sub_I , and mul_I .

Floating point parameters:

 r_F , p_F , $emin_F$, $emax_F$, $denorm_F$, and iec_559_F .

Derived floating point constants:

 $fmax_F$, $fmin_F$, $fminN_F$, $fminD_F$, and $epsilon_F$.

Floating point rounding constant:

 rnd_error_F .

Floating point value sets related to F:

 F^* , F_D , and F_N .

Floating point helper functions:

 e_F , $result_F$, and rnd_F .

Floating point operations:

 neg_F , add_F , sub_F , mul_F , div_F , and ulp_F .

Definitions of terms 4.2

F. OTISOILEC 10961.2:2001

FULL POF OTISOILEC 10961.2:2001 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{R} , or \mathcal{C} .

NOTE 1 - This part specifies requirements for integer and floating point datatypes. Complex numbers are not covered by this part, but will be included in a subsequent part of ISO/IEC 10967 (3)

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 clause 6.1 of part 1.)

denormalisation loss: A larger than normal rounding error caused by the fact that subnormal values have less than full precision. (See clause 5.2 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 (infinitary

(e.g. at a pole), **invalid** (e.g. when the true result is in \mathcal{C} but not in \mathcal{R}), or because the mathematical result cannot, or might not, be representable with sufficient accuracy (**underflow**, **overflow**) or viability (**absolute_precision_underflow**).

NOTE 2 — The term exception is here not used to designate certain methods of handling notifications that fall under the category 'change of control flow'. Such methods of notification handling will be referred to as "[programming language name] exception", when referred to, particularly in annex C.

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

- 3 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.
- 4 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. However, some implementation defined helper functions are required to be documented.
- **implementation** (of this part): The total arithmetic environment presented to a programmer, including hardware, language processors, exception handling facilities, subroutine libraries, other software, and documentation pertinent to this part.
- **literal:** A syntactic entity, that does not have any proper sub-entity that is an expression, denoting a constant value.
- **monotonic approximation:** An approximation helper function $h: ... \times S \times ... \to \mathcal{R}$, where the other arguments are kept constant, and where $S \subseteq \mathcal{R}$, is a monotonic approximation of a predetermined mathematical function $f: \mathcal{R} \to \mathcal{R}$ if, for every $a \in S$ and $b \in S$, where a < b,
 - a) f is monotonic non-decreasing on [a,b] implies $h(...,a,...) \leq h(...,b,...)$,
 - b) f is monotonic nonfinereasing on [a, b] implies $h(..., a, ...) \ge h(..., b, ...)$.
- **monotonic non-decreasing:** A function $f: \mathcal{R} \to \mathcal{R}$ is monotonic non-decreasing on a real interval [a,b] if for every x and y such that $a \leq x \leq y \leq b$, f(x) and f(y) are well-defined and $f(x) \leq f(y)$.
- **monotonic non-increasing:** A function $f: \mathcal{R} \to \mathcal{R}$ is monotonic non-increasing on a real interval [a,b] if for every x and y such that $a \leq x \leq y \leq b$, f(x) and f(y) are well-defined and $f(x) \geq f(y)$.
- **normalised:** The non-zero values of a floating point type F that provide the full precision allowed by that type. (See F_N in clause 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 $\lim_{x\to x_0} f(x)$ 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 result may not exist (see clause 5.2 of part 1).
- **rounding function:** Any function $rnd: \mathcal{R} \to X$ (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} ,

$$x \in X \Rightarrow rnd(x) = x$$

 $x < y \Rightarrow rnd(x) \leqslant rnd(y)$

Thus, if u 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, but so that rnd(-u) = -rnd(u).
- 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. Approximation helper functions may be undefined for some argument values.

The signature $add_I: I \times I \to I \cup \{\text{overflow}\}\$ states that the operation named add_I shall accept any pair of values in I as input, and when given such input shall return either a single value in I as its output or the exceptional value **overflow** possibly accompanied by a continuation value.

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 argument(s). 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 clause 5.2 of part 1 for a full definition. In part 1 and IEC 60559 this concept is called denormal.)

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_F , and radix r_F , is $r_F^{t-p_F}$, and the ulp of a subnormal or zero value is $fminD_F$. (See clause 5.2 of part 1.)

5 Specifications for integer and floating point operations

This clause specifies a number of helper functions and operations for integer and 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 or in part 1), to mathematical functions, and to helper functions (specified in this part or in part 1). They also use special abstract values $(-\infty, +\infty, -0, \mathbf{qNaN}, \mathbf{sNaN})$. For each datatype, two of these abstract values may represent several actual values each: \mathbf{qNaN} and \mathbf{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 Basic integer operations

Clause 5.1 of part 1 specifies integer datatypes and a number of operations on values of an integer datatype. In this clause some additional operations on values of an integer datatype are specified.

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 other than as continuation values.

NOTE—For some integer operations, **infinitary** notifications may occur. For **infinitary** notifications, an infinitary continuation value is recommended. For bounded integer datatypes, $maxint_I$ or $minint_I$ may be used as replacement continuation values as appropriate, if infinitary values are not available in the datatype. For unbounded integer datatypes, however, no $maxint_I$ and $minint_I$ in I are defined, and infinitary values should be used.

5.1.1 The integer result and wrap helper functions

The $result_I$ helper function:

 $result_I: \mathcal{Z} \to I \cup \{\mathbf{overflow}\}\$

$$result_I(x) = x$$
 if $x \in I$
= **overflow** if $x \in \mathcal{Z}$ and $x \notin I$

The $wrap_I$ helper function:

where $n \in \mathcal{Z}$ is chosen such that the result is in I.

NOTES

- 1 $n = \lfloor (x minint_I)/(maxint_I minint_I + 1) \rfloor$ if $x \in \mathcal{Z}$ and $bounded_I = \mathbf{true}$; or equivalently $n = [(x - maxint_I)/(maxint_I - minint_I + 1)]$ if $x \in \mathcal{Z}$ and $bounded_I = \mathbf{true}$.
- 2 For some wrapping basic arithmetic operations this n is computed by the ' $_ov$ ' operations in clause 5.1.9.
- 3 The $wrap_I$ helper function is also used in part 1.

5.1.2 Integer maximum and minimum

$$1 \quad n = \lfloor (x - minint_I) / (maxint_I - minint_I + 1) \rfloor \text{ if } x \in \mathcal{Z} \text{ and } bounded_I = \mathbf{true}; \text{ or equin} \\ n = \lceil (x - maxint_I) / (maxint_I - minint_I + 1) \rceil \text{ if } x \in \mathcal{Z} \text{ and } bounded_I = \mathbf{true}.$$

$$2 \quad \text{For some wrapping basic arithmetic operations this } n \text{ is computed by the `-ov' operation in clause } 5.1.9.$$

$$3 \quad \text{The } wrap_I \text{ helper function is also used in part } 1.$$

$$2 \quad \text{Integer maximum and minimum}$$

$$max_I : I \times I \to I$$

$$max_I(x,y) = \max\{x,y\} \qquad \text{if } x,y \in I$$

$$min_I : I \times I \to I$$

$$min_I(x,y) = \min\{x,y\} \qquad \text{if } x,y \in I$$

$$max_seq_I : [I] \to I \cup \{\text{infinitary}\}$$

$$max_seq_I : [I] \to I \cup \{\text{infinitary}\}$$

$$max_seq_I : [I] \to I \cup \{\text{infinitary}\}$$

$$min_seq_I : [I] \to I \cup \{\text{infinitary}\}$$

5.1.3 Integer diminish

$$dim_I: I \times I \to I \cup \{\mathbf{overflow}\}\$$
 $dim_I(x,y) = result_I(\max\{0, x - y\}) \text{ if } x, y \in I$

NOTE – dim_I cannot be implemented as $max_I(0, sub_I(x, y))$ for bounded integer types, since this latter expression has other overflow properties.

5.1.4 Integer power and arithmetic shift

 $power_I: I \times I \rightarrow I \cup \{ overflow, infinitary, invalid \}$

$$power_I(x,y) = result_I(x^y)$$
 if $x, y \in I$ and $(y > 0 \text{ or } |x| = 1)$

$$= 1$$
 if $x \in I$ and $x \neq 0$ and $y = 0$

$$= invalid(1)$$
 if $x = 0$ and $y = 0$

$$= infinitary(+\infty)$$
 if $x = 0$ and $y \in I$ and $y < 0$

$$= invalid(0)$$
 if $x, y \in I$ and $x \notin \{-1, 0, 1\}$ and $y < 0$

$$\begin{array}{ll} \mathit{shift2}_I: I \times I \to I \cup \{\mathbf{overflow}\} \\ \mathit{shift2}_I(x,y) &= \mathit{result}_I(\lfloor x \cdot 2^y \rfloor) & \text{if } x,y \in I \end{array}$$

$$shift 10_I: I \times I \to I \cup \{\mathbf{overflow}\}$$

 $shift 10_I(x,y) = result_I(\lfloor x \cdot 10^y \rfloor)$ if $x,y \in I$

5.1.5 Integer square root

$$shift 2_I(x,y) = result_I(\lfloor x \cdot 2^y \rfloor)$$
 if $x,y \in I$ $shift 10_I : I \times I \to I \cup \{ \mathbf{overflow} \}$ $shift 10_I(x,y) = result_I(\lfloor x \cdot 10^y \rfloor)$ if $x,y \in I$ Integer square root $sqrt_I : I \to I \cup \{ \mathbf{invalid} \}$ $sqrt_I(x) = \lfloor \sqrt{x} \rfloor$ if $x \in I$ and $x \neq 0$ $= \mathbf{invalid}(\mathbf{qNaN})$ if $x \in I$ and $x \neq 0$

5.1.6 Divisibility tests

$$divides_I: I \times I \to \mathbf{Boolean}$$

 $divides_I(x,y) = \mathbf{true}$
 $= \mathbf{false}$ if $x, y \in I$ and $x|y$
if $x, y \in I$ and not $x|y$

NOTES

- 1 $divides_I(0,0) =$ **false**, since 0 does not divide anything, not even 0.
- 2 divides cannot be implemented as, e.g., $eq_I(0, mod_I(y, x))$, since the remainder functions give notifications for a zero second argument.

if $x \in I$ and 2|x

$$even_I: I \to \mathbf{Boolean}$$
 $even_I(x) = \mathbf{true}$ if $x \in I$ and $2|x$ if $x \in I$ and not $2|x$ $odd_I: I \to \mathbf{Boolean}$ $odd_I(x) = \mathbf{true}$ if $x \in I$ and not $2|x$

= false

5.1.7 Integer division (with floor, round, or ceiling) and remainder

```
\begin{array}{ll} quot_I: I \times I \to I \cup \{\mathbf{overflow}, \mathbf{infinitary}, \mathbf{invalid}\} \\ quot_I(x,y) &= result_I(\lfloor x/y \rfloor) & \text{if } x,y \in I \text{ and } y \neq 0 \\ &= \mathbf{infinitary}(+\infty) & \text{if } x \in I \text{ and } x > 0 \text{ and } y = 0 \\ &= \mathbf{invalid}(\mathbf{qNaN}) & \text{if } x = 0 \text{ and } y = 0 \\ &= \mathbf{infinitary}(-\infty) & \text{if } x \in I \text{ and } x < 0 \text{ and } y = 0 \end{array}
```

NOTE $-quot_I(minint_I, -1)$, for a bounded signed integer datatype where $minint_I = -maxint_I - 1$, is the only case where this operation will overflow.

```
mod_I: I \times I \to I \cup \{\mathbf{invalid}\}
                   = x - (|x/y| \cdot y)
                                                    if x, y \in I and y \neq 0
mod_I(x,y)
                    = invalid(qNaN)
                                                    if x \in I and y = 0
ratio_I: I \times I \to I \cup \{\text{overflow}, \text{infinitary}, \text{invalid}\}
                                                     if x, y \in I and y \neq 0
ratio_I(x,y)
                    = result_I(round(x/y))
                    = infinitary(+\infty)
                                                     if x \in I and x > 0 and y = 0
                    = invalid(qNaN)
                                                     if x = 0 and y = 0
                    = infinitary(-\infty)
                                                     if x \in I and x < 0 and y = 0
residue_I: I \times I \to I \cup \{\mathbf{overflow}, \mathbf{invalid}\}
                  = result_I(x - (round(x/y) \cdot y))
residue_I(x,y)
                                                     if x, y \in I and y \neq 0
                                                    if x \in I and y = 0
group_I: I \times I \to I \cup \{\text{overflow}, \text{infinitary}, \text{invalid}\}
group_I(x,y)
                    = result_I([x]y])
                                                     if x, y \in I and y \neq 0
                    = infinitary(+\infty)
                                                     if x \in I and x > 0 and y = 0
                    = invalid(qNaN)
                                                     if x = 0 and y = 0
                    = \inf_{\mathbf{m}} \operatorname{finitary}(-\infty)
                                                     if x \in I and x < 0 and y = 0
                                                  if x, y \in I and y \neq 0
                                                    if x \in I and y = 0
```

5.1.8 Greatest common divisor and least common positive multiple

$$\begin{split} lcm_I(x,y) &= result_I(\min\{v \in \mathcal{Z} \mid x | v \text{ and } y | v \text{ and } v > 0\}) \\ & \text{if } x,y \in I \text{ and } x \neq 0 \text{ and } y \neq 0 \\ &= 0 & \text{if } x,y \in I \text{ and } (x = 0 \text{ or } y = 0) \end{split}$$

$$gcd_seq_I : [I] \rightarrow I \cup \{ \text{overflow}, \text{infinitary} \}$$

$$gcd_seq_I([x_1,...,x_n]) \\ &= result_I(\max\{v \in \mathcal{Z} \mid v | x_i \text{ for all } i \in \{1,...,n\}\}) \\ &= \inf\{x_1,...,x_n\} \subseteq I \text{ and } \{x_1,...,x_n\} \nsubseteq \{0\} \\ &= \text{infinitary}(+\infty) & \text{if } \{x_1,...,x_n\} \subseteq \{0\} \end{split}$$

$$lcm_seq_I : [I] \rightarrow I \cup \{ \text{overflow} \}$$

$$lcm_seq_I([x_1,...,x_n]) \\ &= result_I(\min\{v \in \mathcal{Z} \mid x_i | v \text{ for all } i \in \{1,...,n\} \text{ and } v > 0\}) \\ &= \inf\{x_1,...,x_n\} \subseteq I \text{ and } 0 \notin \{x_1,...,x_n\} \\ &= 0 & \text{if } \{x_1,...,x_n\} \subseteq I \text{ and } 0 \in \{x_1,...,x_n\} \end{split}$$

NOTE – These specifications imply: $gcd_seq_I([]) = infinitary(+\infty)$ and $lcm_seq_I([]) = 1$.

5.1.9 Support operations for extended integer range

These operations can be used to implement extended range integer datatypes, including unbounded integer datatypes.

```
add\_wrap_I: I \times I \to I
add\_wrap_I(x,y) = wrap_I(x+y)
add\_ov_I: I \times I \to \{-1,0,1\}
add\_ov_I(x,y) = ((x+y) - add\_wrap_I(x,y))/(maxint_I - minint_I + 1)
if \ x,y \in I \ and \ bounded_I = \mathbf{true}
if \ x,y \in I \ and \ bounded_I = \mathbf{false}
sub\_wrap_I: I \times I \to I
sub\_ov_I: I \times I \to \{-1,0,1\}
if \ x,y \in I \ and \ bounded_I = \mathbf{true}
if \ x,y \in I \ and \ bounded_I = \mathbf{true}
if \ x,y \in I \ and \ bounded_I = \mathbf{false}
mul\_wrap_I: I \times I \to I
mul\_wrap_I(x,y) = wrap_I(x \cdot y)
if \ x,y \in I
```

```
mul\_ov_I: I \times I \to I

mul\_ov_I(x,y) = ((x \cdot y) - mul\_wrap_I(x,y))/(maxint_I - minint_I + 1)

if x,y \in I and bounded_I = \mathbf{true}

= 0 if x,y \in I and bounded_I = \mathbf{false}
```

NOTE – The add_ov_I and sub_ov_I will only return -1 (for negative overflow), 0 (no overflow), and 1 (for positive overflow).

5.2 Basic floating point operations

Clause 5.2 of part 1 specifies floating point datatypes and a number of operations on values of a floating point datatype. In this clause some additional operations on values of a floating point datatype are specified.

NOTE – Further operations on values of a floating point datatype, for elementary 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.

5.2.1 The rounding and floating point result helper functions

Floating point rounding helper functions (F^* is defined in part 1):

The floating point helper function

```
down_F: \mathcal{R} \to F^*
```

is the rounding function that rounds towards negative infinity. The floating point helper function

$$up_F:\mathcal{R} o F^*$$

is the rounding function that rounds towards positive infinity. The floating point helper function

```
nearest_F: \mathcal{R} \to F^*
```

is the rounding function that rounds to nearest. $nearest_F$ is partially implementation defined: the handling of ties is implementation defined, but must be sign symmetric. If $iec_559_F = \mathbf{true}$, the semantics of $nearest_F$ is completely defined by IEC 60559: in this case ties are rounded so that the result has an even last digit.

 $result_F$ is a helper function that is partially implementation defined. $result_F$ has a signature:

```
result_F : \mathcal{R} \times (\mathcal{R} \to F^*) \to F \cup \{\text{underflow}, \text{overflow}\}
```

For the overflow cases it is defined as:

```
result_F(x, nearest_F) = \mathbf{overflow}(+\infty) \qquad \text{if } x \in \mathcal{R} \text{ and } nearest_F(x) > fmax_F
result_F(x, nearest_F) = \mathbf{overflow}(-\infty) \qquad \text{if } x \in \mathcal{R} \text{ and } nearest_F(x) < -fmax_F
result_F(x, up_F) \qquad = \mathbf{overflow}(+\infty) \qquad \text{if } x \in \mathcal{R} \text{ and } up_F(x) > fmax_F
result_F(x, up_F) \qquad = \mathbf{overflow}(-fmax_F) \qquad \text{if } x \in \mathcal{R} \text{ and } up_F(x) < -fmax_F
result_F(x, down_F) \qquad = \mathbf{overflow}(fmax_F) \qquad \text{if } x \in \mathcal{R} \text{ and } down_F(x) > fmax_F
result_F(x, down_F) \qquad = \mathbf{overflow}(-\infty) \qquad \text{if } x \in \mathcal{R} \text{ and } down_F(x) < -fmax_F
```

For other cases and for any rounding function rnd in $(\mathcal{R} \to F^*)$, the following shall apply:

```
if x = 0
      result_F(x, rnd) = x
                                                       if x \in \mathcal{R} and fminN_F \leq |x| and |rnd(x)| \leq fmax_F
                         = rnd(x)
                         = rnd(x) or underflow(c)
                                                      if x \in \mathcal{R} and |x| < fminN_F and |rnd(x)| = fminN_F
                                                         and rnd has no denormalisation loss at x
                         = rnd(x) or underflow(c)
                                                       if x \in \mathcal{R} and denorm_F = \mathbf{true} and
                                                         |rnd(x)| < fminN_F \text{ and } x \neq 0
                                                         and rnd has no denormalisation loss at x
                         = underflow(c)
                                                       otherwise
where
                                                       when denorm_F = \mathbf{true} and (rnd(x) \neq 0 \text{ or } x > 0),
                         c = rnd(x)
                                                       when denorm_F = \mathbf{true} and rnd(x) = 0 and x < 0,
                         c = -\mathbf{0}
                         c = 0
                                                       when denorm_F = false and x > 0,
                                                       when denorm_F = false and x < 0
                         c = -0
```

An implementation is allowed to choose between rnd(x) and **underflow** (rnd(x)) in the region between 0 and $fminN_F$. However, a subnormal value without an underflow notification can be chosen only if $denorm_F = \mathbf{true}$ and no denormalisation loss occurs at x.

NOTES

- 1 This differs from the specification of $result_F$ as given in part 1 in the following respects: 1) the continuation values on overflow and underflow are given directly here, and 2) all instances of denormalisation loss must be accompanied with an underflow notification.
- 2 $denorm_F =$ false implies $iec_559_F =$ false, and $iec_559_F =$ true implies $denorm_F =$ true.
- 3 If $iec_559_F = \mathbf{true}$, then subnormal or zero results that have no denormalisation loss do not result in an underflow notification, if the notification is by recording of indicators.

Define the no_result_F , no_result_F , and no_result_F helper functions:

```
no\_result_F : F \rightarrow \{\mathbf{invalid}\}
                                                 if x \in F \cup \{-\infty, -0, +\infty\}
no\_result_F(x)
                   = invalid(qNaN)
                                                 if x is a quiet NaN
                   = qNaN
                   = invalid(qNaN)
                                                 if x is a signalling NaN
                                                 if x, y \in F \cup \{-\infty, -0, +\infty\}
                   = invalid(qNaN)
                   = qNaN
                                                 if at least one of x and y is a quiet NaN and
                                                   neither a signalling NaN
                   = invalid(qNaN)
                                                 if x is a signalling NaN or y is a signalling NaN
no\_result 3_F : F \times F \times F \to \{invalid\}
no\_result \mathcal{J}_F(x,y,z)
                   = invalid(qNaN)
                                                 if x, y, z \in F \cup \{-\infty, -0, +\infty\}
                   = qNaN
                                                 if at least one of x, y, and z is a quiet NaN and
                                                   neither is a signalling NaN
```

= invalid(qNaN) if at least one of x, y, or z 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) is the appropriate result.

NOTE 4 - The handling of other special values, if available, is left unspecified by this part.

5.2.2 Floating point maximum and minimum

The appropriate return value of the maximum and minimum operations given a quiet NaN (qNaN) as one of the input values depends on the circumstances for each point of use. Sometimes qNaN is the appropriate result, sometimes the non-NaN argument is the appropriate result. Therefore, two variants each of the floating point maximum and minimum operations are specified here, and the programmer can decide which one is appropriate to use at each particular place of usage, assuming both variants are included in the binding.

```
max_F: F \times F \to F
max_F(x,y)
                                                         if x, y \in F
                      = \max\{x, y\}
                                                         if x = +\infty and y \in \mathbb{A} \cup \{-\infty, -0\}
                      =+\infty
                                                         if x = -\mathbf{0} and y \in \mathbb{R} and y \geqslant 0
                      = y
                                                         if x = -\mathbf{0} and (y \in F \text{ and } y < 0) or y = -\mathbf{0}
                      = -0
                                                         if x = -\infty and y \in F \cup \{+\infty, -0\}
                      = y
                                                         if y = +\infty and x \in F \cup \{+\infty, -0\}
                      =+\infty
                                                         if y = \mathbf{0} and x \in F and x \ge 0
if y = \mathbf{0} and x \in F and x < 0
                      = x
                      = -0
                                                         if y = -\infty and x \in F \cup \{-\infty, -0\}
                      = x
                                                   - −∝
otherwise
                      = no\_result2_F(x, y)
min_F: F \times F \to F
                                                         if x, y \in F
                                                         if x = +\infty and y \in F \cup \{-\infty, -0\}
                                                         if x = -0 and y \in F and y \ge 0
                                                         if x = -0 and ((y \in F \text{ and } y < 0) \text{ or } y = -0)
                                                         if x = -\infty and y \in F \cup \{+\infty, -0\}
                                                         if y = +\infty and x \in F \cup \{+\infty, -0\}
                                                         if y = -\mathbf{0} and x \in F and x \ge 0
                                                         if y = -0 and x \in F and x < 0
                                                         if y = -\infty and x \in F \cup \{-\infty, -0\}
                                                         otherwise
mmax_F: F \times F \to F
                                                         if x, y \in F \cup \{+\infty, -0, -\infty\}
mmax_F(x,y)
                      = max_F(x,y)
                                                         if x \in F \cup \{+\infty, -0, -\infty\} and y is a quiet NaN
                      = x
                                                         if y \in F \cup \{+\infty, -0, -\infty\} and x is a quiet NaN
                      = no\_result2_F(x, y)
                                                         otherwise
```

 $mmin_F: F \times F \to F$

```
mmin_F(x,y)
                   = min_F(x, y)
                                                   if x, y \in F \cup \{+\infty, -0, -\infty\}
                                                   if x \in F \cup \{+\infty, -0, -\infty\} and y is a quiet NaN
                                                   if y \in F \cup \{+\infty, -0, -\infty\} and x is a quiet NaN
                                                   otherwise
                    = no\_result2_F(x,y)
max\_seq_F : [F] \to F \cup \{ \mathbf{infinitary} \}
max\_seq_F([x_1, ..., x_n])
                   = infinitary(-\infty)
                                                   if n = 0
                                                   if n = 1 and x_1 is not a NaN
                    = max_F(max\_seq_F([x_1, ..., x_{n-1}]), x_n)
                                                   if n \ge 2
                   = no\_result_F(x_1)
                                                   otherwise
min\_seq_F : [F] \to F \cup \{ infinitary \}
min\_seq_F([x_1,...,x_n])
                   = infinitary(+\infty)
                                                   if n = 0
                                                   if n = 1 and x_1 is not a NaN
                    = min_F(min\_seq_F([x_1, ..., x_{n-1}]), x_n)
                                                   if n \geqslant 2
                   = no\_result_F(x_1)
                                                   otherwise
mmax\_seq_F : [F] \to F \cup \{\mathbf{infinitary}\}
mmax\_seq_F([x_1,...,x_n])
                    = infinitary(-\infty)
                                                   \Re n = 1 and x_1 is not a NaN
                    = mmax_F(mmax\_seq_F([x_1, ..., x_{n-1}]), x_n)
                                                   if n \geqslant 2
                                                   otherwise
mmin\_seq_F : [F] \to F \cup \{infinitary\}
mmin\_seq_F([x_1,...,x_n])
                     = infinitary(+\infty)
                                                   if n = 0
                                                   if n = 1 and x_1 is not a NaN
                    = mmin_F(mmin\_seq_F([x_1, ..., x_{n-1}]), x_n)
                                                   if n \geqslant 2
                                                   otherwise
```

5.2.3 Floating point diminish

```
\begin{aligned} dim_F: F \times F &\to F \cup \{\mathbf{underflow}, \mathbf{overflow}\} \\ dim_F(x,y) &= result_F(\max\{0,x-y)\}, rnd_F) \\ &= \mathbf{if} \ x,y \in F \\ &= -\mathbf{0} & \text{if} \ x = -\mathbf{0} \ \text{and} \ y = 0 \\ &= dim_F(0,y) & \text{if} \ x = -\mathbf{0} \ \text{and} \ y \in F \cup \{-\infty, -\mathbf{0}, +\infty\} \ \text{and} \ y \neq 0 \\ &= dim_F(x,0) & \text{if} \ y = -\mathbf{0} \ \text{and} \ x \in F \cup \{-\infty, +\infty\} \end{aligned}
```

$$\begin{array}{ll} = +\infty & \text{if } x = +\infty \text{ and } y \in F \cup \{-\infty\} \\ = 0 & \text{if } x = -\infty \text{ and } y \in F \cup \{+\infty\} \\ = 0 & \text{if } y = +\infty \text{ and } x \in F \\ = +\infty & \text{if } y = -\infty \text{ and } x \in F \\ = no_result2_F(x, y) & \text{otherwise} \end{array}$$

NOTE – dim_F cannot be implemented by $max_F(-\mathbf{0}, sub_F(x, y))$, since this latter expression has other overflow properties.

5.2.4 Floor, round, and ceiling

```
floor_F: F \to F
floor_F(x) = \lfloor x \rfloor
                                                        if x \in F
                                                       if x \in \{-\infty, -0, +\infty\}
                                                        otherwise
                     = no\_result_F(x)
floor\_rest_F: F \to F
floor\_rest_F(x) = result_F(x - \lfloor x \rfloor, rnd_F) \text{ if } x \in F
                                                        if x = -\mathbf{0}
                                                        otherwise
                     = no\_result_F(x)
rounding_F: F \to F \cup \{-\mathbf{0}\}
                                                        if x \in F and (x \ge 0 \text{ or round}(x) \ne 0)
rounding_F(x) = round(x)
                                                        \mathbf{i}(x) \in F \text{ and } x < 0 \text{ and } \text{round}(x) = 0
                     = -0
                                                       \checkmark if x \in \{-\infty, -0, +\infty\}
                     = x
                                                        otherwise
                     = no\_result_F(x)
NOTE – round_F is a different operation specified in part 1.
```

```
rounding\_rest_F: F \to F
rounding\_rest_F(x)
= x - round(x) \qquad \text{if } x \in F
\text{if } x = -\mathbf{0}
\text{otherwise}
ceiling_F(x) = \lceil x \rceil \qquad \text{otherwise}
ceiling_F(x) = \lceil x \rceil \qquad \text{if } x \in F \text{ and } (x \geqslant 0 \text{ or } \lceil x \rceil \neq 0)
= -\mathbf{0} \qquad \text{if } x \in F \text{ and } x < 0 \text{ and } \lceil x \rceil = 0
= x \qquad \text{if } x \in \{-\infty, -\mathbf{0}, +\infty\}
= no\_result_F(x) \qquad \text{otherwise}
```

 $ceiling_rest_F: F \to F$

$$ceiling_rest_F(x)$$

= $result_F(x - \lceil x \rceil, rnd_F)$ if $x \in F$
= 0 if $x = -\mathbf{0}$
= $no_result_F(x)$ otherwise

5.2.5 Remainder after division with round to integer

```
residue_F: F \times F \to F \cup \{ \mathbf{-0}, \mathbf{underflow}, \mathbf{invalid} \}
residue_F(x,y) = result_F(x - (\operatorname{round}(x/y) \cdot y), nearest_F)
\text{if } x, y \in F \text{ and } y \neq 0 \text{ and}
(x \geqslant 0 \text{ or } x - (\operatorname{round}(x/y) \cdot y) \neq 0)
= \mathbf{-0} \qquad \text{if } x, y \in F \text{ and } y \neq 0 \text{ and}
x < 0 \text{ and } x - (\operatorname{round}(x/y) \cdot y) = \mathbf{0}
= \mathbf{-0} \qquad \text{if } x = \mathbf{-0} \text{ and } y \in F \cup \{ \mathbf{-\infty}, \mathbf{+\infty} \} \text{ and } y \neq 0 \text{ and}
= x \qquad \text{if } x \in F \text{ and } y \in \{ \mathbf{-\infty}, \mathbf{+\infty} \} \text{ and } y \neq 0 \text{ otherwise}
```

5.2.6 Square root and reciprocal square root

$$sqrt_F: F \to F \cup \{\mathbf{invalid}\}$$
 $sqrt_F(x) = nearest_F(\sqrt{x}) \quad \text{if } x \in F \text{ and } x \geqslant 0$
 $= x \quad \text{if } x \in \{-0, +\infty\}$
 $= no_result_F(x) \quad \text{otherwise}$
 $rec_sqrt_F: F \to F \cup \{\mathbf{infinitary}, \mathbf{invalid}\}$
 $rec_sqrt_F(x) = rnd_F(1/\sqrt{x}) \quad \text{if } x \in F \text{ and } x > 0$
 $= \mathbf{infinitary}(+\infty) \quad \text{if } x \in \{-\mathbf{0}, 0\}$
 $= 0 \quad \text{if } x = +\infty$
 $= no_result_F(x) \quad \text{otherwise}$

5.2.7 Multiplication to higher precision floating point datatype

For the following operation, F' is a floating point datatype conforming to part 1, where $r_{F'} = r_F$ and $p_{F'} > p_F$.

```
mul_{F \to F'} F \times F \to F' \cup \{-0, \text{underflow}, \text{overflow}\}
mul_{F \to F'}(x, y) = mul_{F'}(convert_{F \to F'}(x), convert_{F \to F'}(y))
NOTES
```

- 1 $convert_{F \to F'}$ is specified in clause 5.4.4.
- 2 F' has the same radix as, but higher precision than F. If the precision is sufficiently much higher, rounding can be avoided. If also $emin_{F'}$ is sufficiently smaller than $emin_F$, underflow can be avoided, and if $emax_{F'}$ is sufficiently greater than $emax_F$, overflow can be avoided.

5.2.8 Support operations for extended floating point precision

These operations are useful when keeping guard digits or implementing extra precision floating point datatypes. The resulting datatypes, e.g. so-called doubled precision, do not necessarily conform to part 1.

```
add\_lo_F : F \times F \to F \cup \{\mathbf{underflow}\}\
add\_lo_F(x,y) = result_F((x+y) - rnd_F(x+y), rnd_F)
                                                      if x, y \in F
                     = -0
                                                      if x = -\mathbf{0} and y \in F \cup \{-\infty, -\mathbf{0}, +\infty\}
                     = -0
                                                      if x \in F \cup \{-\infty, +\infty\} and y = -0
                     = y
                                                      if x = +\infty and y \in F \cup \{+\infty\}
                                                                          711EC 10967.2:
                                                      if x = -\infty and y \in F \cup \{-\infty\}
                     = y
                                                      if x \in F and y \in \{-\infty, +\infty\}
                     = no\_result2_F(x,y)
                                                      otherwise
sub\_lo_F: F \times F \to F \cup \{\mathbf{underflow}\}\
sub\_lo_F(x, y) = add\_lo_F(x, neg_F(y))
```

NOTE 1 – If $rnd_style_F = nearest$, then, in the absence of notifications, add_lo_F and sub_lo_F return exact results.

```
\begin{aligned} mul\_lo_F: F \times F &\to F \cup \{ \mathbf{underflow}, \mathbf{overflow} \} \\ mul\_lo_F(x,y) &= result_F((x \cdot y) - rnd_F(x \cdot y), rnd_F) \\ &= mul\_lo_F(0,y) & \text{if } x,y \in F \\ &= mul\_lo_F(x,0) & \text{if } x = -\mathbf{0} \text{ and } y \in F \cup \{-\infty, -\mathbf{0}, +\infty\} \\ &= mul\_lo_F(x,y) & \text{if } x \in F \cup \{-\infty, +\infty\} \text{ and } y = -\mathbf{0} \\ &= mul_F(x,y) & \text{if } x \in \{-\infty, +\infty\} \text{ and } y \in F \cup \{-\infty, +\infty\} \\ &= mul_F(x,y) & \text{if } x \in F \text{ and } y \in \{-\infty, +\infty\} \\ &= no\_result2_F(x,y) & \text{otherwise} \end{aligned}
```

NOTE 2 – In the absence of notifications, $mul lo_F$ returns an exact result.

```
\begin{array}{ll} div\_rest_F: F\times F\to F\cup \{\mathbf{underflow},\mathbf{invalid}\}\\ div\_rest_F(x,y) &= result_F(x-(y\cdot rnd_F(x/y)),rnd_F)\\ &= fx,y\in F\\ &= div\_rest_F(0,y) & \text{if } x=-\mathbf{0} \text{ and } y\in F\cup \{-\infty,-\mathbf{0},+\infty\}\\ &= x & \text{if } x\in F \text{ and } y\in \{-\infty,+\infty\}\\ &= x & \text{if } x\in \{-\infty,+\infty\} \text{ and } y\in F\\ &= no\_result2_F(x,y) & \text{otherwise} \end{array} sqrt\_rest_F: F\to F\cup \{\mathbf{underflow},\mathbf{invalid}\}
```

```
sqrt\_rest_F(x) = result_F(x - (sqrt_F(x) \cdot sqrt_F(x)), rnd_F)
= \mathbf{-0}
if x \in F and x \ge 0
= \mathbf{-0}
```

$$= +\infty$$
 if $x = +\infty$
= $no_result_F(x)$ otherwise

NOTE 3 - $sqrt_rest_F(x)$ is exact when there is no **underflow**.

5.3 Elementary transcendental floating point operations

The specifications for each of the floating point transcendental operations and the floating point conversion operations (clause 5.4) use an approximation helper function. The approximation helper functions are ideally identical to the true mathematical functions. However, that would imply a maximum error for the corresponding operation of 0.5 ulp (i.e., the minimum value for operations that are not always exact). This part does not require that the maximum error is only 0.5 ulp for the operations specified in clause 5.3, but allows the maximum error to be a bit bigger. To express this, the approximation helper functions need not be identical to the mathematical elementary transcendental functions, but are allowed to be approximate. The approximation helper functions shall be defined for the elements of its given argument signature where the corresponding mathematical function is also defined, unless otherwise noted. The requirements on approximation helper functions apply only where the approximation helper functions are defined.

5.3.1 Maximum error requirements

The approximation helper functions for the individual operations in these subclauses have maximum error parameters that describe the maximum relative error, in ulps, of the helper function composed with $nearest_F$, for non-subnormal and non-zero results. The maximum error parameters also describe the maximum absolute error, in ulps, for $-fminN_F$, $fminN_F$, subnormal, or zero results and underflow continuation values if $denorm_F = \mathbf{true}$. All maximum error parameters shall have a value that is ≥ 0.5 . For the maximum value for the maximum error parameters, see the specification of each of the maximum error parameters. See also Annex A, on partial conformity. The relevant maximum error parameters shall be made available to programs.

When the maximum error for an approximation helper function h_F , approximating f, is $max_error_op_F$, then for all arguments $x, ... \in F \times ...$ the following equation shall hold:

$$|f(x,...) - nearest_F(h_F(x,...))| \le max_error_op_F \cdot r_F^{e_F(f(x,...))-p_F}$$

NOTES

- 1 Partially conforming implementations may have greater values for maximum error parameters than stipulated below. See annex A.
- For most positive (and not too small) return values t, the true result is thus claimed to be in the interval $[t-(max_error_op_F \cdot ulp_F(t)), t+(max_error_op_F \cdot ulp_F(t))]$. But if the return value is exactly r_F^n for some not too small $n \in \mathcal{Z}$, then the true result is claimed to be in the interval $[t-(max_error_op_F \cdot ulp_F(t)/r_F), t+(max_error_op_F \cdot ulp_F(t))]$. Similarly for negative return values.

The results of the approximating helper functions in this clause must be exact for certain arguments as detailed below, and may be exact for all arguments. If the approximating helper function is exact for all arguments, then the corresponding maximum error parameter should have the value 0.5, the minimum value.

5.3.2 Sign requirements

For this part, the approximation helper functions shall be zero exactly at the points where the approximated mathematical function is exactly zero. For this part, at points where the approximation helper functions are not zero, they shall have the same sign as the approximated mathematical function at that point. For the radian trigonometric helper functions, these zero and sign requirements are imposed only for arguments, x, such that $|x| \leq big_angle_r_F$ (see clause 5.3.8).

NOTE – For the operations, the continuation value after an **underflow** may be zero (including negative zero) as given by $result_F^*$ (see below), even though the approximation helper function is not zero at that point. Such zero results are required to be accompanied by an **underflow** notification. When appropriate, zero may also be returned for IEC 60559 infinities arguments. See the individual specifications.

5.3.3 Monotonicity requirements

For this part, each approximation helper function shall be a monotonic approximation to the mathematical function it is approximating, except:

- a) For the radian trigonometric approximation helper functions, the monotonic approximation requirement is imposed only for arguments, x, such that $|x| \leq big_angle_r_F$ (see clause 5.3.8).
- b) The argument angular unit trigonometric and argument angular unit inverse trigonometric approximating helper functions, as well as the angular unit conversion helper functions, are excepted from the monotonic approximation requirement for the angular unit argument(s).

5.3.4 The $result^*$ helper function

The $result_F^*$ helper function is similar to the $result_F$ helper function (see clause 5.2.1), but is simplified compared to $result_F$ concerning underflow: $result_F^*$ always underflows for non-zero arguments that have an absolute value less than $fminN_F - (fminD_F/r_F)$, whereas $result_F$ does not necessarily underflow in that case. This difference from $result_F$ is made since the argument to $result_F^*$ might not be exact. To return underflow or not, for a tiny result, based upon an inexact argument would be misleading. For the operations specified using $result_F^*$ where the specification implies that there can never be any denormalisation loss for certain tiny results, underflow is instead explicitly avoided.

```
result_F^*: \mathcal{R} \times (\mathcal{R} \to F^*) \to F \cup \{\mathbf{underflow}, \mathbf{overflow}\}
result_F^*(x, rnd) = \mathbf{underflow}(c) \qquad \text{if } x \in \mathcal{R} \text{ and } denorm_F = \mathbf{true} \text{ and } |rnd(x)| < fminN_F \text{ and } x \neq 0
= result_F(x, rnd) \qquad \text{otherwise}
where
c = rnd(x) \qquad \text{when } rnd(x) \neq 0 \text{ or } x > 0,
c = -\mathbf{0} \qquad \text{when } rnd(x) = 0 \text{ and } x < 0
```

5.3.5Hypotenuse

There shall be a maximum error parameter for the $hypot_F$ operation:

$$max_error_hypot_F \in F$$

The $max_error_hypot_F$ parameter shall have a value that is ≤ 1 .

The $hypot_F^*$ approximation helper function:

$$hypot_F^*: F \times F \to \mathcal{R}$$

 $hypot_F^*(x,y)$ returns a close approximation to $\sqrt{x^2+y^2}$ in \mathcal{R} , with maximum error $max_error_hypot_F$.

Further requirements on the $hypot_F^*$ approximation helper function are:

```
If x, y \in F if x, y \in F and \sqrt{x^2 + y^2} if x, y \in F and \sqrt{x^2 + y^2}
hypot_F^*(x,y) = hypot_F^*(y,x)
hypot_F^*(-x,y) = hypot_F^*(x,y)
hypot_F^*(x,y) \geqslant \max\{|x|,|y|\}
hypot_F^*(x,y) \leq |x| + |y|
hypot_F^*(x,y) \geqslant 1
hypot_F^*(x,y) \leq 1
```

The $hypot_F$ operation:

```
hypot_F: F \times F \to F \cup \{ underflow, overflow \}
hypot_F(x,y) = result_F^*(hypot_F^*(x,y), nearest_F)
                                                                 if x, y \in \mathbb{R}
                                                                if x = \mathbf{0} and y \in F \cup \{-\infty, -0, +\infty\}
if y = \mathbf{0} and x \in F \cup \{-\infty, +\infty\}
                         = hypot_F(x,0)
                                                                 if x \in \{-\infty, +\infty\} and y \in F \cup \{-\infty, +\infty\}
                         =+\infty
                         = +\infty
= no\_result2_F(x, y)
                                                              Sif y \in \{-\infty, +\infty\} and x \in F
```

5.3.6 Operations for exponentiations and logarithms

There shall be two maximum error parameters for approximate exponentiations and logarithms:

```
max\_error\_exp_F \in F
max\_error\_power_{F} \in F
```

The $max_error_exp_F$ parameter shall have a value that is $\leq 1.5 \cdot rnd_error_F$. The $max_error_power_F$ parameter shall have a value that is $\leq 2 \cdot rnd_error_F$.

5.3.6.1 Integer power of argument base

The $power_{F,I}^*$ approximation helper function:

$$power_{FI}^*: F \times I \to \mathcal{R}$$

 $power_{FI}^*(x,y)$ returns a close approximation to x^y in \mathcal{R} , with maximum error $max_error_power_F$. Further requirements on the $power_{F,I}^*$ approximation helper function are:

```
\begin{array}{ll} power_{F,I}^*(x,y)=x^y & \text{if } x\in\mathcal{Z}\cap F \text{ and } y\in I \text{ and } (|x|=1 \text{ or } y>0) \\ power_{F,I}^*(x,1)=x & \text{if } x\in F \\ power_{F,I}^*(x,0)=1 & \text{if } x\in F \text{ and } x\neq 0 \\ power_{F,I}^*(x,y)< fminD_F/2 & \text{if } x\in F \text{ and } x>0 \text{ and } y\in I \text{ and } x^y< fminD_F/3 \\ power_{F,I}^*(x,y)=power_{F,I}^*(-x,y) & \text{if } x\in F \text{ and } x<0 \text{ and } y\in I \text{ and } 2|y \\ power_{F,I}^*(x,y)=-power_{F,I}^*(-x,y) & \text{if } x\in F \text{ and } x<0 \text{ and } y\in I \text{ and not } 2|y \end{array}
```

The relationship to other $power_{FI'}^*$ helper functions for any $power_{FI'}$ operations in the same library shall be:

$$power_{FI}^*(x,y) = power_{FI'}^*(x,y)$$
 if $x \in F$ and $y \in I \cap I'$

The $power_{FI}$ operation:

```
power_{FI}: F \times I \rightarrow F \cup \{ underflow, overflow, infinitary \}
power_{FI}(x,y) = result_F^*(power_{FI}^*(x,y), nearest_F)
if x \in F and x \neq 0 and y \in I
                                                 if x = -\infty and y \in I and y > 0 and 2|y
                   =+\infty
                                                 if x = -\infty and y \in I and y > 0 and not 2|y
                   = -\infty
                                                 if x = -0 and y \in I and y > 0 and 2|y
                   = 0
                                                 if x = -0 and y \in Y and y > 0 and not 2|y|
                   = -0
                                                 if x = 0 and y \in I and y > 0
                   = 0
                                                 if x = +\infty and y \in I and y > 0
                   =+\infty
                                                 if x \in \{-\infty, -0, 0, +\infty\} and y = 0
                   =1
                   = 0
                                                 if x = -\infty and y \in I and y < 0 and 2|y|
                                                if x = -\infty and y \in I and y < 0 and not 2|y|
                   = -0
                   = infinitary(+\infty)
                                                 if x = -\mathbf{0} and y \in I and y < 0 and 2|y
                   = infinitary(-\infty)
                                                 if x = -\mathbf{0} and y \in I and y < 0 and not 2|y
                                                 if x = 0 and y \in I and y < 0
                                                 if x = +\infty and y \in I and y < 0
                                                 otherwise
```

NOTES

- 1 $power_{FI}(x,y)$ will overflow approximately when $x^y > fmax_F$, i.e., if x > 1, approximately when $y > \log_x(fmax_F)$, and if 0 < x < 1, approximately when $y < \log_x(fmax_F)$ (which is then negative). It will not overflow when x = 0 or when x = 1.
- 2 power_I (in clause 5.1.4) does not allow most negative exponents (unless |x| = 1) since the exact result then is not in \mathcal{Z} unless |x| = 1. power_F (in clause 5.3.6.6) does not allow any negative bases since the (exact) result is not in \mathcal{R} unless the exponent is integer. power_{FI} takes care of this latter case, where all exponents are ensured to be integers that have not arisen from implicit floating point rounding.

5.3.6.2 Natural exponentiation

The exp_F^* approximation helper function:

$$exp_F^*: F \to \mathcal{R}$$

 $exp_F^*(x)$ returns a close approximation to e^x in \mathcal{R} , with maximum error $max_error_exp_F$.

Further requirements on the exp_F^* approximation helper function are:

$$\begin{aligned} exp_F^*(1) &= e \\ exp_F^*(x) &= 1 \end{aligned} & \text{if } x \in F \text{ and } exp_F^*(x) \neq e^x \text{ and } \\ & \ln(1 - (epsilon_F/(2 \cdot r_F))) < x \text{ and } \\ & x < \ln(1 + (epsilon_F/2)) \end{aligned}$$

$$exp_F^*(x) < fminD_F/2$$
 if $x \in F$ and $x < \ln(fminD_F) - 3$

The exp_F operation:

 $exp_F: F \to F \cup \{\mathbf{underflow}, \mathbf{overflow}\}$

$$exp_F: F o F \cup \{ \text{underflow}, \text{overflow} \}$$
 $exp_F(x) = result_F^*(exp_F^*(x), nearest_F)$
 $= 1 \qquad \text{if } x \in F$
 $= 1 \qquad \text{if } x = -\mathbf{0}$
 $= +\infty \qquad \text{if } x = -\infty$
 $= no_result_F(x) \qquad \text{otherwise}$

NOTES

1 $exp_F(1) = nearest_F(e)$.
2 $exp_F(x)$ will overflow approximately when $x > \ln(fmax_F)$.

6.3 Natural exponentiation, minus one
 $expm1_F^*$ approximation helper function:
 $expm1_F^* : F \to \mathcal{R}$
 $m1_F^*(x)$ returns a close approximation to $e^{x} - 1$ in \mathcal{R} , with maximum error max_error_exp are their requirements on the $expm1_F^*$ approximation helper function are:
 $expm1_F^*(1) = e - 1$
 $expm1_F^*(1) = e - 1$
 $expm1_F^*(1) = x$
 $expm1_F^*(1) = x$

5.3.6.3

The $expm1_F^*$ approximation helper function:

$$expm1_F^*: F \to \mathcal{R}$$

 $expm1_F^*(x)$ returns a close approximation to e^x-1 in \mathcal{R} , with maximum error $max_error_exp_F$.

Further requirements on the $expm1_F^*$ approximation helper function are:

$$expm1_F^*(1) = e - 1$$
 $expm1_F^*(x) = x$
if $x \in F$ and $expm1_F^*(x) \neq e^x - 1$ and $-epsilon_F/r_F \leqslant x < 0.5 \cdot epsilon_F/r_F$
if $x \in F$ and $expm1_F^*(x) \neq e^x - 1$ and $x \in F$ and $expm1_F^*(x) \neq e^x - 1$ and $x \in F$ and $expm1_F^*(x) \neq e^x - 1$ and $x \in F$ and x

The relationship to the exp_F^* approximation helper function for the exp_F operation in the same library shall be:

$$expm1_F^*(x) \leqslant exp_F^*(x)$$
 if $x \in F$

The $expmI_F$ operation:

$$expm1_F: F \to F \cup \{\mathbf{overflow}\}$$

$$expm1_F(x) = result_F^*(expm1_F^*(x), nearest_F)$$

$$= x \qquad \qquad \text{if } x \in F \text{ and } |x| \geqslant fminN_F$$

$$= x \qquad \qquad \text{if } x \in F \text{ and } |x| < fminN_F$$

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

$$= +\infty \qquad \qquad \text{if } x = +\infty$$

$$= -1 \qquad \qquad \text{if } x = -\infty$$

$$= no_result_F(x) \qquad \text{otherwise}$$

NOTES

- 1 underflow is explicitly avoided. Part 1 requires that $fminN_F \leq epsilon_F$. This part requires that $fminN_F < 0.5 \cdot epsilon_F/r_F$, so that underflow can be avoided here.
- $2 \quad expm1_F(1) = nearest_F(e-1).$
- 3 $expm1_F(x)$ will overflow approximately when $x > \ln(fmax_F)$.

5.3.6.4 Exponentiation of 2

The $exp2_F^*$ approximation helper function:

$$exp2_F^*: F \to \mathcal{R}$$

 $exp\mathcal{Z}_F^*(x)$ returns a close approximation to 2^x in \mathcal{R} , with maximum error $max_error_exp_F$.

Further requirements on the $exp\mathcal{Z}_F^*$ approximation helper function are:

$$exp2_F^*(x) = 1$$
 if $x \in F$ and $exp2_F^*(x) \neq 2^x$ and
$$\log_2(1 - (epsilon_F/(2 \cdot r_F))) < x \text{ and}$$

$$x < \log_2(1 + (epsilon_F/2))$$

$$exp2_F^*(x) = 2^x$$
 if $x \in F \cap \mathcal{Z}$ and $2^x \in F$
$$exp2_F^*(x) < fminD_F/2$$
 if $x \in F$ and $x < \log_2(fminD_F) - 3$

The $exp2_F$ operation:

$$\mathit{exp2}_F: F \to F \cup \{\mathbf{underflow}, \mathbf{overflow}\}$$

$$exp2_F(x)$$
 = $result_F^*(exp2_F^*(x), nearest_F)$ if $x \in F$
= 1 if $x = -0$
= $+\infty$ if $x = +\infty$
= 0 if $x = -\infty$
= $no_result_F(x)$ otherwise

NOTE – $exp2_F(x)$ will overflow approximately when $x > \log_2(fmax_F)$.

5.3.6.5 Exponentiation of 10

The $exp10_F^*$ approximation helper function:

$$exp10_F^*: F \to \mathcal{R}$$

 $exp10_F^*(x)$ returns a close approximation to 10^x in \mathcal{R} , with maximum error $max_error_exp_F$.

Further requirements on the $exp10_F^*$ approximation helper function are:

$$\exp \mathcal{O}_F^*(x) = 1 \qquad \qquad \text{if } x \in F \text{ and } \exp \mathcal{O}_F^*(x) \neq 10^x \text{ and } \log_{10}(1 - (epsilon_F/(2 \cdot r_F))) < x \text{ and } x < \log_{10}(1 + (epsilon_F/2))$$

$$\exp \mathcal{O}_F^*(x) = 10^x \qquad \qquad \text{if } x \in F \cap \mathcal{Z} \text{ and } 10^x \in F$$

$$\exp \mathcal{O}_F^*(x) < fminD_F/2 \qquad \qquad \text{if } x \in F \text{ and } x < \log_{10}(fminD_F) - 3$$

The $exp10_F$ operation:

```
exp10_F: F \to F \cup \{\mathbf{underflow}, \mathbf{overflow}\}
```

```
exp10_F(x) = result_F^*(exp10_F^*(x), nearest_F) if x \in F if x = -\mathbf{0} if x = +\infty if x = -\infty if x = -\infty otherwise
```

NOTE – $exp10_F(x)$ will overflow approximately when $x > \log_{10}(fmax_F)$.

5.3.6.6 Exponentiation of argument base

The $power_F^*$ approximation helper function:

$$power_F^*: F \times F \to \mathcal{R}$$

 $power_F^*(x,y)$ returns a close approximation to x^y in \mathcal{R} , with maximum error $max_error_power_F$. The $power_F^*$ helper function need be defined only for first arguments that are greater than 0.

Further requirements on the $power_F^*$ approximation helper function are:

```
\begin{array}{ll} power_F^*(1,y) = 1 & \text{if } y \in F \\ power_F^*(x,0) = 1 & \text{if } x \in F \text{ and } x > 0 \\ power_F^*(x,1) = x & \text{if } x \in F \text{ and } x > 0 \\ power_F^*(x,y) < fminD_F/2 & \text{if } x \in F \text{ and } x > 0 \text{ and } y \in F \text{ and } x^y < fminD_F/3 \end{array}
```

The relationship to the $power_{FI}^*$ approximation helper functions for any $power_{FI}$ operations in the same library shall be:

```
if x \in F and x > 0 and y \in I \cap F
   power_F^*(x, y) = power_{F,I}^*(x, y)
The power_F operation:
   power_F: F \times F \to F \cup \{underflow, overflow, infinitary, invalid\}
                        = result_F^*(power_F^*(x, y), nearest_F)
                                                         if x \in F and x > 0 and y \in F
                                                         if x = -\mathbf{0} and y \in F \cup \{-\infty, -\mathbf{0}, +\infty\}
                                                         if y = -\mathbf{0} and x \in F \cup \{-\infty, +\infty\}
                                                         if x = +\infty and ((y \in F \text{ and } y > 0) \text{ or } y = +\infty)
                                                         if x \in F and x > 1 and y = +\infty
                                                         if x \in F and 0 \le x < 1 and y = +\infty
                                                         if x = 0 and y \in F and y > 0
                        = infinitary(+\infty)
                                                         if x = 0 and y \in F and y < 0
                                                         if x \in F and 0 \le x < 1 and y = -\infty
                                                         if x \in F and x > 1 and y = -\infty
                                                         if x = +\infty and ((y \in F \text{ and } y < 0) \text{ or } y = -\infty)
```

NOTE – $power_F(x, y)$ will overflow approximately when $x^y > fmax_F$, i.e., if x > 1, approximately when $y > \log_x(fmax_F)$, and if 0 < x < 1, approximately when $y < \log_x(fmax_F)$ (which is a negative number). It will not overflow when x = 0 or when x = 1.

otherwise

 $= no_result2_F(x,y)$

5.3.6.7 Exponentiation of one plus the argument base, minus one

The $power1pm1_F^*$ approximation helper function:

```
power1pm1_F^*: F \times F \to \mathcal{R}
```

 $power1pm1_F^*(x,y)$ returns a close approximation to $(1+x)^y - 1$ in \mathcal{R} , with maximum error $max_error_power_F$. The $power1pm1_F^*$ helper function need be defined only for first arguments that are greater than or equal to -1.

Further requirements on the $power1pm1_F^*$ approximation helper function are:

```
\begin{array}{ll} power1pm1_F^*(x,y) = (1+x)^y - 1 & \text{if } x,y \in F \cap \mathcal{Z} \text{ and } x \geqslant -1 \text{ and } y > 0 \\ power1pm1_F^*(x,1) = x & \text{if } x,1+x \in F \text{ and } x > -1 \\ power1pm1_F^*(-1,y) = -1 & \text{if } y \in F \text{ and } y > 0 \\ power1pm1_F^*(x,y) = -1 & \text{if } x \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \\ power1pm1_F^*(x,y) \neq (1+x)^y & \text{1 and } x \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } y \in F \text{ and } x > -1 \text{ and } x = -1 \text{ and
```

The relationship to the $power_F^*$ approximation helper function for the $power_F$ operation in the same library shall be:

```
power1pm1_F^*(x,y) \leq power_F^*(1+x,y) if x, 1+x \in F and x = -1 and y \in F

NOTE 1 - power1pm1_F^*(x,y) \approx y \cdot \ln(1+x) if x \in F and x > -1 and y \in F and |y| \cdot \ln(1+x) = -1 and |y| \cdot \ln
```

The $power1pm1_F$ operation:

```
power1pm1_F: F \times F \to F \cup \{-0, underflow, overflow, infinitary, invalid\}
power1pm1_F(x,y)
                   = result_F^*(power1pm1_F^*(x, y), nearest_F)
                                                 \bigcirc if x \in F and x > -1 and x \neq 0 and y \in F and y \neq 0
                                                   if x \in \{-0, 0\} and y \in F and y \neq 0
                                                   if y \in \{-0, 0\} and x \in F and x > -1
                   = mul_F(x, y)
                                                   if x = +\infty and ((y \in F \text{ and } y > 0) \text{ or } y = +\infty)
                                                   if x \in F and x > 0 and y = +\infty
                                                   if x \in F and -1 \leqslant x < 0 and y = +\infty
                                                   if x = -1 and y \in F and y > 0
                    = infinitary(+\infty)
                                                   if x = -1 and y \in F and y < 0
                                                   if x \in F and -1 \leqslant x < 0 and y = -\infty
                                                   if x \in F and x > 0 and y = -\infty
                                                   if x = +\infty and ((y \in F \text{ and } y < 0) \text{ or } y = -\infty)
                    = no\_result2_F(x, y)
                                                   otherwise
```

NOTE 2 – $power1pm1_F(x,y)$ will overflow approximately when $(1+x)^y > fmax_F$, i.e., if x>0, approximately when $y>\log_{1+x}(fmax_F)$, and if -1< x<0, approximately when $y<\log_{1+x}(fmax_F)$. It will not overflow when $x\in\{-1,0\}$.

5.3.6.8 Natural logarithm

The ln_F^* approximation helper function:

$$ln_F^*: F \cup \{e\} \to \mathcal{R}$$

 $ln_F^*(x)$ returns a close approximation to ln(x) in \mathcal{R} , with maximum error $max_error_exp_F$.

A further requirement on the ln_F^* approximation helper function is:

$$ln_F^*(e) = 1$$

The ln_F operation:

 $ln_F: F \to F \cup \{\mathbf{infinitary}, \mathbf{invalid}\}$

$$ln_F(x)$$
 = $result_F^*(ln_F^*(x), nearest_F)$
= $infinitary(-\infty)$ if $x \in F$ and $x > 0$
= $+\infty$ if $x \in \{-0, 0\}$
= $+\infty$ otherwise

5.3.6.9 Natural logarithm of one plus the argument

The $ln1p_F^*$ approximation helper function:

$$ln1p_F^*: F \cup \{e-1\} \to \mathcal{R}$$

 $ln1p_F^*(x)$ returns a close approximation to ln(1+x) in \mathcal{R} , with maximum error $max_error_exp_F$.

Further requirements on the $ln1p_F^*$ approximation helper function are:

$$\begin{aligned} \ln & 1p_F^*(e-1) = 1 \\ & \ln & 1p_F^*(x) = x \end{aligned} \qquad \text{if } x \in F \text{ and } \ln & 1p_F^*(x) \neq \ln(1+x) \text{ and} \\ & -0.5 \cdot epsilon_F/r_F < x \leqslant epsilon_F/r_F \end{aligned}$$

The relationship to the ln_F^* approximation helper function for the ln_F operation in the same library shall be:

$$ln1p_F^*(x) \geqslant ln_F^*(x)$$
 if $x \in F$ and $x > 0$

The $ln1p_F$ operation:

$$ln1p_F: F \rightarrow F \cup \{ \mathbf{infinitary}, \mathbf{invalid} \}$$

$$\begin{array}{ll} ln1p_F(x) & = result_F^*(ln1p_F^\bullet(x), nearest_F) \\ & \qquad \qquad & \text{if } x \in F \text{ and } x > -1 \text{ and } |x| \geqslant fminN_F \\ & = x \qquad \qquad & \text{if } x \in F \text{ and } |x| < fminN_F \\ & = -\mathbf{0} \qquad \qquad & \text{if } x = -\mathbf{0} \\ & = \inf \text{initiary}(-\infty) \qquad & \text{if } x = -1 \\ & = +\infty \qquad \qquad & \text{if } x = +\infty \\ & = no_result_F(x) \qquad & \text{otherwise} \end{array}$$

NOTE underflow is explicitly avoided. Part 1 requires that $fminN_F \leq epsilon_F$. This part requires that $fminN_F < 0.5 \cdot epsilon_F/r_F$, so that underflow can be avoided here.

5.3.6.10 2-logarithm

The $log 2_F^*$ approximation helper function:

$$log2_F^*: F \to \mathcal{R}$$

 $log \mathcal{Z}_F^*(x)$ returns a close approximation to $log_2(x)$ in \mathcal{R} , with maximum error $max_error_exp_F$.

A further requirement on the $log 2_F^*$ approximation helper function is:

```
log 2_F^*(x) = log_2(x) if x \in F and log_2(x) \in \mathcal{Z}

The log 2_F operation:

log 2_F: F \to F \cup \{ \text{infinitary}, \text{invalid} \}

log 2_F(x) = result_F^*(log 2_F^*(x), nearest_F)

if x \in F and x > 0

= \text{infinitary}(-\infty) if x \in \{-\mathbf{0}, 0\}

= +\infty if x = +\infty

= no\_result_F(x) otherwise
```

5.3.6.11 10-logarithm

The $log10_F^*$ approximation helper function:

$$log10_F^*: F \to \mathcal{R}$$

 $log10_F^*(x)$ returns a close approximation to $log_{10}(x)$ in \mathcal{R} , with maximum error $max_error_exp_F$.

A further requirement on the $log10_F^*$ approximation helper function is

$$log 10_F^*(x) = log_{10}(x)$$
 if $x \in F$ and $log_{10}(x) \in \mathcal{Z}$

The $log10_F$ operation:

$$\begin{split} log10_F : F \to F \cup & \{ \mathbf{infinitary}, \mathbf{invalid} \} \\ log10_F(x) &= result_F^*(log10_F^*(x), nearest_F) \\ &= \mathbf{infinitary}(-\infty) & \text{if } x \in F \text{ and } x > 0 \\ &= +\infty & \text{if } x = +\infty \\ &= no_result_F(x) & \text{otherwise} \end{split}$$

5.3.6.12 Argument base logarithm

The $logbase_F^*$ approximation helper function:

$$logbase_F^*: F \times F \to \mathcal{R}$$

 $logbase_F^*(x,y)$ returns a close approximation to $log_x(y)$ in \mathcal{R} , with maximum error $max_error_power_F$.

A further requirement on the $logbase_F^*$ approximation helper function is:

$$logbase_F^*(x, x) = 1$$
 if $x \in F$ and $x > 0$ and $x \neq 1$

The *logbase* poperation:

$$\begin{aligned} logbase_F : F \times F \to F \cup \{-\mathbf{0}, \mathbf{infinitary}, \mathbf{invalid}\} \\ logbase_F (x,y) &= result_F^*(logbase_F^*(x,y), nearest_F) \\ &= infinitary(+\infty) \\ &= infinitary(-\infty) \\ &= 0 \end{aligned} \quad \begin{aligned} &\text{if } x \in F \text{ and } x > 0 \text{ and } x \neq 1 \text{ and } y \in F \text{ and } y > 0 \\ &= logbase_F(0,y) \\ &= infinitary(+\infty) \\ &= \mathbf{infinitary}(+\infty) \end{aligned} \quad \begin{aligned} &\text{if } x = -\mathbf{0} \text{ and } x \neq 1 \text{ and } y \in F \text{ and } y > 0 \\ &= \mathbf{0} \end{aligned} \quad \end{aligned}$$

if $x \in F$ and $1 \leq x$ and $y = +\infty$

 $=+\infty$

```
= -\infty
= -0
= 0
if x \in F and 0 < x < 1 and y = +\infty
if x = 0 and y \in F and y \ge 1
= 0
if x = 0 and y \in F and 0 < y < 1
= infinitary(+\infty)
if x \in F and 0 < x < 1 and y = 0
= infinitary(-\infty)
if x \in F and 1 < x and y = 0
= -0
if x = +\infty and y \in F and 0 < y < 1
= no\_result2_F(x, y)
otherwise
```

5.3.6.13 Argument base logarithm of one plus each argument

The $logbase1p1p_F^*$ approximation helper function:

 $logbase1p1p_F^*: F \times F \to \mathcal{R}$

 $logbase1p1p_F^*(x,y)$ returns a close approximation to $log_{(1+x)}(1+y)$ in \mathcal{R} , with maximum error $max_error_power_F$.

A further requirements on $logbase1p1p_F^*$ approximation helper function is

$$logbase1p1p_F^*(x,x) = 1$$
 if $x \in F$ and $x > -1$ and $x \neq 0$

The $logbase1p1p_F$ operation:

 $logbase1p1p_F: F \times F \rightarrow F \cup \{ \mathbf{-0}, \mathbf{underflow}, \mathbf{infinitary}, \mathbf{invalid} \}$ $logbase1p1p_F(x,y)$ $= result_F^*(logbase1p1p_F^*(x, y), nearest_F)$ if $x \in F$ and x > -1 and $x \neq 0$ and $y \in F$ and y > -1 and $y \neq 0$ if $x \in \{-0, 0\}$ and $((y \in F \text{ and } y > -1 \text{ and } y \neq 0) \text{ or } y = +\infty)$ if $y \in \{-\mathbf{0}, 0\}$ and $((x \in F \text{ and } x > -1) \text{ or } x = +\infty)$ if $x = +\infty$ and $y \in F$ and $y \geqslant 0$ if $x \in F$ and 0 < x and $y = +\infty$ if $x \in F$ and -1 < x < 0 and $y = +\infty$ if x = -1 and $y \in F$ and $y \ge 0$ if x = -1 and $y \in F$ and -1 < y < 0if $x \in F$ and -1 < x < 0 and y = -1= infinitary $(+\infty)$ $= \mathbf{infinitary}(-\infty)$ if $x \in F$ and 0 < x and y = -1if $x = +\infty$ and $y \in F$ and -1 < y < 0

5.3.7 Introduction to operations for trigonometric elementary functions

Two different operations for each of sin, cos, tan, cot, sec, csc, arcsin, arccos, arctan, arccot, arccotc, arcsec, and arccsc are specified. One version for radians and one version where the angular unit is given as a parameter.

otherwise

For use in the specifications below, define the following mathematical functions:

$$rad: \mathcal{R} \to \mathcal{R}$$

 $axis_rad: \mathcal{R} \to \{(1,0), (0,1), (-1,0), (0,-1)\} \times \mathcal{R}$
 $arc: \mathcal{R} \times \mathcal{R} \to \mathcal{R}$

The rad, angular value normalisation, function is defined by

$$rad(x)$$
 = $x - round(x/(2 \cdot \pi)) \cdot 2 \cdot \pi$

The axis_rad function is defined by

$$axis_rad(x)$$
 = $((1,0), \arcsin(\sin(x)))$ if $\cos(x) \ge 1/\sqrt{2}$
 = $((0,1), \arcsin(\cos(x)))$ if $\sin(x) > 1/\sqrt{2}$
 = $((-1,0), \arcsin(\sin(x)))$ if $\cos(x) \le -1/\sqrt{2}$
 = $((0,-1), \arcsin(\cos(x)))$ if $\sin(x) < -1/\sqrt{2}$

The arc, angle, function is defined by

the
$$arc$$
, angle, function is defined by $arc(x,y) = -\arccos(x/\sqrt{x^2+y^2})$ if $y < 0$ $= \arccos(x/\sqrt{x^2+y^2})$ if $y \geqslant 0$

Operations for radian trigonometric elementary functions a shall be one radian big-angle parameter: $big_angle_r_F \in F$ and have the following default value:

There shall be one radian big-angle parameter:

$$big_angle_r_F \in F$$

It should have the following default value:

$$big_angle_r_F = r_F^{\lceil p_F/2 \rceil}$$

A binding or implementation can include a method to change the value of the radian big-angle parameter. This method should only allow the value of this parameter to be set to a value greater than $2 \cdot \pi$ and such that $ulp_F(big_angle_r_F) < \pi/1000$.

NOTE – Part 1 requires only that $p_F \ge 2$, but see also A.5.2.0.2 in part 1. This part requires that $p_F \ge 2 \cdot \max\{1, \lceil \log_{r_F}(2 \cdot \pi) \rceil\}$, in order to allow at least the first two cycles (plus and minus) to be in the interval $[-big_angle_r_F]$ big_angle_r_F]. In order to allow $ulp_F(big_angle_r_F) <$ $\pi/1000, p_F \ge 2 + \lceil \log_{r_F}(1000) \rceil$ should hold.

For use in the approximation belper function's signatures, define

$$F^{2\cdot\pi} = (F \cup \{n \cdot \pi/4, n - \pi/6 \mid n \in \mathcal{Z}\}) \cap [-big_angle_r_F, big_angle_r_F]$$

There shall be three maximum error parameters for radian trigonometric operations:

```
max\_error\_rad_F \in F
max\_error\_sin_F \in F
max\_error\_tan_F \in F
```

The $max_error_rad_F$ parameter shall have a value that is 0.5 (ulp). The $max_error_sin_F$ parameter shall have a value that is $\leq 1.5 \cdot rnd_error_F$. The $max_error_tan_F$ parameter shall have a value that is $2 \cdot rnd_error_F$. If the binding standard requires that the $max_error_rad_F$ parameter has the value 0.5, that parameter need not be made available for programs.

5.3.8.1 Radian angle normalisation

The rad_F^* approximation helper function:

$$rad_F^*: F^{2\cdot \pi} \to \mathcal{R}$$

 $rad_F^*(x)$ returns a close approximation to rad(x) in \mathcal{R} , if $|x| \leq big_angle_r_F$, with maximum error $max_error_rad_F$.

The $axis_rad_F^*$ approximation helper function:

$$axis_rad_F^*: F^{2\cdot\pi} \to \{(1,0), (0,1), (-1,0), (0,-1)\} \times \mathcal{R}$$

 $axis_rad_F^*(x)$ returns a close approximation to $axis_rad(x)$, if $|x| \leq big_angle_r_F$, with maximum error max_error_rad_F for the second part of the result. The approximation consists of that the second part of the result (the offset from the indicated axis) is approximate. The first part (the nearest axis indication) shall be exact if $|x| \leq big_angle_r_F$.

NOTE 1 - With the maximum error 0.5 ulp, these helper functions are not really needed. However, Annex A allows for partial conformity, such that the maximum error for these two helper functions may be greater than 0.5 ulp.

Further requirements on the rad_F^* and $axis_rad_F^*$ approximation helper functions are:

$$rad_F^*(x) = x$$
 if $x \in F^{2 \cdot \pi}$ and $|x| = x$ $snd(axis_rad_F^*(x)) = rad_F^*(x)$ if $x \in F^{2 \cdot \pi}$ and $fst(axis_rad_F^*(x)) = (1,0)$ are rad_F operation:

The rad_F operation:

 $rad_F: F \to F \cup \{\text{underflow}, \text{absolute_precision_underflow}\}$

$$rad_{F}(x)$$
 = $result_{F}^{*}(rad_{F}^{*}(x), nearest_{F})$ if $x \in F$ and $|x| > fminN_{F}$ and $|x| \leqslant big_angle_r_{F}$ if $x \in F$ and $|x| \leqslant fminN_{F}$) or $x = -\mathbf{0}$

= absolute_precision_underflow(qNaN)
if
$$x \in F$$
 and $|x| > big_angle_r_F$ otherwise

The $axis_rad_F$ operation:

$$axis_rad_F: F \rightarrow (\{(1,0),(0,1),(-1,0),(0,-1)\} \times F) \cup \{absolute_precision_underflow\}$$

$$axis_rad_F(x) = (fst(axis_rad_F^*(x)), result_F^*(snd(axis_rad_F^*(x)), nearest_F))$$

$$if \ x \in F \ \text{and} \ |x| > fminN_F \ \text{and} \ |x| \leqslant big_angle_r_F$$

$$= ((1,0),x) \qquad if \ (x \in F \ \text{and} \ |x| > big_angle_r_F$$

$$= ((qNaN, qNaN), qNaN)$$

$$if \ x \in F \ \text{and} \ |x| > big_angle_r_F$$

$$= ((qNaN, qNaN), qNaN)$$

if x is a quiet NaN = invalid((qNaN, qNaN), qNaN)

otherwise

NOTE 2 - rad_F is simpler, easier to use, but less accurate than $axis_rad_F$. The latter may still not be sufficient for implementing the radian trigonometric operations to less than the maximum error stated by the parameters.

5.3.8.2 Radian sine

The sin_F^* approximation helper function:

$$sin_F^*: F^{2\cdot \pi} \to \mathcal{R}$$

 $sin_F^*(x)$ returns a close approximation to sin(x) in \mathcal{R} if $|x| \leq big_angle_r_F$, with maximum error $max_error_sin_F$.

Further requirements on the sin_F^* approximation helper function are:

$$\begin{array}{ll} sin_F^*(n \cdot 2 \cdot \pi + \pi/6) = 1/2 & \text{if } n \in \mathcal{Z} \text{ and } |n \cdot 2 \cdot \pi + \pi/6| \leqslant big_angle_r_F \\ sin_F^*(n \cdot 2 \cdot \pi + \pi/2) = 1 & \text{if } n \in \mathcal{Z} \text{ and } |n \cdot 2 \cdot \pi + \pi/2| \leqslant big_angle_r_F \\ sin_F^*(n \cdot 2 \cdot \pi + 5 \cdot \pi/6) = 1/2 & \text{if } n \in \mathcal{Z} \text{ and } |n \cdot 2 \cdot \pi + \pi/2| \leqslant big_angle_r_F \\ sin_F^*(x) = x & \text{if } x \in F^{2 \cdot \pi} \text{ and } sin_F^*(x) \neq sin(x) \text{ and } \\ |x| \leqslant \sqrt{3 \cdot epsilon_F/r_F} \\ sin_F^*(-x) = -sin_F^*(x) & \text{if } x \in F^{2 \cdot \pi} \end{array}$$

The sin_F operation:

 $sin_F: F \rightarrow F \cup \{ \mathbf{underflow}, \mathbf{absolute_precision_underflow} \}$

$$sin_F(x)$$
 = $result_F^*(sin_F^*(x), nearest_F)$
if $x \in F$ and $fmin_F < |x|$ and $|x| \le big_angle_r_F$
= $rad_F(x)$ otherwise

NOTE – **underflow** is here explicitly avoided for subnormal arguments, but the operation may **underflow** for other arguments.

5.3.8.3 Radian cosine

The cos_F^* approximation helper function:

$$cos_F^*: F^{2\cdot \pi} \to \mathcal{R}$$

 $cos_F^*(x)$ returns a close approximation to cos(x) in \mathcal{R} if $|x| \leq big_angle_r_F$, with maximum error $max_error_sin_F$.

Further requirements on the \cos_F^* approximation helper function are:

$$\begin{array}{lll} \cos_F^*(n\cdot 2\cdot \pi)=1 & \text{if } n\in\mathcal{Z} \text{ and } |n\cdot 2\cdot \pi|\leqslant big_angle_r_F\\ \cos_F^*(n\cdot 2\cdot \pi+\pi/3)=1/2 & \text{if } n\in\mathcal{Z} \text{ and } |n\cdot 2\cdot \pi+\pi/3|\leqslant big_angle_r_F\\ \cos_F^*(n\cdot 2\cdot \pi+2\cdot \pi/3)=-1/2 & \text{if } n\in\mathcal{Z} \text{ and } |n\cdot 2\cdot \pi+\pi/3|\leqslant big_angle_r_F\\ \cos_F^*(n\cdot 2\cdot \pi+\pi)=-1 & \text{if } n\in\mathcal{Z} \text{ and } |n\cdot 2\cdot \pi+2\cdot \pi/3|\leqslant big_angle_r_F\\ \cos_F^*(x)=1 & \text{if } n\in\mathcal{Z} \text{ and } |n\cdot 2\cdot \pi+\pi|\leqslant big_angle_r_F\\ \cos_F^*(x)=1 & \text{if } x\in F^{2\cdot \pi} \text{ and } \cos_F^*(x)\neq \cos(x) \text{ and } \\ |x|<\sqrt{epsilon_F/r_F}\\ \cos_F^*(x)=1 & \text{if } x\in F^{2\cdot \pi} \end{array}$$

The cos_F operation:

 $cos_F: F \to F \cup \{underflow, absolute_precision_underflow\}$

$$cos_F(x)$$
 = $result_F^*(cos_F^*(x), nearest_F)$
if $x \in F$ and $|x| \leq big_angle_r_F$
= 1 if $x = -\mathbf{0}$
= $rad_F(x)$ otherwise

5.3.8.4 Radian tangent

The tan_F^* approximation helper function:

$$tan_F^*: F^{2\cdot \pi} \to \mathcal{R}$$

 $tan_F^*(x)$ returns a close approximation to tan(x) in \mathcal{R} if $|x| \leq big_angle_r_F$, with maximum error $max_error_tan_F$.

Further requirements on the tan_F^* approximation helper function are:

$$tan_F^*(n \cdot 2 \cdot \pi + \pi/4) = 1$$
 if $n \in \mathbb{Z}$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$
$$tan_F^*(n \cdot 2 \cdot \pi + 3 \cdot \pi/4) = -1$$
 if $n \in \mathbb{Z}$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ if $n \in \mathbb{Z}$ and $|n \cdot 2 \cdot \pi + 3 \cdot \pi/4| \leqslant big_angle_r_F$ if $n \in \mathbb{Z}$ and $|n \cdot 2 \cdot \pi + 3 \cdot \pi/4| \leqslant big_angle_r_F$ if $n \in \mathbb{Z}$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ if $n \in \mathbb{Z}$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ if $n \in \mathbb{Z}$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ if $n \in \mathbb{Z}$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ if $n \in \mathbb{Z}$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ if $n \in \mathbb{Z}$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ if $n \in \mathbb{Z}$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ if $n \in \mathbb{Z}$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ if $n \in \mathbb{Z}$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ if $n \in \mathbb{Z}$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ if $n \in \mathbb{Z}$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ if $n \in \mathbb{Z}$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ if $n \in \mathbb{Z}$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ if $n \in \mathbb{Z}$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ if $n \in \mathbb{Z}$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ if $n \in \mathbb{Z}$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ if $n \in \mathbb{Z}$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ if $n \in \mathbb{Z}$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ if $n \in \mathbb{Z}$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ if $n \in \mathbb{Z}$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ if $n \in \mathbb{Z}$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ if $n \in \mathbb{Z}$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ if $n \in \mathbb{Z}$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ if $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ if $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle$

NOTE 1 – tan has a smallest period of π , but the above expresses a period of $2 \cdot \pi$, which is more in line with the other operations. The desired points of extra accuracy are still covered.

The tan_F operation:

$$tan_F: F \rightarrow F \cup \{ ext{underflow}, ext{overflow}, ext{absolute_precision_underflow} \}$$

$$tan_F(x)$$
 = $result_F^*(tan_F^*(x), nearest_F)$
if $x \in F$ and $fminN_F < |x|$ and $|x| \le big_angle_r_F$
= $rad_F(x)$ otherwise

NOTE 2 - **underflow** is explicitly avoided for subnormal arguments, but the operation may **underflow** for other arguments.

5.3.8.5 Radian cotangent

The cot_F^* approximation helper function:

$$cot_F^*: F^{2\cdot \pi} \to \mathcal{R}$$

 $cot_F^*(x)$ returns a close approximation to cot(x) in \mathcal{R} if $|x| \leq big_angle_r_F$, with maximum error $max_error_tan_F$.

Further requirements on the cot_F^* approximation helper function are:

$$\cot_F^*(n \cdot 2 \cdot \pi + \pi/4) = 1$$
 if $n \in \mathcal{Z}$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$
$$\cot_F^*(n \cdot 2 \cdot \pi + 3 \cdot \pi/4) = -1$$
 if $n \in \mathcal{Z}$ and $|n \cdot 2 \cdot \pi + \pi/4| \leqslant big_angle_r_F$
$$\cot_F^*(-x) = \cot_F^*(x)$$
 if $x \in F^{2 \cdot \pi}$

NOTE – cot has a smallest period of π , but the above expresses a period of $2 \cdot \pi$, which is more in line with the other operations. The desired points of extra accuracy are still covered.

The cot_E operation:

 $cot_F: F \to F \cup \{\text{underflow}, \text{overflow}, \text{infinitary}, \text{absolute_precision_underflow}\}$

$$cot_F(x)$$
 = $result_F^*(cot_F^*(x), nearest_F)$
if $x \in F$ and $x \neq 0$ and $|x| \leq big_angle_r_F$
= $infinitary(+\infty)$ if $x = 0$
= $infinitary(-\infty)$ if $x = -0$
= $rad_F(x)$ otherwise

5.3.8.6 Radian secant

The sec_F^* approximation helper function:

$$sec_F^*: F^{2\cdot \pi} \to \mathcal{R}$$

 $sec_F^*(x)$ returns a close approximation to sec(x) in \mathcal{R} if $|x| \leq big_angle_r_F$, with maximum error $max_error_tan_F$.

Further requirements on the sec_F^* approximation helper function are:

$$\begin{array}{lll} \sec c_F^*(n \cdot 2 \cdot \pi) = 1 & \text{if } n \in \mathcal{Z} \text{ and } |n \cdot 2 \cdot \pi| \leqslant big_angle_r_F \\ \sec c_F^*(n \cdot 2 \cdot \pi + \pi/3) = 2 & \text{if } n \in \mathcal{Z} \text{ and } |n \cdot 2 \cdot \pi + \pi/3| \leqslant big_angle_r_F \\ \sec c_F^*(n \cdot 2 \cdot \pi + 2 \cdot \pi/3) = -2 & \text{if } n \in \mathcal{Z} \text{ and } |n \cdot 2 \cdot \pi + 2 \cdot \pi/3| \leqslant big_angle_r_F \\ \sec c_F^*(n \cdot 2 \cdot \pi + \pi) = -1 & \text{if } n \in \mathcal{Z} \text{ and } |n \cdot 2 \cdot \pi + \pi| \leqslant big_angle_r_F \\ \sec c_F^*(x) = 1 & \text{if } x \in F^{2 \cdot \pi} \text{ and } \sec c_F^*(x) \neq \sec(x) \text{ and } \\ \sec c_F^*(-x) = \sec c_F^*(x) & \text{if } x \in F^{2 \cdot \pi} \end{array}$$

The sec_F operation:

 $sec_F: F \to F \cup \{ overflow, absolute_precision_underflow \}$

$$sec_F(x)$$
 = $result_F^*(sec_F^*(x), nearest_F)$ if $x \in F$ and $|x| \geqslant big_angle_r_F$ = 1 if $x = -\mathbf{0}$ otherwise

5.3.8.7 Radian cosecant

The csc_F^* approximation helper function:

$$csc_F^*: F^{2\cdot \pi} \to \mathcal{R}$$

 $csc_F^*(x)$ returns a close approximation $\operatorname{ccsc}(x)$ in \mathcal{R} if $|x| \leq big_angle_r_F$, with maximum error $max_error_tan_F$.

Further requirements on the csc_F^* approximation helper function are:

$$csc_F^*(n \cdot 2 \cdot \pi + \pi/6) = 2$$
 if $n \in \mathbb{Z}$ and $|n \cdot 2 \cdot \pi + \pi/6| \leqslant big_angle_r_F$
$$csc_F^*(n \cdot 2 \cdot \pi + \pi/2) = 1$$
 if $n \in \mathbb{Z}$ and $|n \cdot 2 \cdot \pi + \pi/2| \leqslant big_angle_r_F$
$$csc_F^*(n \cdot 2 \cdot \pi + 5 \cdot \pi/6) = 2$$
 if $n \in \mathbb{Z}$ and $|n \cdot 2 \cdot \pi + 5 \cdot \pi/6| \leqslant big_angle_r_F$
$$csc_F^*(-x) = -csc_F^*(x)$$
 if $x \in F^{2 \cdot \pi}$

The csc_F operation:

 $csc_F : \mathcal{F} \to F \cup \{\text{overflow}, \text{infinitary}, \text{absolute_precision_underflow}\}$

$$csc_F(x)$$
 = $result_F^*(csc_F^*(x), nearest_F)$
if $x \in F$ and $x \neq 0$ and $|x| \leq big_angle_r_F$
= $infinitary(+\infty)$ if $x = 0$
= $infinitary(-\infty)$ if $x = -0$
= $rad_F(x)$ otherwise

5.3.8.8 Radian cosine with sine

 $cossin_F: F \rightarrow (F \times F) \cup \{ underflow, absolute_precision_underflow \}$ $=(cos_F(x), sin_F(x))$

NOTES

- 1 If there is an absolute_precision_underflow notification, then both result parts suffer from the absolute_precision_underflow and the continuation values for both parts are qNaN. Similarly for NaN and infinitary arguments.
- 2 If there is an **underflow** notification, only one of the result parts suffer from the underflow, and the other part has an absolute value greater than $fminN_F$.

5.3.8.9 Radian arc sine

The $arcsin_F^*$ approximation helper function:

$$arcsin_F^*: F \to \mathcal{R}$$

 $arcsin_F^*(x)$ returns a close approximation to arcsin(x) in \mathcal{R} , with maximum error $max_error_sin_F$.

Further requirements on the $arcsin_F^*$ approximation helper function are:

$$arcsin_F^*(1/2) = \pi/6$$

 $arcsin_F^*(1) = \pi/2$
 $arcsin_F^*(x) = x$ if $x \in F$ and $arcsin_F^*(x) \neq arcsin(x)$ and $|x| < \sqrt{2 \cdot epsilon_F/r_F}$
 $arcsin_F^*(-x) = -arcsin_F^*(x)$ if $x \in F$
The $arcsin_F^\#$ range limitation helper function (for $x \in F$):
 $arcsin_F^\#(x) = \max\{un_F(-\pi/2) \min\{arcsin_F^*(x) \mid down_F(\pi/2)\}\}$

$$arcsin_F^{\#}(x) = \max\{up_F(-\pi/2), \min\{arcsin_F^{*}(x), down_F(\pi/2)\}\}$$

The $arcsin_F$ operation:

$$arcsin_F: F \to F \cup \{\mathbf{invalid}\}$$

$$arcsin_F(x) = result_F^*(arcsin_F^\#(x), nearest_F)$$

$$if \ x \in F \ \text{and} \ fminN_F < |x| \leqslant 1$$

$$= x \qquad \qquad \text{if} \ (x \in F \ \text{and} \ |x| \leqslant fminN_F) \ \text{or} \ x = -\mathbf{0}$$

$$= not result_F(x) \qquad \text{otherwise}$$

NOTE - underflow is explicitly avoided.

5.3.8.10 Radian arc cosine

The $arccos_F^*$ approximation helper function:

$$arccos_F^*: F \to \mathcal{R}$$

 $arccos_F^*(x)$ returns a close approximation to arccos(x) in \mathcal{R} , with maximum error $max_error_sin_F$.

Further requirements on the $arccos_F^*$ approximation helper function are:

```
arccos_F^*(1/2) = \pi/3
arccos_F^*(0) = \pi/2
arccos_F^*(-1/2) = 2 \cdot \pi/3
arccos_F^*(-1) = \pi
```

The $arccos_F^{\#}$ range limitation helper function (for $x \in F$):

$$arccos_F^{\#}(x) = \min\{arccos_F^*(x), down_F(\pi)\}$$

The $arccos_F$ operation:

```
arccos_F: F \to F \cup \{\mathbf{invalid}\}

arccos_F(x) = result_F^*(arccos_F^\#(x), nearest_F)

if \ x \in F \ \text{and} \ |x| \leqslant 1

= arccos_F(0) \qquad \text{if } x = -\mathbf{0}

= no\_result_F(x) \qquad \text{otherwise}
```

5.3.8.11 Radian arc tangent

The $arctan_F^*$ approximation helper function:

$$arctan_F^*: F \to \mathcal{R}$$

 $arctan_F^*(x)$ returns a close approximation to arctan(x) in \mathcal{R} , with maximum error $max_error_tan_F$.

Further requirements on the $arctan_F^*$ approximation helper function are:

$$arctan_F^*(1) = \pi/4$$

$$arctan_F^*(x) = x$$
 if $x \in F$ and $arctan_F^*(x) \neq \arctan(x)$ and
$$x \mid \leqslant \sqrt{1.5 \cdot epsilon_F/r_F}$$

$$arctan_F^*(x) = \pi/2$$
 if $x \in F$ and $arctan_F^*(x) \neq \arctan(x)$ and
$$x > 3 \cdot r_F/epsilon_F$$
 if $x \in F$

The $arctan_F^{\#}$ range limitation helper function (for $x \in F$):

$$arctan_F^{\#}(x) = \max\{up_F(-\pi/2), \min\{arctan_F^*(x), down_F(\pi/2)\}\}$$

The $arctan_F$ operation:

$$arctan_F: F \to F$$

$$arctan_F(x)$$
 $Presult_F^*(arctan_F^\#(x), nearest_F)$

if
$$x \in F$$
 and $fminN_F < |x|$
 $= x$ if $(x \in F \text{ and } |x| \leqslant fminN_F)$ or $x = -0$
 $= up_F(-\pi/2)$ if $x = -\infty$
 $= down_F(\pi/2)$ if $x = +\infty$
 $= no_result_F(x)$ otherwise

NOTES

- 1 $arctan_F(x) \approx arc_F(1,x)$. (arc_F is specified in subclause 5.3.8.15 below.)
- 2 **underflow** is explicitly avoided.

5.3.8.12 Radian arc cotangent

This clause specifies two inverse cotangent operations. One approximating the sign symmetric (but discontinuous at 0) arccot, the other approximating the continuous (but not sign symmetric) arccotc.

The $arccot_F^*$ approximation helper function:

$$arccot_F^*: F \to \mathcal{R}$$

 $arccot_F^*(x)$ returns a close approximation to arccot(x) in \mathcal{R} , with maximum error $max_error_tan_F$.

The $arccotc_F^*$ approximation helper function:

$$arccotc_F^*: F \to \mathcal{R}$$

 $arccotc_F^*(x)$ returns a close approximation to arccotc(x) in \mathcal{R} , with maximum error $max_e cror_tan_F$.

Further requirements on the $arccot_F^*$ and $arccot_F^*$ approximation helper functions are:

$$\begin{array}{ll} \operatorname{arccot}_F^*(1) = \pi/4 \\ \operatorname{arccot}_F^*(0) = \pi/2 \\ \operatorname{arccot}_F^*(-x) = -\operatorname{arccot}_F^*(x) & \text{if } x \in F \text{ and } x \neq 0 \\ \operatorname{arccotc}_F^*(x) = \operatorname{arccot}_F^*(x) & \text{if } x \in F \text{ and } x \neq 0 \\ \operatorname{arccotc}_F^*(-1) = 3 \cdot \pi/4 & \text{if } x \in F \text{ and } \operatorname{arccotc}_F^*(x) \neq \operatorname{arccotc}(x) \text{ and } x \in F \text{ and$$

The $arccot_F^{\#}$ and $arccot_F^{\#}$ range limitation helper functions (for $x \in F$):

$$arccot_F^{\#}(x) = \max\{up_F(-\pi/2), \min\{arccot_F^{*}(x), down_F(\pi/2)\}\}$$
$$arccote_F^{\#}(x) = \min\{arccote_F^{*}(x), down_F(\pi)\}$$

The $arccot_F$ operation:

$$arccot_F: F \to F \cup \{\mathbf{underflow}\}\$$

$$arccot_F(x)$$
 = $result_F^*(arccot_F^\#(x), nearest_F)$ if $x \in F$ if $x = -0$ if $x = -\infty$ if $x = +\infty$ otherwise

NOTES

- 1 $arccot_F(neg_F(x)) = neg_F(arccot_F(x)).$
- 2 Due to the range limitation, $arccot_F(0)$ need not equal $arccotc_F(0)$.

The $arccotc_F$ operation:

$$arccotc_F: F \to F \cup \{ underflow \}$$

$$arccotc_F(x) = result* (arccote$$

$$arccotc_F(x)$$
 = $result_F^*(arccotc_F^\#(x))$ if $x \in F$
 = $nearest_F(\pi/2)$ if $x = -\mathbf{0}$
 = $down_F(\pi)$ if $x = -\infty$
 = 0 if $x = +\infty$
 = $no_result_F(x)$ otherwise

NOTE 3 - $arccotc_F(x) \approx arc_F(x,1)$. (arc_F is specified in subclause 5.3.8.15 below.)

5.3.8.13 Radian arc secant

The $arcsec_F^*$ approximation helper function:

$$arcsec_F^*: F \to \mathcal{R}$$

 $arcsec_F^*(x)$ returns a close approximation to arcsec(x) in \mathcal{R} , with maximum error $max_error_tan_F$.

Further requirements on the $arcsec_F^*$ approximation helper function are:

```
arcsec_F^*(2) = \pi/3
arcsec_F^*(-2) = 2 \cdot \pi/3
arcsec_F^*(-1) = \pi
arcsec_F^*(x) \leqslant \pi/2
arcsec_F^*(x) \geqslant \pi/2
arcsec_F^*(x) \geqslant \pi/2
arcsec_F^*(x) = \pi/2
if x \in F and x > 0
if x \in F and x < 0
if x \in F and x < 0
arcsec_F^*(x) = \pi/2
if x \in F and x < 0
arcsec_F^*(x) \neq arcsec(x) \text{ and } x > 0
|x| > 3 \cdot r_F / epsilon_F
```

The $arcsec_F^{\#}$ range limitation helper function (for $x \in F$):

$$arcsec_F^{\#}(x) = \min\{arcsec_F^{*}(x), down_F(\pi/2)\}$$

$$if x \ge 1$$

$$= \max\{up_F(\pi/2), \min\{arcsec_F^{*}(x), down_F(\pi)\}\}$$

$$if x \le 1$$

The $arcsec_F$ operation:

$$arcsec_F: F \to F \cup \{\mathbf{invalid}\}$$

$$arcsec_F(x) = result_F^*(arcsec_F^\#(x), nearest_F)$$

$$if \ x \in F \ \text{and} \ 1 \leqslant |x|$$

$$= up_F(\pi/2) \qquad \text{if} \ x = -\infty$$

$$= down_F(\pi/2) \qquad \text{if} \ x = +\infty$$

$$= no_result_F(x) \qquad \text{otherwise}$$

5.3.8.14 Radian arc cosecant

The $arccsc_F^*$ approximation helper function:

$$arccsc_F^* : F \to \mathcal{R}$$

 $arccsc_F^*(x)$ returns a close approximation to arccsc(x) in \mathcal{R} , with maximum error $max_error_tan_F$.

Further requirements on the $arccsc_F^*$ approximation helper function are:

```
arccsc_F^*(2) = \pi/6

arccsc_F^*(1) = \pi/2

arccsc_F^*(-x) = -arccsc_F^*(x) if x \in F
```

The $arccsc_F^{\#}$ range limitation helper function (for $x \in F$):

$$\operatorname{arccsc}_F^\#(x) = \max\{\operatorname{up}_F(-\pi/2), \min\{\operatorname{arccsc}_F^*(x), \operatorname{down}_F(\pi/2)\}\}$$

The $arccsc_F$ operation:

```
arccsc_F : F \to F \cup \{\mathbf{underflow}, \mathbf{invalid}\}
```

```
arccsc_F(x) = result_F^*(arccsc_F^\#(x), nearest_F)

if x \in F and 1 \leq |x|

= -\mathbf{0} if x = -\infty

= 0 if x = +\infty

= no\_result_F(x) otherwise
```

5.3.8.15 Radian angle from Cartesian co-ordinates

The arc_F^* approximation helper function:

$$arc_F^*: F \times F \to \mathcal{R}$$

 $arc_F^*(x,y)$ returns a close approximation to arc(x,y) in \mathcal{R} , with maximum error $max_error_tan_F$.

NOTE 1 - The arc operations, with the arguments swapped, are often called arctan20

Further requirements on the arc_F^* approximation helper function are:

```
\begin{array}{ll} arc_F^*(x,0)=0 & \text{if } x\in F \text{ and } x>0 \\ arc_F^*(x,x)=\pi/4 & \text{if } x\in F \text{ and } x>0 \\ arc_F^*(0,y)=\pi/2 & \text{if } y\in F \text{ and } y>0 \\ arc_F^*(x,-x)=3\cdot\pi/4 & \text{if } x\in F \text{ and } x<0 \\ arc_F^*(x,0)=\pi & \text{if } x\in F \text{ and } x<0 \\ arc_F^*(x,-y)=-arc_F^*(x,y) & \text{if } x,y\in F \text{ and } y\neq 0 \text{ or } x>0) \end{array}
```

The $arc_F^{\#}$ range limitation helper function (for $x, y \in F$):

$$arc_F^{\#}(x,y) = \max\{up_F(-\pi), \min\{arc_F^{*}(x,y), down_F(\pi)\}\}$$

The arc_F operation:

```
arc_F: F \times F \to F \cup \{\mathbf{underflow}\}\
                    = result_F^*(arc_F^\#(x, y), nearest_F)
                                                     if x, y \in F and (x \neq 0 \text{ or } y \neq 0)
                                                     if x = 0 and y = 0
                    = down_F(\pi)
                                                     if x = -0 and y = 0
                                                     if x = -\mathbf{0} and y \in F \cup \{-\infty, +\infty\} and y \neq 0
                    = arc_F(0, y)
                    = neg_F(arc_F(x,0))
                                                     if y = -\mathbf{0} and x \in F \cup \{-\infty, -\mathbf{0}, +\infty\}
                                                     if x = +\infty and y \in F and y \ge 0
                                                     if x = +\infty and y \in F and y < 0
                     = nearest_F(\pi/4)
                                                     if x = +\infty and y = +\infty
                    = nearest_F(\pi/2)
                                                     if x \in F and y = +\infty
                    = nearest_F(3 \cdot \pi/4)
                                                    if x = -\infty and y = +\infty
                    = down_F(\pi)
                                                     if x = -\infty and y \in F and y \ge 0
                                                     if x = -\infty and y \in F and y < 0
                    = nearest_F(-3 \cdot \pi/4)
                                                     if x = -\infty and y = -\infty
                                                     if x \in F and y = -\infty
                    = nearest_F(-\pi/2)
                    = nearest_F(-\pi/4)
                                                     if x = +\infty and y = -\infty
                    = no\_result2_F(x, y)
                                                     otherwise
```

NOTE 2 – Note that the arc operations do *not* return an **invalid** notification at the origin (both arguments in $\{-0,0\}$). See B.5.3.8. Bindings may choose to alter this behaviour.

5.3.9 Operations for trigonometrics with given angular unit

There shall be one big-angle parameter for argument angular-unit trigonometric operations:

```
big\_angle\_u_F \in F
```

It should have the following default value:

$$big_angle_u_F = \lceil r_F^{\lceil p_F/2 \rceil}/6 \rceil$$

A binding or implementation can include a method to change the value of this parameter. This method should only allow the value of this parameter to be set to a value greater than or equal to 1 and such that $ulp_F(big_angle_u_F) \leq 1/2000$.

```
NOTE 1 – In order to allow ulp_F(big\_angle\_u_F) \leqslant 1/2000, p_F \geqslant 2 + \log_{r_F}(1000) should hold.
```

There shall be a derived parameter signifying the minimum allowed angular unit:

```
min\_angular\_unit_F = r_F \cdot fminN_F/epsilon_F
```

```
NOTE 2 – That is, min\_angular\_unit_F = r_F^{(emin_F - 1 + p_F)}
```

For use in the approximation helper function's signatures, define

```
F^{u} = (F \cup \{n \cdot u/8, n \cdot u/12 \mid n \in \mathcal{Z}\}) \cap [-big\_angle\_u_F \cdot |u|, big\_angle\_u_F \cdot |u|]
```

Note that u is a parameter here, a parameter which is the value of the first argument to the approximation helper function. To signify this, the notation (y:E) is used below.

To make the specifications below a bit easier to express, let

```
G_F = \{x \in F \mid min\_angular\_unit_F \leq |x|\}.
```

Let $T = \{1, 2, 360, 400, 6400\}$. T consists of angle values for exactly one revolution for some common non-radian angular units: cycles, half-cycles, arc degrees, grades, and mils.

There shall be two parameterised maximum error parameters for argument angular-unit trigonometric operations:

```
max\_error\_sinu_F : F \to F \cup \{invalid\}

max\_error\_tanu_F : F \to F \cup \{invalid\}
```

For $u \in G_F$, the $max_error_sinu_F(u)$ parameter shall have a value that is $\leq 2 \cdot max_error_sin_F$. The $max_error_sinu_F(u)$ parameter shall have the value of $max_error_sin_F$ if $|u| \in T$. For $u \in G_F$, the $max_error_tanu_F(u)$ parameter shall have a value that is $\leq 2 \cdot max_error_tan_F$. The $max_error_tanu_F(u)$ parameter shall have the value of $max_error_tan_F$ if $|u| \in T$. The $max_error_sinu_F(u)$ and $max_error_tanu_F(u)$ parameters return invalid(qNaN) if $u \notin G_F$.

5.3.9.1 Argument angular-unit angle normalisation

The argument angular-unit normalisation computes exactly $rad(2 \cdot \pi \cdot x/u) \cdot u/(2 \cdot \pi)$, where x is the angular value, and u is the angular unit.

The $cycle_F$ operation:

```
\begin{aligned} cycle_F: F \times F \to F \cup \{ \mathbf{-0}, \mathbf{absolute\_precision\_underflow}, \mathbf{invalid} \} \\ cycle_F(u,x) &= residue_F(x,u) & \text{if } u \in G_F \text{ and } (x = \mathbf{-0} \text{ or } \\ & (x \in F \text{ and } |x/u| \leqslant big\_angle\_u_F)) \\ &= \mathbf{absolute\_precision\_underflow}(\mathbf{qNaN}) \\ &= no\_result2_F(u,x) & \text{otherwise} \end{aligned}
```

The $axis_cycle_F$ operation:

$$axis_cycle_F: F \times F \rightarrow (\{(1,0),(0,1),(-1,0),(0,-1)\} \times (F \cup \{-\mathbf{0}\})) \cup \\ \quad \{absolute_precision_underflow,invalid\} \\ axis_cycle_F(u,x) \\ = (axis(u,x),result_F(x-(\operatorname{round}(x\cdot 4/u)\cdot u/4),nearest_F)) \\ \quad \text{if } u \in G_F \text{ and } x \in F \text{ and } |x/u| \leqslant big_angle_u_F \text{ and } \\ \quad (x/u \geqslant 0 \text{ or } x-(\operatorname{round}(x\cdot 4/u)\cdot u/4) \neq 0) \\ = (axis(u,x),-\mathbf{0}) \qquad \text{if } u \in G_F \text{ and } x \in F \text{ and } |x/u| \leqslant big_angle_u_F \text{ and } \\ \quad x/u < 0 \text{ and } x-(\operatorname{round}(x\cdot 4/u)\cdot u/4) = 0 \\ = ((1,0),-\mathbf{0}) \qquad \text{if } u \in G_F \text{ and } x = -\mathbf{0} \\ = absolute_precision_underflow((\mathbf{qNaN},\mathbf{qNaN}),\mathbf{qNaN}) \\ \quad \text{if } u \in G_F \text{ and } x \in F \text{ and } |x/u| > big_angle_u_F \\ = ((\mathbf{qNaN},\mathbf{qNaN}),\mathbf{qNaN}) \\ \quad \text{if at least one of } x \text{ and } u \text{ is a quiet NaN and}$$

neither is a signalling NaN = invalid((qNaN, qNaN), qNaN)otherwise

where

$$axis(u, x)$$
 = (1,0) if round($x \cdot 4/u$) = $4 \cdot n$
= (0,1) if round($x \cdot 4/u$) = $4 \cdot n + 1$
= (-1,0) if round($x \cdot 4/u$) = $4 \cdot n + 2$
= (0,-1) if round($x \cdot 4/u$) = $4 \cdot n + 3$

for some $n \in \mathcal{Z}$.

NOTES

- 1 $axis_cycle_F(u, x)$ is exact when $div_F(u, 4) = u/4$.
- $2 \quad cycle_F$ is an exact operation.
- 3 $cycle_F(u,x)$ is $-\mathbf{0}$ or has a result in the interval [-|u/2|,|u/2|] if there is no notification.
- 4 A zero resulting angle is negative if the original angle value is negative.
- 5 The $cycle_F$ operation is used also in the specifications of the unit argument trigonometric operations. This does not imply that the implementation has to use the $cycle_F$ operation, when implementing the operations. It only implies that the results (including notifications) must be as if the $cycle_F$ operation was used.

Argument angular-unit sine

The $sinu_F^*$ approximation helper function: $sinu_F^*: (u:F) \times F^u \to \mathcal{R}$

$$sinu_F^*: (u:F) \times F^u \to \mathcal{R}$$

 $sinu_F^*(u,x)$ returns a close approximation to $sin(x\cdot 2\cdot \pi/u)$ in \mathcal{R} if $u\neq 0$, with maximum error $max_error_sinu_F(u)$.

Further requirements on the $sinu_F^*$ approximation helper function are:

$$sinu_F^*(u, n \cdot u + x) = sinu_F^*(u, x)$$
 if $n \in \mathbb{Z}$ and $u \in F$ and $u \neq 0$ and $x \in F^u$
 $sinu_F^*(u, u/12) = 1/2$ if $u \in F$ and $u \neq 0$
 $sinu_F^*(u, u/4) = 1$ if $u \in F$ and $u \neq 0$

```
\begin{aligned} sinu_F^*(u,5\cdot u/12) &= 1/2 & \text{if } u \in F \text{ and } u \neq 0 \\ sinu_F^*(u,-x) &= -sinu_F^*(u,x) & \text{if } u \in F \text{ and } u \neq 0 \text{ and } x \in F^u \\ sinu_F^*(-u,x) &= -sinu_F^*(u,x) & \text{if } u \in F \text{ and } u \neq 0 \text{ and } x \in F^u \\ \text{NOTE } - sinu_F^*(u,x) &\approx x \cdot 2 \cdot \pi/u \text{ if } |x \cdot 2 \cdot \pi/u| < fminN_F. \end{aligned}
```

The $sinu_F$ operation:

 $sinu_F: F \times F \to F \cup \{-0, underflow, absolute_precision_underflow, invalid\}$

$$sinu_F(u,x) = result_F^*(sinu_F^*(u,x), nearest_F)$$
 if $cycle_F(u,x) \in F$ and $cycle_F(u,x) \notin \{-u/2, 0, u/2\}$
$$= div_F(0,u) \qquad \text{if } cycle_F(u,x) \in \{0, u/2\}$$

$$= div_F(-\mathbf{0},u) \qquad \text{if } cycle_F(u,x) \in \{-u/2, -\mathbf{0}\}$$

$$= cycle_F(u,x) \qquad \text{otherwise}$$

5.3.9.3 Argument angular-unit cosine

The $cosu_F^*$ approximation helper function:

$$cosu_F^*: (u:F) \times F^u \to \mathcal{R}$$

 $cosu_F^*(u,x)$ returns a close approximation to $cos(x \cdot 2 \cdot \pi/u)$ in \mathcal{R} if $w \neq 0$, with maximum error $max_error_sinu_F(u)$.

Further requirements on the $cosu_F^*$ approximation helper function are:

```
\begin{aligned} \cos u_F^*(u,n\cdot u+x) &= \cos u_F^*(u,x) & \text{if } n\in \mathcal{Z} \text{ and } u\in F \text{ and } u\neq 0 \text{ and } x\in F^u \\ \cos u_F^*(u,u/6) &= 1/2 & \text{if } u\in F \text{ and } u\neq 0 \\ \cos u_F^*(u,u/3) &= -1/2 & \text{if } u\in F \text{ and } u\neq 0 \\ \cos u_F^*(u,u/2) &= -1 & \text{if } u\in F \text{ and } u\neq 0 \\ \cos u_F^*(u,-x) &= \cos u_F^*(u,x) & \text{if } u\in F \text{ and } u\neq 0 \\ \cos u_F^*(-u,x) &= \cos u_F^*(u,x) & \text{if } u\in F \text{ and } u\neq 0 \text{ and } x\in F^u \\ \cos u_F^*(-u,x) &= \cos u_F^*(u,x) & \text{if } u\in F \text{ and } u\neq 0 \text{ and } x\in F^u \\ \text{NOTE } - \cos u_F^*(u,x) &= 1 \text{ should hold if } |x\cdot 2\cdot \pi/u| < \sqrt{epsilon_F/r_F} \end{aligned}
```

The $cosu_F$ operation:

 $cosu_F: F \times F \to F \cup \{underflow, absolute_precision_underflow, invalid\}$

5.3.9.4 Argument angular-unit tangent

The $tanu_F^*$ approximation helper function:

$$tanu_F^*: (u:F) \times F^u \to \mathcal{R}$$

 $tanu_F^*(u,x)$ returns a close approximation to $tan(x \cdot 2 \cdot \pi/u)$ in \mathcal{R} if $u \neq 0$, with maximum error $max_error_tanu_F(u)$.

Further requirements on the $tanu_F^*$ approximation helper function are:

```
tanu_F^*(u, n \cdot u + x) = tanu_F^*(u, x) \qquad \text{if } n \in \mathcal{Z} \text{ and } u \in F \text{ and } u \neq 0 \text{ and } x \in F^u
tanu_F^*(u, u/8) = 1 \qquad \text{if } u \in F \text{ and } u \neq 0
tanu_F^*(u, 3 \cdot u/8) = -1 \qquad \text{if } u \in F \text{ and } u \neq 0
tanu_F^*(u, -x) = -tanu_F^*(u, x) \qquad \text{if } u \in F \text{ and } u \neq 0 \text{ and } x \in F^u
tanu_F^*(-u, x) = -tanu_F^*(u, x) \qquad \text{if } u \in F \text{ and } u \neq 0 \text{ and } x \in F^u
NOTE \ 1 - tanu_F^*(u, x) \approx x \cdot 2 \cdot \pi/u \text{ if } |x \cdot 2 \cdot \pi/u| < fminN_F.
```

The $tanu_F$ operation:

 $tanu_F: F \times F \rightarrow F \cup \{-0, underflow, overflow, infinitary, \}$

absolute_precision_underflow, invalid}

```
tanu_{F}(u,x) = result_{F}^{*}(tanu_{F}^{*}(u,x), nearest_{F})
if \ cycle_{F}(u,x) \in F \text{ and}
cycle_{F}(u,x) \notin \{-u/2, -u/4, 0, u/4, u/2\}\}
= div_{F}(0,u) \qquad \text{if } cycle_{F}(u,x) \in \{-u/2, 0\}
= div_{F}(-\mathbf{0},u) \qquad \text{if } cycle_{F}(u,x) \in \{-\mathbf{0}, u/2\}\}
= \mathbf{infinitary}(+\infty) \qquad \text{if } cycle_{F}(u,x) = u/4
= \mathbf{infinitary}(-\infty) \qquad \text{if } cycle_{F}(u,x) = -u/4
= cycle_{F}(u,x) \qquad \text{otherwise}
```

NOTE 2 – The **infinitary** notification can arise for $tanu_F(u,x)$ only when u/4 is in F.

5.3.9.5 Argument angular-unit cotangent

The $cotu_F^*$ approximation helper function:

$$cotu_F^*: (u:F) \times F^u \to \mathcal{R}$$

 $cotu_F^*(u,x)$ returns a close approximation to $cot(x^*2 \cdot \pi/u)$ in \mathcal{R} if $u \neq 0$, with maximum error $max_error_tanu_F(u)$.

Further requirements on the $\cot u_F^*$ approximation helper function are:

```
\cot u_F^*(u, n \cdot u + x) = \cot u_F^*(u, x) \qquad \text{if } n \in \mathcal{Z} \text{ and } u \in F \text{ and } u \neq 0 \text{ and } x \in F^u
\cot u_F^*(u, u/8) = 1 \qquad \text{if } u \in F \text{ and } u \neq 0
\cot u_F^*(u, x) = -\cot u_F^*(u, x) \qquad \text{if } u \in F \text{ and } u \neq 0
\cot u_F^*(u, -x) = -\cot u_F^*(u, x) \qquad \text{if } u \in F \text{ and } u \neq 0 \text{ and } x \in F^u
\cot u_F^*(-u, x) = -\cot u_F^*(u, x) \qquad \text{if } u \in F \text{ and } u \neq 0 \text{ and } x \in F^u
```

The $cotu_F$ operation.

 $cotu_F: F \times F \mapsto F \cup \{-0, \text{underflow}, \text{overflow}, \text{infinitary}, \}$

$absolute_precision_underflow, invalid\}$

```
cotu_{F}(u,x) = result_{F}^{*}(cotu_{F}^{*}(u,x), nearest_{F})
if \ cycle_{F}(u,x) \in F \text{ and}
cycle_{F}(u,x) \notin \{-u/2, -u/4, 0, u/2\}
= -\mathbf{0} \qquad \text{if } cycle_{F}(u,x) = -u/4
= div_{F}(u, tanu_{F}(u,x)) \qquad \text{if } cycle_{F}(u,x) \in \{-u/2, -\mathbf{0}, 0, u/2\}
= cycle_{F}(u,x) \qquad \text{otherwise}
```

5.3.9.6 Argument angular-unit secant

The $secu_F^*$ approximation helper function:

```
secu_F^*: (u:F) \times F^u \to \mathcal{R}
```

 $secu_F^*(u,x)$ returns a close approximation to $sec(x \cdot 2 \cdot \pi/u)$ in \mathcal{R} if $u \neq 0$, with maximum error $max_error_tanu_F(u)$.

Further requirements on the $secu_F^*$ approximation helper function are:

```
\begin{split} & secu_F^*(u,n\cdot u+x) = secu_F^*(u,x) & \text{if } n\in\mathcal{Z} \text{ and } u\in F \text{ and } u\neq 0 \text{ and } x\in F^u \\ & secu_F^*(u,0)=1 & \text{if } u\in F \text{ and } u\neq 0 \\ & secu_F^*(u,u/6)=2 & \text{if } u\in F \text{ and } u\neq 0 \\ & secu_F^*(u,u/3)=-2 & \text{if } u\in F \text{ and } u\neq 0 \\ & secu_F^*(u,u/2)=-1 & \text{if } u\in F \text{ and } u\neq 0 \\ & secu_F^*(u,-x)=secu_F^*(u,x) & \text{if } u\in F \text{ and } u\neq 0 \text{ and } x\in F^u \\ & secu_F^*(-u,x)=secu_F^*(u,x) & \text{if } u\in F \text{ and } u\neq 0 \text{ and } x\in F^u \\ & secu_F^*(-u,x)=secu_F^*(u,x) & \text{if } u\in F \text{ and } u\neq 0 \text{ and } x\in F^u \\ & secu_F^*(-u,x)=secu_F^*(u,x) & \text{if } u\in F \text{ and } u\neq 0 \text{ and } x\in F^u \\ & secu_F^*(-u,x)=secu_F^*(-u,x) & \text{if } u\in F \text{ and } u\neq 0 \text{ and } x\in F^u \\ & secu_F^*(-u,x)=secu_F^*(-u,x) & \text{if } u\in F \text{ and } u\neq 0 \text{ and } x\in F^u \\ & secu_F^*(-u,x)=secu_F^*(-u,x) & \text{if } u\in F \text{ and } u\neq 0 \text{ and } x\in F^u \\ & secu_F^*(-u,x)=secu_F^*(-u,x) & \text{if } u\in F \text{ and } u\neq 0 \text{ and } x\in F^u \\ & secu_F^*(-u,x)=secu_F^*(-u,x) & \text{if } u\in F \text{ and } u\neq 0 \text{ and } x\in F^u \\ & secu_F^*(-u,x)=secu_F^*(-u,x) & \text{if } u\in F \text{ and } u\neq 0 \text{ and } x\in F^u \\ & secu_F^*(-u,x)=secu_F^*(-u,x) & \text{if } u\in F \text{ and } u\neq 0 \text{ and } x\in F^u \\ & secu_F^*(-u,x)=secu_F^*(-u,x) & \text{if } u\in F \text{ and } u\neq 0 \text{ and } x\in F^u \\ & secu_F^*(-u,x)=secu_F^*(-u,x) & \text{if } u\in F \text{ and } u\neq 0 \text{ and } x\in F^u \\ & secu_F^*(-u,x)=secu_F^*(-u,x) & \text{if } u\in F \text{ and } u\neq 0 \text{ and } x\in F^u \\ & secu_F^*(-u,x)=secu_F^*(-u,x) & \text{if } u\in F \text{ and } u\neq 0 \text{ and } x\in F^u \\ & secu_F^*(-u,x)=secu_F^*(-u,x) & \text{if } u\in F \text{ and } u\neq 0 \text{ and } x\in F^u \\ & secu_F^*(-u,x)=secu_F^*(-u,x) & \text{if } u\in F \text{ and } u\neq 0 \text{ and } x\in F^u \\ & secu_F^*(-u,x)=secu_F^*(-u,x) & \text{if } u\in F \text{ and } u\neq 0 \text{ and } x\in F^u \\ & secu_F^*(-u,x)=secu_F^*(-u,x) & \text{if } u\in F \text{ and } u\neq 0 \text{ and } x\in F^u \\ & secu_F^*(-u,x)=secu_F^*(-u,x) & \text{if } u\in F \text{ and } u\neq 0 \text{ and } x\in F^u \\ & secu_F^*(-u,x)=secu_F^*(-u,x) & \text{if } u\in F \text{ and } u\neq 0 \text{ and } x\in F^u \\ & secu_F^*(-u,x)=secu_F^*(-u,x) & \text{if } u\in F \text{ and } u\neq 0 \text{ and } x\in F^u \\ & secu_F^*(-u,x)=secu_F^*(-u,x)
```

The $secu_F$ operation:

 $secu_F: F \times F \to F \cup \{\text{overflow}, \text{infinitary}, \text{absolute_precision_underflow}, \text{invalid}\}$

```
secu_F(u,x) = result_F^*(secu_F^*(u,x), nearest_F)
if \ cycle_F(u,x) \notin F \text{ and }
cycle_F(u,x) \notin \{-u/4, u/4\}
= div_F(1, cosu_F(u,x)) \quad \text{if } cycle_F(u,x) \in \{-u/4, -\mathbf{0}, u/4\}
= cycle_F(u,x) \quad \text{otherwise}
```

5.3.9.7 Argument angular-unit cosecant

The $cscu_F^*$ approximation helper function:

```
cscu_F^*: (u:F) \times F^u \to \mathcal{R}
```

 $cscu_F^*(u,x)$ returns a close approximation to $csc(x \cdot 2 \cdot \pi/u)$ in \mathcal{R} if $u \neq 0$, with maximum error $max_error_tanu_F(u)$.

Further requirements on the $cscu_F^*$ approximation helper function are:

```
\begin{array}{lll} cscu_F^*(u,n\cdot u+x)=cscu_F^*(u,x) & \text{if } n\in\mathcal{Z} \text{ and } u\in F \text{ and } u\in 0 \text{ and } x\in F^u\\ cscu_F^*(u,u/12)=2 & \text{if } u\in F \text{ and } u\neq 0\\ cscu_F^*(u,u/4)=1 & \text{if } u\in F \text{ and } u\neq 0\\ cscu_F^*(u,5\cdot u/12)=2 & \text{if } u\in F \text{ and } u\neq 0\\ cscu_F^*(u,x)=-cscu_F^*(u,x) & \text{if } u\in F \text{ and } u\neq 0 \text{ and } x\in F^u\\ cscu_F^*(u,x)=-cscu_F^*(u,x) & \text{if } u\in F \text{ and } u\neq 0 \text{ and } x\in F^u\\ \end{array}
```

The $cscl_F$ operation:

 $cscu_F: F \times F \to F \cup \{\text{overflow}, \text{infinitary}, \text{absolute_precision_underflow}, \text{invalid}\}$

```
cscu_F(u,x) = result_F^*(cscu_F^*(u,x), nearest_F) if cycle_F(u,x) \in F and cycle_F(u,x) \notin \{-u/2, 0, u/2\} = div_F(1, sinu_F(u,x)) if cycle_F(u,x) \in \{-u/2, -\mathbf{0}, 0, u/2\} = cycle_F(u,x) otherwise
```

5.3.9.8 Argument angular-unit cosine with sine

 $cossinu_F: F \times F \rightarrow (F \times (F \cup \{-0\})) \cup \{ \text{underflow}, \text{absolute_precision_underflow}, \\ \text{invalid} \}$

$$cossinu_F(u, x) = (cosu_F(u, x), sinu_F(u, x))$$

NOTES

- 1 If there is an absolute_precision_underflow notification, then both result parts suffer from the absolute_precision_underflow and the continuation values for both parts are qNaN. Similarly for NaN and infinitary arguments, as well as an angular unit with too small absolute value.
- 2 If there is an **underflow** notification, only one of the result parts suffer from the underflow, and the other part has an absolute value greater than $fminN_F$.

5.3.9.9 Argument angular-unit arc sine

The $arcsinu_F^*$ approximation helper function:

$$arcsinu_F^*: F \times F \to \mathcal{R}$$

 $arcsinu_F^*(u,x)$ returns a close approximation to $arcsin(x) \cdot u/(2\pi)$ in \mathcal{R} , with maximum error $max_error_sinu_F(u)$.

Further requirements on the $arcsinu_F^*$ approximation helper function are:

$$\begin{aligned} & \arcsin u_F^*(u,1/2) = u/12 & \text{if } u \in F \\ & \arcsin u_F^*(u,1) = u/4 & \text{if } u \in F \\ & \arcsin u_F^*(u,-x) = -\arcsin u_F^*(u,x) & \text{if } u,x \in F \\ & \arcsin u_F^*(-u,x) = -\arcsin u_F^*(u,x) & \text{if } u,x \in F \\ & \text{NOTE } - & \arcsin u_F^*(u,x) \approx u/(2 \cdot \pi) \text{ if } |x| < fminN_F. \end{aligned}$$

The $arcsinu_F^{\#}$ range limitation helper function (for $u, x \in F$):

$$arcsinu_F^{\#}(u,x) = \max\{up_F(-|u\mathcal{A}|), \min\{arcsinu_F^{*}(u,x), down_F(|u/4|)\}\}$$

The $arcsinu_F$ operation:

$$arcsinu_F: F \times F \to F \cup \{-0, \text{underflow}, \text{invalid}\}$$

$$arcsinu_F(u,x) = result_F^*(arcsinu_F^\#(u,x), nearest_F)$$
 if $u \in G_F$ and $x \in F$ and $|x| \le 1$ and $x \ne 0$ if $u \in G_F$ and $x \in \{-\mathbf{0}, 0\}$ otherwise

5.3.9.10 Argument angular-unit arc cosine

The $arccosu_F^*$ approximation helper function:

$$arccosu_F^*: F \times F \to \mathcal{R}$$

 $arccosu_F^*(u,x)$ returns a close approximation to $arccos(x) \cdot u/(2 \cdot \pi)$ in \mathcal{R} , with maximum error $max_error_sinu_F(u)$.

Further requirements on the $arccosu_F^*$ approximation helper function are:

```
arccosu_F^*(u, 1/2) = u/6
                                               if u \in F
arccosu_F^*(u,0) = u/4
                                               if u \in F
arccosu_F^*(u, -1/2) = u/3
                                               if u \in F
arccosu_F^*(u, -1) = u/2
                                               if u \in F
arccosu_F^*(-u, x) = -arccosu_F^*(u, x)
                                               if u, x \in F
```

The $arccosu_F^{\#}$ range limitation helper function (for $u, x \in F$):

```
arccosu_F^{\#}(u, x) = \max\{up_F(-|u/2|), \min\{arccosu_F^*(u, x), down_F(|u/2|)\}\}
```

The $arccosu_F$ operation:

```
arccosu_F: F \times F \rightarrow F \cup \{ \mathbf{underflow}, \mathbf{invalid} \}
```

```
arccosu_F(u, x) = result_F^*(arccosu_F^\#(u, x), nearest_F)
                                              if u \in G_F and x \in F and |x| \leq 1
                                              if u \in G_F and x = -0
                  = nearest_F(u/4)
                  = no\_result2_F(u, x)
                                              otherwise
```

5.3.9.11Argument angular-unit arc tangent

The $arctanu_F^*$ approximation helper function:

```
arctanu_F^*: F \times F \to \mathcal{R}
```

 $arctanu_F^*(u,x)$ returns a close approximation to $arctan(x) : u/(2 \cdot \pi)$ in \mathcal{R} , with maximum error $max_error_tanu_F(u)$.

Further requirements on the $arctanu_F^*$ approximation helper function are:

```
arctanu_F^*(u,1) = u/8
                                                 if u, x \in F and arctanu_F^*(u, x) \neq \arctan(x) \cdot u/(2 \cdot \pi)
arctanu_F^*(u,x) = u/4
                                                     and x > 3 \cdot r_F/epsilon_F
arctanu_F^*(u, -x) = -arctanu_F^*(u, x)
arctanu_F^*(-u, x) = -arctanu_F^*(u, x)
NOTE 1 - arctanu_F^*(u, x) \approx u/(2 \cdot \pi) if |x| < fminN_F
```

The $arctanu_F^{\#}$ range limitation helper function (for $u, x \in F$):

```
arctanu_F^{\#}(u, x) = \max\{up_F(-|u/4|), \min\{arctanu_F^*(u, x), down_F(|u/4|)\}\}
```

The $arctanu_F$ operation:

 $arctanu_F: F \hookrightarrow F \cup \{-0, \mathbf{underflow}, \mathbf{invalid}\}$

if $u \in G_F$ and $x \in F$ and $x \neq 0$

```
arctanu_F(u,x) = result_F^*(arctanu_F^\#(u,x), nearest_F)
= mul_F(x,u) \quad \text{if } u \in G_F \text{ a}
= up_F(-u/4) \quad \text{if } u \in G_F \text{ a}
                                                           if u \in G_F and x \in \{-0, 0\}
                                                           if u \in G_F and x = -\infty and u > 0
                       = down_F(u/4)
                                                           if u \in G_F and x = +\infty and u > 0
                       = down_F(-u/4)
                                                           if u \in G_F and x = -\infty and u < 0
                                                           if u \in G_F and x = +\infty and u < 0
                       = up_F(u/4)
                                                           otherwise
                       = no\_result2_F(u, x)
```

NOTE 2 - $arctanu_F(u, x) \approx arcu_F(u, 1, x)$. ($arcu_F$ is specified in subclause 5.3.9.15 below.)

5.3.9.12 Argument angular-unit arc cotangent

This clause specifies two inverse cotangent operations. One approximating the sign symmetric (but discontinuous at 0) arccot, the other approximating the continuous (but not sign symmetric) arccotc (both for non-radian angular units).

The $arccotu_F^*$ approximation helper function:

```
arccotu_F^*: F \times F \to \mathcal{R}
```

 $arccotu_F^*(u,x)$ returns a close approximation to $arccot(x) \cdot u/(2 \cdot \pi)$ in \mathcal{R} , with maximum error $max_error_tanu_F(u)$.

The $arccotcu_F^*$ approximation helper function:

$$arccotcu_F^*: F \times F \to \mathcal{R}$$

 $arccotcu_F^*(u,x)$ returns a close approximation to $arccotc(x) \cdot u/(2 \cdot \pi)$ in \mathcal{R} , with maximum error $max_error_tanu_F(u)$.

Further requirements on the $arccotu_F^*$ and $arccotcu_F^*$ approximation helper functions are:

```
\begin{aligned} &\operatorname{arccotu}_F^*(u,1) = u/8 & \text{if } u \in F \\ &\operatorname{arccotu}_F^*(u,0) = u/4 & \text{if } u \in F \\ &\operatorname{arccotu}_F^*(u,-x) = -\operatorname{arccotu}_F^*(u,x) & \text{if } u,x \in F \text{ and } x \neq 0 \\ &\operatorname{arccotcu}_F^*(u,x) = \operatorname{arccotu}_F^*(u,x) & \text{if } u,x \in F \text{ and } x \geqslant 0 \\ &\operatorname{arccotcu}_F^*(u,-1) = 3 \cdot u/8 & \text{if } u \in F \\ &\operatorname{arccotcu}_F^*(u,x) = u/2 & \text{if } u,x \in F \text{ and } \operatorname{arccotcu}_F^*(u,x) \neq \operatorname{arccotc}(x) \cdot u/(2 \cdot \pi) \\ &\operatorname{arccotcu}_F^*(-u,x) = -\operatorname{arccotcu}_F^*(u,x) & \text{if } u,x \in F \text{ and } \operatorname{arccotcu}_F^*(u,x) \neq \operatorname{arccotc}(x) \cdot u/(2 \cdot \pi) \\ &\operatorname{arccotcu}_F^*(-u,x) = -\operatorname{arccotcu}_F^*(u,x) & \text{if } u,x \in F \end{aligned}
```

The $arccotu_F^{\#}$ and $arccotcu_F^{\#}$ range limitation helper functions (for $u, x \in F$):

$$\begin{aligned} & \operatorname{arccot} u_F^\#(u,x) = \max\{up_F(-|u/4|), \min\{\operatorname{arccot} u_F^*(u,x), \operatorname{down}_F(|u/4|)\}\} \\ & \operatorname{arccot} u_F^\#(u,x) = \max\{up_F(-|u/2|), \min\{\operatorname{arccot} u_F^*(u,x), \operatorname{down}_F(|u/2|)\}\} \end{aligned}$$

The $arccotu_F$ operation:

 $arccotu_F: F \times F \to F \cup \{\mathbf{underflow}, \mathbf{invalid}\}$

```
arccotu_F(u,x) = result_F^*(arccotu_F^\#(u,x), nearest_F) if u \in G_F and x \in F = div_F(u,x) if u \in G_F and x \in -\mathbf{0} if u \in G_F and x \in \{-\infty, +\infty\} = no\_result2_F(u,x) otherwise
```

NOTES

- 1 $arccotu_F(u, neg_F(x)) = neg_F(arccotu_F(u, x)).$
- 2 Due to the range limitation, $arccotu_F(u,0)$ need not equal $arccotcu_F(u,0)$.

The $arccotcu_F$ operation:

```
\begin{split} arccotcu_F: F \times F \to F \cup \{\mathbf{underflow}, \mathbf{invalid}\} \\ arccotcu_F(u, x) &= result_F^*(arccotcu_F^\#(u, x), nearest_F) \\ &\qquad \qquad \text{if } u \in G_F \text{ and } x \in F \\ &= nearest_F(u/4) \qquad \qquad \text{if } u \in G_F \text{ and } x = -\mathbf{0} \\ &= down_F(u/2) \qquad \qquad \text{if } u \in G_F \text{ and } x = -\mathbf{\infty} \text{ and } u > 0 \end{split}
```

```
= up_F(u/2)
                            if u \in G_F and x = -\infty and u < 0
= div_F(u, x)
                            if u \in G_F and x = +\infty
= no\_result 2_F(u, x)
                            otherwise
```

NOTE 3 - $arccotcu_F(u, x) \approx arcu_F(u, x, 1)$. ($arcu_F$ is specified in subclause 5.3.9.15 below.)

5.3.9.13 Argument angular-unit arc secant

The $arcsecu_F^*$ approximation helper function:

```
arcsecu_F^*: F \times F \to \mathcal{R}
```

 $arcsecu_F^*(u,x)$ returns a close approximation to $arcsec(x) \cdot u/(2 \cdot \pi)$ in \mathcal{R} , with maximum error $max_error_tanu_F(u)$.

Further requirements on the $arcsecu_F^*$ approximation helper function are:

```
arcsecu_F^*(u,2) = u/6
                                                   if u \in F
arcsecu_F^*(u, -2) = u/3
                                                   if u \in F
arcsecu_F^*(u, -1) = u/2
                                                   if u \in F
\operatorname{arcsecu}_F^*(u,x)\leqslant u/4
                                                   if u, x \in F and x > 0 and u > 0
                                                   if u, x \in F and x < 0 and u > 0
arcsecu_F^*(u, x) \geqslant u/4
arcsecu_F^*(u, x) = u/4
                                                   if u, x \in F and arcsecu_F^*(u, x) \neq arcsec(x) \cdot u/(2 \cdot \pi)
                                                      and |x| > 3C_F/epsilon_F
                                                  if u, x \in F
arcsecu_F^*(-u, x) = -arcsecu_F^*(u, x)
```

The $arcsecu_F^{\#}$ range limitation helper function (for $x \in F$):

$$\begin{split} arcsecu_F^\#(u,x) &= \max\{up_F(-|u/4|), \min\{arcsecu_F^*(u,x), down_F(|u/4|)\}\} \\ &= \max\{up_F(u/4), \min\{arcsecu_F^*(u,x), down_F(u/2)\}\} \\ &= \max\{up_F(u/4), \min\{arcsecu_F^*(u,x), down_F(u/2)\}\} \\ &= \max\{up_F(u/2), \min\{arcsecu_F^*(u,x), down_F(u/4)\}\} \\ &= \inf x \leqslant -1 \text{ and } u < 0 \end{split}$$

The $arcsecu_F$ operation:

```
arcsecu_F: F \times F \rightarrow F underflow, invalid} arcsecu_F(u, x) = result_F^*(arcsecu_F^\#(u, x), nearest_F) if u \in G_F a
                                                                         if \ u \in G_F \text{ and } x \in F \text{ and } 1 \leq |x|
= down_F(u/4) \qquad if \ u \in G_F \text{ and } x = -\infty \text{ and } u > 0
= up_F(u/4) \qquad if \ u \in G_F \text{ and } x = +\infty \text{ and } u > 0
= up_F(u/4) \qquad if \ u \in G_F \text{ and } x = -\infty \text{ and } u < 0
= down_F(u/4) \qquad if \ u \in G_F \text{ and } x = +\infty \text{ and } u < 0
= no\_result2_F(u, x) \qquad \text{otherwise}
```

5.3.9.14 Argument angular-unit arc cosecant

The $arccscu_F^*$ approximation helper function:

```
arccscu_F^*: F \times F \to \mathcal{R}
```

 $arccscu_F^*(u,x)$ returns a close approximation to $arccsc(x) \cdot u/(2 \cdot \pi)$ in \mathcal{R} , with maximum error $max_error_tanu_F(u)$.

Further requirements on the $arccscu_F^*$ approximation helper function are:

```
if u \in F
arccscu_F^*(u,2) = u/12
                                                if u \in F
arccscu_F^*(u,1) = u/4
arccscu_F^*(u, -x) = -arccscu_F^*(u, x)
                                                if u, x \in F
                                                if u, x \in F
arccscu_F^*(-u, x) = -arccscu_F^*(u, x)
```

The $arccscu_F^{\#}$ range limitation helper function (for $u, x \in F$):

```
arccscu_F^{\#}(u, x) = \max\{up_F(-|u/4|), \min\{arccscu_F^{*}(u, x), down_F(|u/4|)\}\}
```

The $arccscu_F$ operation:

```
arccscu_F: F \times F \to F \cup \{\mathbf{underflow}, \mathbf{invalid}\}
```

```
(x, wearest_F)
if u \in G_F and x \in F and x \in F and x \in F and x \in G_F and x \in G_F and x \in G_F and x \in G_F
arccscu_F(u, x) = result_F^*(arccscu_F^\#(u, x), nearest_F)
                         = mul_F(-u,0)
                         = mul_F(u,0)
                         = no\_result2_F(u, x)
                                                                 otherwise
```

Argument angular-unit angle from Cartesian co-ordinates 5.3.9.15

The $arcu_F^*$ approximation helper function:

```
arcu_F^*: F \times F \times F \to \mathcal{R}
```

 $arcu_F^*(u,x,y)$ returns a close approximation to $arc(x,y) \cdot u/(2 \cdot \pi)$ in \mathcal{R} , with maximum error $max_error_tanu_F(u)$.

Further requirements on the $arcu_F^*$ approximation helper function are:

```
arcu_F^*(u, x, x) = u/8
                                                   if u, x \in F and x > 0
arcu_F^*(u, 0, y) = u/4
                                                   if u, y \in F and y > 0
arcu_F^*(u, x, -x) = 3 \cdot u/8
                                                   if u, x \in F and x < 0
arcu_F^*(u, x, 0) = u/2
                                                   if u, x \in F and x < 0
arcu_F^*(u, x, -y) = -arcu_F^*(u, x, y)
                                                   if u, x, y \in F and (y \neq 0 \text{ or } x > 0)
arcu_F^*(-u, x, y) = -arcu_F^*(u, x, y)
```

The $arcu_F^{\#}$ range limitation helper function (for $u, x, y \in F$):

```
arcu_F^{\#}(u, x, y) = \max\{up_F(-|u/2|), \min\{arcu_F^{*}(u, x, y), down_F(|u/2|)\}\}
```

The $arcu_F$ operation:

```
arcu_F : F \times F \times F \to F \cup \{-0, underflow, invalid\}
```

```
arcu_F(u, x, y) = result_F^*(arcu_F^\#(u, x, y), nearest_F)
                                                    if u \in G_F and x, y \in F and (x < 0 \text{ or } y \neq 0)
                                                    if u \in G_F and x \in F and x \ge 0 and y = 0
                    = mul_F(u,0)
                    = down_F(u/2)
                                                    if u \in G_F and x = -\mathbf{0} and y = 0 and u > 0
                                                    if u \in G_F and x = -\mathbf{0} and y = 0 and u < 0
                    = up_F(u/2)
                                                   if u \in G_F and x = -0 and y \in F \cup \{-\infty, +\infty\} and
                    = arcu_F(u, 0, y)
                                                       y \neq 0
                                                   if u \in G_F and y = -\mathbf{0} and x \in F \cup \{-\infty, -\mathbf{0}, +\infty\}
                    = neg_F(arcu_F(u, x, 0))
```

```
if u \in G_F and x = +\infty and y \in F and y \ge 0
= mul_F(0, u)
= mul_F(0, -u)
                             if u \in G_F and x = +\infty and y \in F and y < 0
= nearest_F(u/8)
                             if u \in G_F and x = +\infty and y = +\infty
= nearest_F(u/4)
                             if u \in G_F and x \in F and y = +\infty
= nearest_F(3 \cdot u/8)
                             if u \in G_F and x = -\infty and y = +\infty
= down_F(u/2)
                             if u \in G_F and x = -\infty and y \in F and
                                y \geqslant 0 and u > 0
= up_F(-u/2)
                             if u \in G_F and x = -\infty and y \in F and
                                y < 0 and u > 0
= up_F(u/2)
                             if u \in G_F and x = -\infty and y \in F and
                                y > 0 and u < 0
= down_F(-u/2)
                             if u \in G_F and x = -\infty and y \in F and
                                y \leq 0 and u < 0
                             if u \in G_F and x = -\infty and y = -\infty
= nearest_F(-3 \cdot u/8)
                             if u \in G_F and x \in F and y \leftarrow \infty
= nearest_F(-u/4)
= nearest_F(-u/8)
                             if u \in G_F and x = +\infty and y = -\infty
                             otherwise
= no\_result \Im_F(u, x, y)
```

NOTE – Note that the arc operations do *not* return an **invalid** notification at the origin (both second and third arguments in $\{-0,0\}$). See B.5.3.8 and B.5.3.9. Bindings may choose to alter this behaviour.

5.3.10 Operations for angular-unit conversions

5.3.10.1 Converting radian angle to argument angular-unit angle

Define the mathematical function:

```
rad\_to\_cycle : \mathcal{R} \times \mathcal{R} \to \mathcal{R} rad\_to\_cycle(x, w) = \arccos(\cos(x)) \cdot w/(2 \cdot \pi) if \sin(x) \ge 0 and w \ne 0 = \arccos(\cos(x)) \cdot w/(2 \cdot \pi) if \sin(x) < 0 and w \ne 0
```

The $rad_to_cycle_F^*$ approximation helper function:

$$rad_to_cycle_F^* : F^{2\pi} \times F \to \mathcal{R}$$

 $rad_to_cycle_F^*(x, w)$ returns a close approximation to $rad_to_cycle(x, w)$ in \mathcal{R} , with maximum error $max_error_rad_{\mathcal{R}}$, if $|x| \leq big_angle_r_F$.

Further requirements on the $rad_to_cycle_F^*$ approximation helper function are (for $w \in F$):

```
 \begin{array}{ll} rad\_to\_cycle_F^*(n \cdot 2 \cdot \pi + \pi/6, w) = w/12 & \text{if } n \in \mathcal{Z} \text{ and } |n \cdot 2 \cdot \pi + \pi/6| \leqslant big\_angle\_r_F \\ rad\_to\_cycle_F^*(n \cdot 2 \cdot \pi + \pi/4, w) = w/8 & \text{if } n \in \mathcal{Z} \text{ and } |n \cdot 2 \cdot \pi + \pi/4| \leqslant big\_angle\_r_F \\ rad\_to\_cycle_F^*(n \cdot 2 \cdot \pi + \pi/3, w) = w/6 & \text{if } n \in \mathcal{Z} \text{ and } |n \cdot 2 \cdot \pi + \pi/3| \leqslant big\_angle\_r_F \\ rad\_to\_cycle_F^*(n \cdot 2 \cdot \pi + \pi/2, w) = w/4 & \text{if } n \in \mathcal{Z} \text{ and } |n \cdot 2 \cdot \pi + \pi/2| \leqslant big\_angle\_r_F \\ rad\_to\_cycle_F^*(n \cdot 2 \cdot \pi + 2 \cdot \pi/3, w) = w/3 & \text{if } n \in \mathcal{Z} \text{ and } |n \cdot 2 \cdot \pi + 2 \cdot \pi/3| \leqslant big\_angle\_r_F \\ rad\_to\_cycle_F^*(n \cdot 2 \cdot \pi + 3 \cdot \pi/4, w) = 3 \cdot w/8 & \text{if } n \in \mathcal{Z} \text{ and } |n \cdot 2 \cdot \pi + 3 \cdot \pi/4| \leqslant big\_angle\_r_F \\ \end{array}
```

```
rad\_to\_cycle_F^*(n \cdot 2 \cdot \pi + 5 \cdot \pi/6, w) = 5 \cdot w/12
                                                           if n \in \mathcal{Z} and |n \cdot 2 \cdot \pi + 5 \cdot \pi/6| \leq big\_angle\_r_F
   rad\_to\_cycle_F^*(n \cdot 2 \cdot \pi + \pi, w) = w/2
                                                           if n \in \mathcal{Z} and |n \cdot 2 \cdot \pi + \pi| \leq big\_angle\_r_F
   rad\_to\_cycle_F^*(-x, w) = -rad\_to\_cycle_F^*(x, w)
                                                           if x \in F^{2 \cdot \pi} and rad\_to\_cycle(x, w) \neq w/2
   rad\_to\_cycle_F^*(x, -w) = -rad\_to\_cycle_F^*(x, w)
                                                           if x \in F^{2 \cdot \pi} and rad\_to\_cycle(x, w) \neq w/2
The rad\_to\_cycle_F operation:
   rad\_to\_cycle_F : F \times F \to F \cup \{underflow, absolute\_precision\_underflow, invalid\}
   rad\_to\_cycle_F(x, w)
                         = result_F^*(rad\_to\_cycle_F^*(x, w), nearest_F)
                                                           if w \in G_F and x \in F and |x| \leq big\_argle\_r_F and
                                                           x \neq 0 if w \in G_F and x \in \{-\mathbf{0}, 0\}
                         = mul_F(w, x)
                         = absolute_precision_underflow(qNaN)
                                                           if w \in G_F and x \in F and x \in F and x \in F
                         = no\_result2_F(x, w)
                                                           otherwise
```

5.3.10.2 Converting argument angular-unit angle to radian angle

Define the mathematical function:

```
cycle\_to\_rad : \mathcal{R} \times \mathcal{R} \to \mathcal{R}
cycle\_to\_rad(u, x) = \arccos(\cos(x \cdot 2 \cdot \pi/u)) \quad \text{if } \sin(x \cdot 2 \cdot \pi/u) \geqslant 0
= -\arccos(\cos(x \cdot 2 \cdot \pi/u)) \quad \text{if } \sin(x \cdot 2 \cdot \pi/u) < 0
```

The $cycle_to_rad_F^*$ approximation helper function:

```
cycle\_to\_rad_F^*: (u:F) \times F^u \to \mathcal{R}
```

 $cycle_to_rad_F^*(u, x)$ returns a close approximation to $cycle_to_rad(u, x)$ in \mathcal{R} , if $u \neq 0$, with maximum error $max_error_rad_F$.

Further requirements on the $vicle_to_rad_F^*$ approximation helper function are (for $u \in F$, $u \neq 0$):

```
cycle\_to\_rad_F^*(u,n\cdot u+x) = cycle\_to\_rad_F^*(u,x) if n\in\mathcal{Z} and x\in F^u cycle\_to\_rad_F^*(u,u/12) = \pi/6 cycle\_to\_rad_F^*(u,u/8) = \pi/4 cycle\_to\_rad_F^*(u,u/6) = \pi/3 cycle\_to\_rad_F^*(u,u/4) = \pi/2 cycle\_to\_rad_F^*(u,u/3) = 2\cdot \pi/3 cycle\_to\_rad_F^*(u,3\cdot u/8) = 3\cdot \pi/4 cycle\_to\_rad_F^*(u,5\cdot u/12) = 5\cdot \pi/6 cycle\_to\_rad_F^*(u,u/2) = \pi cycle\_to\_rad_F^
```

The $cycle_to_rad_F$ operation:

```
\begin{aligned} & cycle\_to\_rad_F: F \times F \to F \cup \{\textbf{-0}, \textbf{underflow}, \textbf{absolute\_precision\_underflow}, \textbf{invalid}\} \\ & cycle\_to\_rad_F(u, x) \\ & = result_F^*(cycle\_to\_rad_F^*(u, x), nearest_F) \\ & \qquad \qquad & \text{if } cycle_F(u, x) \in F \text{ and } cycle_F(u, x) \neq 0 \\ & = mul_F(cycle_F(u, x), u) & \text{if } cycle_F(u, x) \in \{\textbf{-0}, 0\} \\ & = cycle_F(u, x) & \text{otherwise} \end{aligned}
```

5.3.10.3 Converting argument angular-unit angle to (another) argument angular-unit angle

Define the mathematical function:

```
\begin{aligned} cycle\_to\_cycle : \mathcal{R} \times \mathcal{R} \times \mathcal{R} \to \mathcal{R} \\ cycle\_to\_cycle(u, x, w) \\ &= \arccos(\cos(x \cdot 2 \cdot \pi/u)) \cdot w/(2 \cdot \pi) \\ &\quad \text{if } u \neq 0 \text{ and } w \neq 0 \text{ and } \sin(x \cdot 2 \cdot \pi/u) \geqslant 0 \\ &= -\arccos(\cos(x \cdot 2 \cdot \pi/u)) \cdot w/(2 \cdot \pi) \\ &\quad \text{if } u \neq 0 \text{ and } w \neq 0 \text{ and } \sin(x \cdot 2 \cdot \pi/u) < 0 \end{aligned}
```

The $cycle_to_cycle_F^*$ approximation helper function:

```
cycle\_to\_cycle_F^*: (u:F) \times F^u \times F \to \mathcal{R}
```

 $cycle_to_cycle_F^*(u, x, w)$ returns a close approximation $cycle_to_cycle(u, x, w)$ in \mathcal{R} if $u \neq 0$ and $|x/u| \leq big_angle_u_F$, with maximum error $max_error_rad_F$.

Further requirements on the $cycle_to_cycle_F^*$ approximation helper function are (for $u, w \in F$, $u \neq 0$):

```
cycle\_to\_cycle_F^*(u, n \cdot u + x, w) = cycle\_to\_cycle_F^*(u, x, w)
                                                    if n \in \mathcal{Z} and x \in F^u
cycle\_to\_cycle_F^*(u, u/12, w) = w/12
cycle\_to\_cycle_F^*(u, u/8, w) = w/8
cycle\_to\_cycle_F^*(u, u/6, w) = w/6
cycle\_to\_cycle_F^*(u, u/4, w) = w/4
cycle\_to\_cycle_F^*(u,u)3,w) = w/3
cycle\_to\_cycle_E^*(u_*) \cdot u/8, w) = 3 \cdot w/8
cycle\_to\_cycle^*_F(u, 5 \cdot u/12, w) = 5 \cdot w/12
cycle\_to\_cycle_F^*(u, u/2, w) = w/2
cycle\_to\_cycle_F^*(u, -x, w) = -cycle\_to\_cycle_F^*(u, x, w)
                                                    if x \in F^u and cycle\_to\_cycle(u, x, w) \neq w/2
cuele\_to\_cycle_F^*(-u, x, w) = -cycle\_to\_cycle_F^*(u, x, w)
                                                    if x \in F^u and cycle\_to\_cycle(u, x, w) \neq w/2
cycle\_to\_cycle_F^*(u,x,-w) = -cycle\_to\_cycle_F^*(u,x,w)
                                                    if x \in F^u and cycle\_to\_cycle(u, x, w) \neq w/2
```

The $cycle_to_cycle_F$ operation:

```
cycle\_to\_cycle_F: F \times F \times F \to F \cup \{\textbf{-0}, \textbf{underflow}, \textbf{absolute\_precision\_underflow}, \\ \textbf{invalid}\}
```

```
cycle\_to\_cycle_F(u,x,w)
= result_F^*(cycle\_to\_cycle_F^*(u,x,w), nearest_F)
\text{if } w \in G_F \text{ and } cycle_F(u,x) \in F \text{ and } cycle_F(u,x) \neq 0
= mul_F(w, cycle_F(u,x)) \quad \text{if } w \in G_F \text{ and } cycle_F(u,x) \in \{-\mathbf{0},0\}
= \mathbf{absolute\_precision\_underflow}(\mathbf{qNaN})
\text{if } w \in G_F \text{ and }
cycle_F(u,x) = \mathbf{absolute\_precision\_underflow}
= no\_result\beta_F(u,x,w) \quad \text{otherwise}
```

5.3.11 Operations for hyperbolic elementary functions

There shall be two maximum error parameters for operations corresponding to the hyperbolic and inverse hyperbolic functions:

```
max\_error\_sinh_F \in F

max\_error\_tanh_F \in F
```

The $max_error_sinh_F$ parameter shall have a value that is $\leq 2 \cdot rnd_error_F$. The $max_error_tanh_F$ parameter shall have a value that is $\leq 2 \cdot rnd_error_F$.

5.3.11.1 Hyperbolic sine

The $sinh_F^*$ approximation helper function:

$$sinh_F^*: F \to \mathcal{R}$$

 $sinh_F^*(x)$ returns a close approximation to sinh(x) in \mathcal{R} , with maximum error $max_error_sinh_F$. Further requirements on the $sinh_F^*$ approximation helper function are:

```
sinh_F^*(x) = x
sinh_F^*(-x) = -sinh_F^*(x)
sinh_F^*(-x) = -sinh_F^*(x)
if x \in F and sinh_F^*(x) \neq sinh(x) and |x| < \sqrt{2 \cdot epsilon_F/r_F}
if x \in F
```

The $sinh_F$ operation:

$$\begin{aligned} sinh_F: F \to F \cup \{ \mathbf{overflow} \} \\ sinh_F(x) &= result_F^*(sinh_F^*(x), nearest_F) \\ & \text{if } x \in F \text{ and } |x| > fminN_F \\ & \text{if } x \in F \text{ and } |x| \leqslant fminN_F \\ & \text{if } x \in \{-\infty, -\mathbf{0}, +\infty\} \\ &= no_result_F(x) \end{aligned}$$

NOTES

- 1 underflow is explicitly avoided.
- 2 $sinh_F(x)$ will overflow approximately when $|x| > \ln(2 \cdot fmax_F)$.

5.3.11.2 Hyperbolic cosine

The $cosh_F^*$ approximation helper function:

$$cosh_F^*: F \to \mathcal{R}$$

 $cosh_F^*(x)$ returns a close approximation to cosh(x) in \mathcal{R} , with maximum error $max_error_sinh_F$.

Further requirements on the $cosh_F^*$ approximation helper function are:

$$\cosh_F^*(x) = 1$$
 if $x \in F$ and $\cosh_F^*(x) \neq \cosh(x)$ and $|x| < \sqrt{epsilon_F}$ $\cosh_F^*(-x) = \cosh_F^*(x)$ if $x \in F$

The relationship to the $sinh_F^*$ approximation helper function for the $sinh_F$ operation in the same library shall be:

$$cosh_F^*(x) \geqslant sinh_F^*(x)$$
 if $x \in F$

The $cosh_F$ operation:

$$cosh_F: F \to F \cup \{\mathbf{overflow}\}$$

$$cosh_F(x) = result_F^*(cosh_F^*(x), nearest_F)$$

$$= 1 \qquad \qquad \text{if } x \in F$$

$$= 1 \qquad \qquad \text{if } x = -\mathbf{0}$$

$$= +\infty \qquad \qquad \text{if } x \in \{-\infty, +\infty\}$$

$$= no_result_F(x) \qquad \text{otherwise}$$

NOTE – $cosh_F(x)$ will overflow approximately when $|x| > \ln(2 \cdot max)$ PDFofis

Hyperbolic tangent 5.3.11.3

The $tanh_F^*$ approximation helper function:

$$tanh_F^*: F \to \mathcal{R}$$

 $tanh_F^*(x)$ returns a close approximation to tanh(x) in \mathcal{R} , with maximum error $max_error_tanh_F$.

Further requirements on the $tanh_F^*$ approximation helper function are:

$$tanh_F^*(x) = x$$
 if $x \in F$ and $tanh_F^*(x) \neq tanh(x)$ and
$$|x| \leq \sqrt{1.5 \cdot epsilon_F/r_F}$$
 if $x \in F$ and $tanh_F^*(x) \neq tanh(x)$ and
$$x > \operatorname{arctanh}(1 - (epsilon_F/(3 \cdot r_F)))$$
 if $x \in F$
The $tanh_F$ operation

$$tanh_F(x)$$
 = $result_F^*(tanh_F^*(x), nearest_F)$
 $= x$ if $x \in F$ and $|x| > fminN_F$
 $= -\mathbf{0}$ if $x = -\mathbf{0}$
 $= -1$ if $x = -\infty$
 $= 1$ if $x = +\infty$
 $= no_result_F(x)$ otherwise

NOTE - **underflow** is explicitly avoided.

5.3.11.4 Hyperbolic cotangent

The $coth_F^*$ approximation helper function:

$$coth_F^*: F \to \mathcal{R}$$

 $coth_F^*(x)$ returns a close approximation to coth(x) in \mathcal{R} , with maximum error $max_error_tanh_F$.

Further requirements on the $coth_F^*$ approximation helper function are:

$$coth_F^*(x) = 1$$
 if $x \in F$ and $coth_F^*(x) \neq coth(x)$ and $x > \operatorname{arccoth}(1 + (epsilon_F/4))$ if $x \in F$

The $coth_F$ operation:

$$coth_F: F \to F \cup \{ \mathbf{infinitary}, \mathbf{overflow} \}$$

$$coth_F(x) = result_F^*(coth_F^*(x), nearest_F)$$

$$= \mathbf{infinitary}(+\infty) \qquad \text{if } x \in F \text{ and } x \neq 0$$

$$= \mathbf{infinitary}(-\infty) \qquad \text{if } x = -\mathbf{0}$$

$$= -1 \qquad \text{if } x = -\infty$$

$$= 1 \qquad \text{if } x = +\infty$$

$$= no_result_F(x) \qquad \text{otherwise}$$

NOTE – $coth_F(x)$ will overflow approximately when $|1/x| \geqslant fmax_F$. the full PD

5.3.11.5 Hyperbolic secant

The $sech_F^*$ approximation helper function:

$$sech_F^*: F \to \mathcal{R}$$

 $sech_F^*(x)$ returns a close approximation $\operatorname{tosech}(x)$ in \mathcal{R} , with maximum error $\operatorname{max_error_tanh}_F$.

Further requirements on the $sech_{\kappa}^*$ approximation helper function are:

```
sech_F^*(x) = 1
                                                                         if x \in F and \operatorname{sech}_F^*(x) \neq \operatorname{sech}(x) and |x| < \sqrt{\operatorname{epsilon}_F/r_F}
    sech_F^*(-x) = sech_F^*(x)
sech_F^*(x) < fminD_F/2
                                                                         if x \in F and x > 2 - \ln(fminD_F/4)
The sech_F operation.
    sech_F: F \longrightarrow F \cup \{\mathbf{underflow}\}\
                               = result_F^*(sech_F^*(x), nearest_F)
                                                                         if x \in F
                                                                         if x = -0
                                                                         if x \in \{-\infty, +\infty\}
                               = no\_result_F(x)
                                                                         otherwise
```

5.3.11.6 Hyperbolic cosecant

The $csch_F^*$ approximation helper function:

$$csch_F^*: F \to \mathcal{R}$$

 $csch_F^*(x)$ returns a close approximation to csch(x) in \mathcal{R} , with maximum error $max_error_tanh_F$.

Further requirements on the $csch_F^*$ approximation helper function are:

$$csch_F^*(-x) = -csch_F^*(x)$$
 if $x \in F$
 $csch_F^*(x) < fminD_F/2$ if $x \in F$ and $x > 2 - \ln(fminD_F/4)$

The relationship to the $sech_F^*$ approximation helper function for the $sech_F$ operation in the same library shall be:

$$csch_F^*(x) \geqslant sech_F^*(x)$$
 if $x \in F$ and $x > 0$

The $csch_F$ operation:

$$csch_F: F \to F \cup \{ ext{underflow}, ext{overflow}, ext{infinitary} \}$$

$$csch_F(x) = result_F^*(csch_F^*(x), nearest_F)$$

$$= div_F(1, x) \qquad \text{if } x \in F \text{ and } x \neq 0$$

$$= div_F(1, x) \qquad \text{if } x \in \{-\infty, -\mathbf{0}, 0\} + \infty \}$$

$$= no_result_F(x) \qquad \text{otherwise}$$

NOTE – $csch_F(x)$ will overflow approximately when $|x| > fmax_F$.

5.3.11.7 Inverse hyperbolic sine

The $arcsinh_F^*$ approximation helper function

$$arcsinh_F^*: F \to \mathcal{R}$$

 $arcsinh_F^*(x)$ returns a close approximation to arcsinh(x) in \mathcal{R} , with maximum error $max_error_sinh_F$.

Further requirements on the $arcsinh_F^*$ approximation helper function are:

$$arcsinh_F^*(x) = x \qquad \text{if } x \in F \text{ and } arcsinh_F^*(x) \neq \arcsin(x) \text{ and } |x| \leqslant \sqrt{3 \cdot epsilon_F/r_F}$$

$$arcsinh_F^*(-x) = arcsinh_F^*(x) \qquad \text{if } x \in F$$

The $arcsinh_F$ operation:

$$\begin{array}{ll} \operatorname{arcsinh}_F : F \to F \\ \operatorname{arcsinh}_F(x) &= \operatorname{result}_F^*(\operatorname{arcsinh}_F^*(x), \operatorname{nearest}_F) \\ &= \operatorname{if} \ x \in F \ \operatorname{and} \ |x| > \operatorname{fmin} N_F \\ &= x \qquad \qquad \operatorname{if} \ x \in F \ \operatorname{and} \ |x| \leqslant \operatorname{fmin} N_F \\ &= x \qquad \qquad \operatorname{if} \ x \in \{-\infty, -\mathbf{0}, +\infty\} \\ &= \operatorname{no_result}_F(x) \qquad \operatorname{otherwise} \end{array}$$

NOTE - **underflow** is explicitly avoided.

5.3.11.8 Inverse hyperbolic cosine

The $arccosh_F^*$ approximation helper function:

$$arccosh_F^*: F \to \mathcal{R}$$

 $arccosh_F^*(x)$ returns a close approximation to arccosh(x) in \mathcal{R} , with maximum error $max_error_sinh_F$.

The relationship to the $arcsinh_F^*$ approximation helper function for the $arcsinh_F$ operation in the same library shall be:

$$arccosh_F^*(x) \leqslant arcsinh_F^*(x)$$
 if $x \in F$

The $arccosh_F$ operation:

$$arccosh_F: F \to F \cup \{\mathbf{invalid}\}$$

 $arccosh_F(x) = result_F^*(arccosh_F^*(x), nearest_F)$
 $if \ x \in F \text{ and } x \geqslant 1$
 $= +\infty$
 $= no_result_F(x)$ otherwise

Inverse hyperbolic tangent 5.3.11.9

The $arctanh_F^*$ approximation helper function:

$$arctanh_F^*: F \to \mathcal{R}$$

 $arctanh_F^*(x)$ returns a close approximation to arctanh(x) in \mathbb{R} , with maximum error $max_error_tanh_F$.

Further requirements on the $arctanh_F^*$ approximation helper function are:

arctanh
$$_F^*(x) = x$$
 if $x \in F$ and $arctanh_F^*(x) \neq arctanh(x)$ and $arctanh_F^*(-x) = -arctanh_F^*(x)$ if $x \in F$ if $x \in F$ and $arctanh_F^*(x) \neq arctanh(x)$ and $arctanh_F^*(-x) = -arctanh_F^*(x)$

The $arctanh_F$ operation:

$$arctanh_F: F \to F \cup \{ \mathbf{infinitary, invalid} \}$$
 $arctanh_F(x) = result_F^*(arctanh_F^*(x), nearest_F)$
 $if \ x \in F \ and \ fminN_F < |x| < 1$
 $= x \qquad \qquad if \ x \in F \ and \ |x| \leqslant fminN_F$
 $= 0 \qquad \qquad if \ x = -0$
 $= \mathbf{infinitary}(+\infty) \qquad \text{if } x = 1$
 $= \mathbf{infinitary}(-\infty) \qquad \text{if } x = -1$
 $= no_result_F(x) \qquad \text{otherwise}$

underflow is explicitly avoided.

5.3.11.10 Inverse hyperbolic cotangent

The $arccoth_F^*$ approximation helper function:

$$arccoth_F^*: F \to \mathcal{R}$$

 $arccoth_F^*(x)$ returns a close approximation to arccoth(x) in \mathcal{R} , with maximum error $max_error_tanh_F$.

A further requirement on the $arccoth_F^*$ approximation helper function is:

$$arccoth_F^*(-x) = -arccoth_F^*(x)$$
 if $x \in F$

The $arccoth_F$ operation:

 $arccoth_F: F \to F \cup \{ underflow, infinitary, invalid \}$

$$arccoth_F(x)$$
 = $result_F^*(arccoth_F^*(x), nearest_F)$
if $x \in F$ and $|x| > 1$
= $infinitary(+\infty)$ if $x = 1$
= $infinitary(-\infty)$ if $x = -1$
= -0 if $x = -\infty$
= 0 if $x = +\infty$
= $no_result_F(x)$ otherwise

NOTE - There is no underflow for this operation for most kinds of floating point types OllEC 10961.2. e.g. IEC 60559 ones.

5.3.11.11 Inverse hyperbolic secant

The $arcsech_F^*$ approximation helper function:

$$arcsech_F^*: F \to \mathcal{R}$$

 $arcsech_F^*(x)$ returns a close approximation to arcsech(x) in \mathcal{R} , with maximum error $max_error_tanh_F$.

The $arcsech_F$ operation:

$$arcsech_F: F \to F \cup \{ \mathbf{infinitary}, \mathbf{invalid} \}$$

$$arcsech_F(x) = result_F^*(arcsech_F^*(x), nearest_F)$$

$$\mathbf{if} \ x \in F \ \text{and} \ 0 < x \leqslant 1$$

$$= \mathbf{infinitary}(+\infty) \quad \text{if} \ x \in \{-\mathbf{0}, 0\}$$

$$= no_result_F(x) \quad \text{otherwise}$$

5.3.11.12 Inverse hyperbolic cosecant

The $arccsch_F^*$ approximation helper function:

$$arccsch_F^*: F \to \mathcal{R}$$

 $arccsch_F^*(x)$ returns a close approximation to arccsch(x) in \mathcal{R} , with maximum error $max_error_tanh_F$.

A further requirement on the $arccsch_F^*$ approximation helper function is:

$$arccsch_F^*(x) = -arccsch_F^*(x)$$
 if $x \in F$

 $= no_result_F(x)$

The $arccsch_F$ operation:

$$arccsch_F: F \to F \cup \{\mathbf{underflow}, \mathbf{infinitary}\}\$$
 $arccsch_F(x) = result_F^*(arccsch_F^*(x), nearest_F)$ if $x \in F$ and $x \neq 0$ $= div_F(1, x)$ if $x \in \{-\infty, -\mathbf{0}, 0, +\infty\}$

NOTE - There is no underflow for this operation for most kinds of floating point types, e.g. IEC 60559 ones.

otherwise

5.4 Operations for conversion between numeric datatypes

Numeric conversion between different representation forms for integer and fractional values can take place under a number of different circumstances. E.g.:

- a) explicit or implicit conversion between different numeric datatypes conforming to part 1;
- b) explicit or implicit conversion between different numeric datatypes only one of which conforms to part 1;
- c) explicit or implicit conversion between a character string and a numeric datatype.

The latter includes outputting a numeric value as a character string, inputting a numeric value from a character string source, and converting a numeral in the source program to a value in a numeric datatype (see clause 5.5). This part covers only the cases where at least one of the source and target is a numeric datatype conforming to part 1.

When a character string is involved as either source or target of a conversion, this part does not specify the lexical syntax for the numerals parsed or formed. A binding standard should specify the lexical syntax or syntaxes for these numerals, and, when appropriate, how the lexical syntax for the numerals can be altered. This could include which script for the digits to use in a position system (Latin-Arabic digits, Arabic-Indic digits, traditional Thai digits, etc.). With the exception of the radix used in numerals expressing fractional values, differences in lexical syntactic details that do not affect the value in $\mathcal R$ denoted by the numerals should not affect the result of the conversion.

Character string representations for integer values can include representations for -0, $+\infty$, $-\infty$, and quiet NaNs. Character string representations for floating point and fixed point values should have formats for -0, $+\infty$, $-\infty$, and quiet NaNs. For both integer and floating point values, character strings that are not numerals nor special values according to the lexical syntax used, shall be regarded as signalling NaNs when used as source of a numerical conversion.

For the cases where one of the datatypes involved in the conversion does not conform to part 1, the values of some numeric datatype parameters need to be inferred. For integers, one need to infer the value for bounded, and if that is **true** then also values for maxint and minint, and for string formats also the radix. For floating point values, one need to infer the values for r, p, and emax or emin. In case a precise determination is not possible, values that are 'safe' for that instance should be used. 'Safe' values for otherwise undetermined inferred parameters are such that

- a) monotonicity of the conversion function is not affected,
- b) the error in the conversion does not exceed that specified by the maximum error parameter (see below).
- c) if the value resulting from the conversion is converted back to the source datatype by a conversion conforming to this part, the original value should be regenerated if possible, and
- d) overflow and underflow are avoided if possible.

If, and only if, a specified infinite special value result cannot be represented in the target datatype, the infinity result shall be interpreted as implying the **infinitary** notification. If, and only if, a specified NaN special value result cannot be represented in the target datatype, the NaN result shall be interpreted as implying the **invalid** notification. If, and only if, a specified

-0 special value result cannot be represented in the target datatype, the -0 result shall be interpreted as 0.

5.4.1 Integer to integer conversions

Let I and I' be non-special value sets for integer datatypes. At least one of the datatypes corresponding to I and I' conforms to part 1.

```
convert_{I \to I'} : I \to I' \cup \{ \mathbf{overflow} \}
convert_{I \to I'}(x) = result_{I'}(x) \qquad \text{if } x \in I
= x \qquad \qquad \text{if } x \in \{ -\infty, -0, +\infty \}
= \mathbf{qNaN} \qquad \qquad \text{if } x \text{ is a quiet NaN}
= \mathbf{invalid}(\mathbf{qNaN}) \qquad \qquad \text{if } x \text{ is a signalling NaN}
```

NOTE – If both I and I' are conforming to part 1, then this conversion is covered by part 1. This operation generalises the $cvt_{I \to I'}$ of part 1, since only one of the integer datatypes in the conversion need be conforming to part 1.

5.4.2 Floating point to integer conversions

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.

```
floor_{F \to I} : F \to I \cup \{\mathbf{overflow}\}\
floor_{F \to I}(x) = result_I(\lfloor x \rfloor)
                                                         if x \in \{-\infty, -0, +\infty\}
                                                       \bigcircif x is a quiet NaN
                                                         if x is a signalling NaN
rounding_{F \to I} : F \to I \cup \{-0, \mathbf{overflow}\}\
rounding_{F \to I}(x)
                                                         if x \in F and (x \ge 0 \text{ or round}(x) \ne 0)
                                                         if x \in F and x < 0 and round(x) = 0
                                                         if x \in \{-\infty, -0, +\infty\}
                                                         if x is a quiet NaN
                                                         if x is a signalling NaN
                                                         if x \in F and (x \ge 0 \text{ or } \lceil x \rceil \ne 0)
                                                         if x \in F and x < 0 and \lceil x \rceil = 0
                                                         if x \in \{-\infty, -0, +\infty\}
                      = qNaN
                                                         if x is a quiet NaN
                      = invalid(qNaN)
                                                         if x is a signalling NaN
```

Integer to floating point conversions

Let I be the non-special value set for an integer datatype. Let F be the non-special value set for a floating point datatype. At least one of the source and target datatypes is conforming to part 1.

```
convert_{I \to F} : I \to F \cup \{\mathbf{overflow}\}\
convert_{I \to F}(x) = result_F(x, nearest_F) if x \in I
                                                 if x \in \{-\infty, -0, +\infty\}
                   = qNaN
                                                 if x is a quiet NaN
                   = invalid(qNaN)
                                                 if x is a signalling NaN
```

NOTE – When both I and F conform to part 1, integer to nearest floating point conversions are covered by part 1. In this case the operations $cvt_{I\to F}$ and $convert_{I\to F}$ are identical,

5.4.4 Floating point to floating point conversions

Define the least radix function, lb, defined for arguments that are greater than 0:

```
lb: \mathcal{Z} \to \mathcal{Z}
lb(r) = \min\{n \in \mathcal{Z} \mid n \geqslant 1 \text{ and there is an } m \in \mathcal{Z} \text{ such that } r \neq m''\}
```

Let F, F', and F'' be non-special value sets for floating point datatypes. At least one of the source and target datatypes in the conversion conforms to part 1.

There shall be a $max_error_convert_{F'}$ parameter that gives the maximum error when converting from F to F' and $lb(r_F) \neq lb(r_{F'})$. The max_error convert_{F'} parameter shall have the value 0.5. If the binding standard requires that this parameter has the value 0.5 (see annex A), this parameter need not be made available for programs.

If $lb(r_F) = lb(r_{F'})$, the maximum error shall be 0.5 ulp when converting from F to F', even when the implementation is only partially conforming (see Annex A), but this is not reflected in any parameter.

The $convert_{F\to F'}^*$ approximation helper functions:

 $convert_{F \to F'}^* : F \to \mathcal{R}$ $convert_{F \to F'}^* (x)$ returns a close approximation to x in \mathcal{R} , with maximum error $max_error_convert_{F'}$.

NOTE 1 – With the maximum error 0.5 ulp, this and the below conversion helper functions are not really needed. However, Annex A allows for partial conformity, such that the maximum error for these helper functions may be greater than 0.5 ulp.

Further requirements on the $convert^*_{F\to F'}$ approximation helper functions are:

```
\begin{array}{ll} convert_{F \to F'}^*(x) = x & \text{if } x \in \mathcal{Z} \cap F \\ convert_{F \to F'}^*(x) > 0 & \text{if } x \in F \text{ and } x > 0 \\ convert_{F \to F'}^*(-x) = -convert_{F \to F'}^*(x) & \text{if } x \in F \\ convert_{F \to F'}^*(x) \leqslant convert_{F \to F'}^*(y) & \text{if } x, y \in F \text{ and } x < y \end{array}
```

Relationship to other floating point to floating point conversion approximation helper functions for conversion operations in the same library shall be:

```
convert^*_{F \to F'}(x) = convert^*_{F'' \to F'}(x) if lb(r_{F''}) = lb(r_F) and x \in F \cap F''
```

The $convert_{F \to F'}$ operation:

```
convert_{F \to F'} : F \to F' \cup \{ overflow, underflow \}
```

```
convert_{F 	o F'}(x) = result_{F'}(x, nearest_{F'}) if x \in F and lb(r_F) = lb(r_{F'})

= result_{F'}^*(convert_{F 	o F'}^*(x), nearest_{F'}) if x \in F and lb(r_F) \neq lb(r_{F'})

= x if x \in \{-\infty, -0, +\infty\}

= \mathbf{qNaN} if x is a quiet NaN

= \mathbf{invalid}(\mathbf{qNaN}) if x is a signalling NaN
```

NOTES

- 2 Modern techniques allow, on the average, efficient conversion with a maximum error of 0.5 ulp even when the radices differ. C99 [17], for instance, requires that all floating point value conversion is done with a maximum error of 0.5 ulp.
- 3 IEC 60559 requirements imply that the $max_error_convert_{F'}$ parameter has a value ≤ 0.97 . Such a large maximum error for the conversion is only partially conforming. See Annex A.
- 4 When the maximum error is 0.5, the conversion helper function above can be the identity function.
- 5 When both datatypes conform to part 1, and the radices for both of these floating point datatypes are the same, floating point to nearest floating point conversions are covered by part 1. In this case the operations $cvt_{F\to F'}$ and $convert_{F\to F'}$ are identical.

5.4.5 Floating point to fixed point conversions

Let F be the non-special value set for a floating point datatype conforming to part 1. Let D be the non-special value set for a fixed point datatype.

A fixed point datatype D is a subset of \mathcal{R} , characterised by a radix, $r_D \in \mathcal{Z} \ (\geqslant 2)$, a density, $d_D \in \mathcal{Z} \ (\geqslant 0)$, and if it is bounded, a maximum positive value, $dmax_D \in D^* \ (\geqslant 1)$. Given these values, the following sets are defined:

$$D^* = \{n/(r_D^{d_D}) \mid n \in Z\}$$

$$D = D^*$$

$$D = D^* \cap [-dmax_D, dmax_D]$$
if D is not bounded if D is bounded

NOTE 1 – D corresponds to scaled (r_D, d_D) in ISO/IEC 11404 Language independent datatypes (LID) [10]. LID has no parameter corresponding to $dmax_D$ even when the datatype is bounded.

The fixed point rounding helper function:

```
nearest_D : \mathcal{R} \leftrightarrow \mathcal{D}^*
```

is the rounding function that rounds to nearest, ties round to even last digit.

The fixed point result helper function, $result_D$, is like $result_F$, but for a fixed point datatype. It will return **overflow** if the rounded result is not representable:

```
result_D: \mathcal{R} \times (\mathcal{R} \to D^*) \to D \cup \{\mathbf{-0}, \mathbf{overflow}\}
result_D(x, rnd) = rnd(x) \qquad \text{if } rnd(x) \in D \text{ and } (rnd(x) \neq 0 \text{ or } x \geqslant 0)
= \mathbf{-0} \qquad \text{if } rnd(x) = 0 \text{ and } x < 0
= \mathbf{overflow} \qquad \text{if } x \in \mathcal{R} \text{ and } rnd(x) \notin D
```

There shall be a $max_error_convert_D$ parameter that gives the maximum error when converting from F to D and $lb(r_F) \neq lb(r_D)$. The $max_error_convert_D$ parameter shall have the value 0.5. If the binding standard requires that this parameter has the value 0.5 (see annex A), this parameter need not be made available for programs.

If $lb(r_F) = lb(r_D)$, the maximum error shall be 0.5 ulp when converting from F to D, even when the implementation is only partially conforming (see Annex A), but this is not reflected in any parameter.

The $convert_{F\to D}^*$ approximation helper function:

```
convert^*_{F \to D} : F \to \mathcal{R}
```

 $convert_{F\to D}^*(x)$ returns a close approximation to x in \mathcal{R} , with maximum error $max_error_convert_D$.

Further requirements on the $convert_{F\to D}^*$ approximation helper functions are:

```
if x \in \mathcal{Z} \cap F
convert_{F\to D}^*(x) = x
                                                                                       if x \in F and x > 0
convert^*_{F \to D}(x) > 0
convert^*_{F \to D}(x) = -convert^*_{F \to D}(x) \qquad \text{if } x \in F
convert^*_{F \to D}(x) \leqslant convert^*_{F \to D}(y) \qquad \text{if } x, y \in F \text{ and } x < y
```

Relationship to other floating point to fixed point conversion approximation helper functions for conversion operations in the same library shall be:

```
if lb(r_{F''}) = lb(r_F) and x \in FF''
convert_{F\to D}^*(x) = convert_{F''\to D}^*(x)
```

The $convert_{F\to D}$ operation:

```
convert_{F\to D}: F\to D\cup \{-0, \mathbf{overflow}\}
convert_{F\to D}(x) = result_D(x, nearest_D) if x \in F and lb(r_F)
                  = result_D(convert^*_{F \to D}(x), nearest_D)
                                                if x \in F and b(r_F) \neq lb(r_D)
                                                if x \in \{-\infty, -0, +\infty\}
                                                if x is a quiet NaN
                   = qNaN
                   = invalid(qNaN)
                                                if x is a signalling NaN
```

NOTES

- 2 The datatype D need not be visible in the programming language. D may be a subtype of strings, according to some format. Even so, no datatype for strings need be present in the programming language.
- 3 This covers, among other things O'output" of floating point datatype values, to fixed point string formats. E.g. a binding may say that float_to_fixed_string(x, m, n) is bound to $convert_{F\to S_{m,n}}(x)$ where $S_{m,n}$ is strings of length m, representing fixed point values in radix 10 with n decimals! The binding should also detail how NaNs, signed zeroes and infinities are represented in $S_{m,n}$, as well as the precise format of the strings representing ordinary values. (Note that if the length of the target string is limited, the conversion may overflow.)
- 4 IEC 60559 requirements imply that the $max_error_convert_D$ parameter has a value ≤ 0.97 . Such a large maximum error for the conversion is only partially conforming. See Annex A.
- 5 When the maximum error is 0.5, the conversion helper function above can be the identity function.

Fixed point to floating point conversions

Let F be the non-special value set for a floating point datatype conforming to part 1. Let D and D' be the non-special value set for fixed point datatypes.

The $convert_{D\to F}^*$ approximation helper function:

$$convert_{D \to F}^* : D \to \mathcal{R}$$

 $convert_{D\to F}^*(x)$ returns a close approximation to x in \mathcal{R} , with maximum error $max_error_convert_F$.

Further requirements on the $convert_{D\to F}^*$ approximation helper functions are:

```
convert_{D\to F}^*(x) = x
                                                                      if x \in \mathcal{Z} \cap D
convert_{D\to F}^*(x) > 0
                                                                      if x \in D and x > 0
convert^*_{D \to F}(-x) = -convert^*_{D \to F}(x)convert^*_{D \to F}(x) \leqslant convert^*_{D \to F}(y)
                                                                     if x \in D
                                                                     if x, y \in D and x < y
```

Relationship to other floating point and fixed point to floating point conversion approximation helper functions for conversion operations in the same library shall be:

```
if lb(r_{D'}) = lb(r_D) and x \in D \cap D'
convert_{D \to F}^*(x) = convert_{D' \to F}^*(x)
                                                       if lb(r_{F'}) = lb(r_D) and x \in D \cap F'
convert_{D\to F}^*(x) = convert_{F'\to F}^*(x)
```

The $convert_{D\to F}$ operation:

```
convert_{D \to F}(x) = result_F(x, nearest_F) \quad \text{if } x \in D \text{ and } lb(r_D) = lb(r_F)
= result_F^*(convert_{D \to F}^*(x), nearest_F)
= result_F^*(convert_{D \to F}^*(x), nearest_F)
                                                                       if x \in \{-\infty, -0, +\infty\}
                            = x
                            = qNaN
                                                                       if x is a quiet NaN
                            = invalid(qNaN)
                                                                       if x is a signalling NaN
```

NOTES

- 1 This covers, among other things, "input" of floating point datatype values, from fixed point string formats. E.g. a binding may say that string to float(s) is bound to $convert_{S_{m,n}\to F}(s)$ where $S_{m,n}$ is strings of length m, where m is the length of s, and n is the number of digits after the "decimal symbol" represented in $S_{m,n}$ as well as the precise format of the strings representing ordinary values.
- When the maximum error is 0.5, the conversion helper function above can be the identity function.

Numerals as operations in a programming language

NOTE - Numerals in strings, or input, is covered by the conversion operations in clause 5.4.

Each numeral is a parameterless operation. Thus, this clause introduces a very large number of operations, since the number of numerals is in principle infinite.

Numerals for integer datatypes

Let I' be a non-special value set for integer numerals for the datatype corresponding to I.

An integer numeral, denoting an abstract value n in $I' \cup \{-0, +\infty, -\infty, qNaN, sNaN\}$, for an integer datatype with non-special value set I, shall result in

```
convert_{I' \to I}(n)
```

For each integer datatype conforming to part 1 and made directly available, there shall be integer numerals with radix 10.

For each radix for numerals made available for a bounded integer datatype with non-special value set I, there shall be integer numerals for all non-negative values of I. For each radix for numerals made available for an unbounded integer datatype, there shall be integer numerals for all non-negative integer values smaller than 10^{20} .

For each integer datatype made directly available and that may have special values:

- a) There should be a numeral for positive infinity. There shall be a numeral for positive infinity if there is a positive infinity in the integer datatype.
- b) There should be numerals for quiet and signalling NaNs.

5.5.2 Numerals for floating point datatypes

Let D be a non-special value set for fixed point numerals for the datatype corresponding to F. Let F' be a non-special value set for floating point numerals for the datatype corresponding to F.

A fixed point numeral, denoting an abstract value x in $D \cup \{-0, +\infty, -\infty, \mathbf{qNaN}, \mathbf{sNaN}\}$, for a floating point datatype with non-special value set F, shall result in

$$convert_{D\to F}(x)$$

A floating point numeral, denoting an abstract value x in $F' \cup \{-0, +\infty, -\infty, \mathbf{qNaN}, \mathbf{sNaN}\}$, for a floating point datatype with non-special value set F, shall result in

```
convert_{F' \to F}(x)
```

For each floating point datatype conforming to part 1 and made directly available, there should be radix 10 floating point numerals, and there shall be radix 10 fixed point numerals.

For each radix for fixed point numerals made available for a floating point datatype, there shall be numerals for all bounded precision and bounded range expressible non-negative values of \mathcal{R} . At least a precision (d_D) of 20 should be available. At least a range $(dmax_D)$ of 10^{20} should be available.

For each radix for floating point numerals made available for a floating point datatype with non-special value set F, there shall be numerals for all bounded precision and bounded range expressible non-negative values of \mathcal{R} . The precision and range bounds for the numerals shall be large enough to allow all non-negative values of F to be reachable.

For each floating point datatype made directly available and that may have special values:

- a) There should be a numeral for positive infinity. There shall be a numeral for positive infinity if there is a positive infinity in the floating point datatype.
- b) There should be numerals for quiet and signalling NaNs.

The conversion operations used for numerals as operations should be the same as those used by default for converting strings to values in conforming integer or floating point datatypes.

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.

68 Notification

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 or F (as appropriate), or be special values $(-0, -\infty, +\infty, \text{ or a } \mathbf{qNaN})$.

Floating point datatypes that satisfy the requirements of IEC 60559 have special values in addition to the values in F. These are: $-\mathbf{0}$, $+\infty$, $-\infty$, signalling NaNs (sNaN), and quiet NaNs (qNaN). Such values may be 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 $-\mathbf{0}$, $+\infty$, $-\infty$, or NaN.

Continuation values of -0, $+\infty$, $-\infty$, and NaN are required only if the parameter iec_559_F has the value **true**. If the implementation can represent such special values in the result datatype, they should be used according to the specifications in this part.

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 radian arc sine operation for an argument x ($arcsin_F(x)$) might be invoked as

```
\begin{array}{lll} \arcsin(x) & \text{in Pascal [27] and Ada [11]} \\ \text{asin}(x) & \text{in C [17] and Fortran [22]} \\ (\text{asin } x) & \text{in Common Lisp [42] and ISLisp [24]} \\ \text{function asin}(x) & \text{in COBOL [19]} \end{array}
```

with a suitable expression for 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.

NOTE 3 - An arithmetic expression might not be executed as written.

For example, if x is declared to be single precision (SP) floating point, and calculation is done in single precision, then the expression

```
arcsin(x)
might translate to
```

```
arcsin_{SP}(\mathbf{x})
```

If the language in question did all computations in double precision (DP) floating point, the above expression might translate to

```
arcsin_{DP}(convert_{SP\to DP}(\mathbf{x}))
```

Alternatively, if x was declared to be an integer, and the expected result datatype is single precision float, the above expression might translate to

```
convert_{DP \to SP}(arcsin_{DP}(convert_{I \to DP}(\mathbf{x})))
```

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 " $\cos(-x)$ " " $\cos(x)$ " (exactly the same result) or " $pi \arccos(x)$ " " $\arccos(-x)$ " (more accurate result) or " $\exp(x)-1$ " " $\exp(x)$ " (more accurate result if x > -1, less accurate result if x < -1, different notification behaviour).
- 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.

NOTE 4 – It is highly desirable that programming languages intended for numerical use provide means for limiting the transformations applied to particular arithmetic expressions.

8 Documentation requirements

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

NOTE – 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 $big_angle_r_F$ and $big_angle_u_F$. Only big angle parameters that are relevant to the provided operations need be given.
- d) For the $nearest_F$ function, the rule used for rounding halfway cases, unless iec_559_F 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.
- 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 $big_angle_r_F$ and $big_angle_u_F$
- 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.

- 1) 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**.

ECHORN.COM. Click to view the full politic agost 2.200 h.

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

```
max\_error\_sinh_F \in [0.5, 2 \cdot rnd\_error_F]
     max\_error\_tanh_F \in [0.5, 2 \cdot rnd\_error_F]
    max\_error\_convert_F = 0.5

max\_error\_convert_D = 0.5
```

For $u \in G_F$, the $max_error_sinu_F(u)$ parameter shall have a value in the interval [0.5, 2]. $max_error_sin_F$], and the $max_error_tanu_F(u)$ parameter shall have a value in the interval $[0.5, 2 \cdot max_error_tan_F]$. For $u \in T$, the $max_error_sinu_F(u)$ parameter shall be equal to $max_error_sin_F$, and the $max_error_tanu_F(u)$ parameter shall be equal to $max_error_tan_F$, for the same library

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.

The extra accuracy requirements together with the sign and monotonicity requirements imply a number of requirements that are not stated explicitly, due to that they are implied. Removing one or more of thee given requirements may thus remove some weaker requirements that were not intended to be removed. Some of the remaining weaker requirements may need to be stated explicitly if a stronger requirement is removed.

A partially conforming implementation shall document which extra accuracy requirements are not fulfilled by the implementation, and which weaker requirements that are still fulfilled.

A.3 Relationships to other operations relaxation

This part has a number of requirements giving relations to other operations. These are detailed in the paragraphs beginning with wordings like "Relationship to the op_F^* approximation helper function for operations in the same library shall be:".

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

A partially conforming implementation shall document which operation relationships are not fulfilled by the implementation.

A.4 Very-close-to-axis angular normalisation relaxation

This part requires, explicitly and by implication, that angular normalisation (sometimes called argument reduction) is done so that the (intermediate or explicit) result is accurate within less than an ulp. For angular values, especially in radians, that denote an angle very close to an axis, that requires extra high precision in the calculation of the normalised value.

In a partially conforming implementation the accuracy requirements for angular normalisation for angles that are very close to an axis need not be fulfilled.

A partially conforming implementation shall document which trigonometric operations and for which (small) intervals near axes angular values, that are not so large that **absolute_precision_underflow** notifications would be the result, the angular normalisation accuracy requirements are not fulfilled by the implementation. The implementation shall also document how large the absolute error for angular normalisation is also for angles that are in those intervals very near an axis. It may be appropriate for a binding to specify one or more parameters describing this relaxation if this relaxation is allowed by a binding. The maximum error parameter values given by an implementation shall still adequately reflect the accuracy of the relevant trigonometric operations for angular values outside of those very-near-axis intervals, if a claim of partial conformity is made.

A.5 Part 1 requirements relaxation

Part 2 depends on the datatypes and operations specified in part 1. Part 1 allows for partial conformity. Part 2 operations may thus be only partially conforming if a relevant datatype or part 1 operation is only partially conforming to part 1.

ECHORM. Com. Cick to view the full poly of Isonitic hope of the control of the co

ECMORM. Com. Click to view the full policy of isonitic agost 2:200 h.

76 Partial conformity

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). This allows the standard itself to be more concise. The clause numbering matches that of the standard, although additional clauses have been added.

B.1 Scope

The scope of LIA-2 includes the traditional arithmetic operations, that are not already covered by LIA-1, usually provided in programming languages. This includes operations that are numeric approximations to real elementary functions. Even though these operations usually are implemented in software rather than hardware they are still to be regarded as *atomic* in the sense that they are never (as seen by the user) interrupted by an intermediate notification.

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 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. In order to include also fixed point string formats for floating point values, fixed point datatypes are also involved in some of the LIA-2 conversion operations.

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. Annex C also surveys which operations are already covered by various programming languages.

LIA-2 includes some operations that are not (yet) common in programming languages. Like operations to normalise angular values, and to convert angular values between different angular units. These operations are closely related to the other operations included in LIA-2, and these operations are non-trivial to implement with high accuracy. The angular normalisation operations are useful to keep high accuracy in the angular values used when increasing angular values are used.

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

B. Rationale 77

arithmetic datatypes, like rational number or fixed point datatypes. Some of these may become the topic of standardization 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 operation 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* standardization, it standardized.

LIA-2 does not cover how to represent numeric values, internally (as bit patterns) or externally (as character strings).

B.2 Conformity

Conformity to this standard is dependent on the existence of language binding standards. Each programming language committee (or other organization 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.

A conforming system consists of an implementation (which obeys the requirements) together with documentation which shows how the implementation conforms to the standard. This documentation is vital since it gives crucial characteristics of the system, such as the range for trigonometric operations, and the accuracy of the operations.

The binding of IIA-2 facilities to a particular programming language should be as natural as possible. Existing language syntax and features should be used for operations, parameters, notification, and so on. For example, if a language expresses application of cosine as " $\cos(x)$," then LIA-2 cosine operations \cos_F should be bound to (overloaded) " \cos " functions.

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.

LIA-2 has fairly strict accuracy requirements. Annex A deals with the case that an implementation (or binding standard) conforms to most aspects of LIA-2, but not necessarily all of the accuracy requirements.

Some implementations, or binding standards, may wish to conform to most of the requirements in LIA-2, but make exceptions from the specifications given by LIA-2 in certain cases. Some of the bindings examples in annex C also exemplify, in different ways, such changes of specification. Real bindings are expected to elaborate such differences much more than in the examples given in annex C.

B.2.1 Validation

LIA-2 gives a very precise description of the operations included. This will expedite the construction of conformity tests. It is important that objective tests are available.

LIA-2 does not define any process for validating conformity.

Independent assurance of conformity to LIA-2 could be by spot checks of products by a validation suite. Alternatively, checking could be regarded as the responsibility of the vendor, who would then document the evidence supporting any claim to conformity.

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

LIA-2 uses the same specification mechanisms as LIA-1. LIA-2, however, uses helper functions to a much higher degree, in particular for the specification of the operations approximating elementary transcendental functions.

As in LIA-1, operations specified in LIA-2 are done so by cases, and in some of the cases helper functions are used. In contrast to LIA-1, LIA-2 also cover cases that involve "special values" for the floating point operations. The specification of how to handle these "special values" as arguments and results for the included operations is one of the major added-values of LIA-2.

The cases in each operation specification are non-overlapping, though there is an "otherwise" case at the end of many lists of cases.

B.4.1 Symbols

B.4.1. Sets and intervals

The interval notation is in common use. It has been chosen over the other commonly used interval notation (with brackets and round parentheses mixed) because the chosen notation has no risk of confusion with the pair notation.

B.2.1 Validation 79

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 program excerpts which are clearly marked, any programming language. For example, x/u stands for the mathematically exact result of dividing x by u, independently of whether that value is representable in any floating point datatype or not, and $x/u \neq div_F(x,u)$ is often the case. Likewise, = is the mathematical equality, not the eq_F operation: $0 \neq -\mathbf{0}$, while $eq_F(0, -\mathbf{0}) = \mathbf{true}$.

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 (or, more precisely, partially specified) 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 that may be bound to LIA-2 operations. sin and cos are thus very different from sin and cos

B.4.1.4 Exceptional values

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 distinctions of this kind, if desired.

infinitary is used for integer operations, when the operation rightfully should return an infinitary value, but no infinitary value occurs among the arguments. **infinitary** is also used for floating point operations for the same circumstances. That includes when the approximated real-valued function has a pole at the argument point.

LIA allows for three methods for handing notifications: recording of indicators, change of control flow (returnable or not), and termination of program. The LIA-2 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. For example, **invalid** without a specified continuation value may cause change of control flow (like an Ada [11] exception), while **invalid** with a specified continuation value may 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.

Most bindings are expected to be such that **underflow** and **infinitary** are "quietly handled. If infinities are guaranteed to be representable, **infinitary** may even be disregarded completely, quietly returning the infinitary result without even any setting of any indicator."

B.4.1.5 Datatypes

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 sets, multi-sets (bags), and tuples also qualify.

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 that LIA-1 used the term denormal for what IEC 60559 and LIA-2 refer to as subnormal.

B.5 Specifications for the numerical functions

The abstract values used in the specifications are independent of datatype, just like the mathematical numbers are. That they are represented differently in, say, single precision and in double precision is out of scope for LIA-2.

The specifications in LIA-2 for floating point operations give details about certain special values (they are 'special' in that they are not in \mathcal{R}). These special values are commonly representable in floating point datatypes, in particular all floating point datatypes conforming to IEC 60559.

B.5.1 Basic integer operations

Integer datatypes can have infinity values as well as NaN values, and also may have a -0. A corresponding I must, however, be a subset of \mathcal{Z} . -0 is commonly available when the integer datatype is represented using radix-minus-1-complement, e.g., 1's complement. When using, e.g., 2's complement, the representation that would otherwise represent the most negative value can be used as a NaN. Especially for unbounded integer types, the inclusion of infinities is advisable, not for overflow, since these do not occur, but in order to have a smallest and a largest value in the type.

B.5.1.1 The integer result and wrap helper functions

The $result_I$ helper function notifies overflow when the result cannot be represented in T. When an overflow occurs, and recording of indicators is the method for handling (integer) overflows, a continuation value must be given. For bounded integer datatypes, $maxint_F$ and $minint_F$ can be suitable continuation values, if infinities are not representable. In some instances a wrapped result, see below, may be used as continuation value on overflow. Few integer datatypes offer representations for positive and negative infinity. In case such representations are offered, they can be used as continuation values on overflow, similar to their use in floating point datatypes. LIA does not specify the continuation value in this case, that is left to bindings or implementations, but LIA does require that the continuation value(s) be documented.

The $wrap_I$ helper function wraps the result into a value that can be represented in I. The result is wrapped in such a way that the value returned can be used to implement extended range integer arithmetic.

B.5.1.2 Integer maximum and minimum

The operations for integer maximum and minimum are trivial, except taking the maximum or minimum of an empty sequence (empty array, empty list, zero number of parameters, or similar). The case for zero number of parameters is often syntactically excluded (as in Fortran, Common Lisp, and ISLisp), while an empty array or empty list given as a single argument must usually be possible to handle at 'runtime' LIA specifies an **infinitary** notification for this case. **infinitary** is to be interpreted as "exact infinite result from finite operands", in this case an empty list of numbers. The **infinitary** notification is not specified if any of the arguments is an infinity.

If infinity values are required to be available for a particular integer datatype, a binding may require the continuation values specified to be returned without any **infinitary** notification. When the specified continuation value, $+\infty$ or $-\infty$, is not available, other suitable continuation values may be used, and if so they must be documented. If the integer datatype is bounded, but without infinities, $maxint_F$ may be used in place of $+\infty$ and $minint_F$ may be used instead of $-\infty$.

Infinities as arguments are not specified for these operations, since infinities are rarely available in integer datatypes. However, compare the specification for max and min operations for floating point datatypes (clause 5.2.2).

B.5.1.3 Integer diminish

Integer diminish is sometimes called 'monus'.

B.5.1.4 Integer power and arithmetic shift

The integer arithmetic shift operations can be used to implement integer multiplication and integer division more quickly in special cases (assuming the shift operation is supported by the hardware, and that support is used).

The shift operations shift either 'right' or 'left' depending on the sign of the second argument. 'Right' shift is done with a positive second argument, and 'left' shift with a negative second argument.

Any continuation value used on overflow here must be documented, either by the binding standard or by the implementation.

B.5.1.5 Integer square root

B.5.1.6 Divisibility tests

Even and odd are simple special cases of the divisibility test offered as separately named operations in several programming languages.

B.5.1.7 Integer division (with floor, round, or ceiling) and remainder

When the result of a division between integers is not an integer, but the final result is required to be an integer, the quotient must be rounded. There are several ways of doing this; floor, ceiling, and unbiased round to nearest being the most important. Truncating, rounding towards zero, is often provided, which, however, may introduce subtle program errors. Integer division, and remainder, is often used for grouping into groups of n items, it is natural to put item i into group divide(i, n). If i can be negative, and truncation is used, group 0 will get $2 \cdot n - 1$ items, rather than the desired n.

 pad_I returns the *negative* of the remainder after division and ceiling. The reason for this is twofold: 1) for unsigned integer datatypes the remainder is ≤ 0 , and would thus often not be representable unless negated, and 2) it is intuitively easier to think of the "places left in the last unfilled group of equi-sized and packed groups" as a positive entity, a padding.

 $residue_I$ can overflow only for unsigned integer datatypes $(minint_I = 0)$, and does so for too many cases, and negating it does not change this. $residue_I$ should therefore not be provided for unsigned integer datatypes. $residue_I$ rounds in the same way as $residue_F$. $residue_F$ is often referred to as IEEE remainder.

When there is no exception, for $n \in \mathcal{Z}$ these operations fulfill:

```
\begin{aligned} quot_I(x+n\cdot y,y) &= quot_I(x,y) + n, \\ ratio_I(x+2\cdot n\cdot y,y) &= ratio_I(x,y) + 2\cdot n, \\ group_I(x+n\cdot y,y) &= group_I(x,y) + n, \\ mod_I(x+n\cdot y,y) &= mod_I(x,y), \\ residue_I(x+2\cdot n\cdot y,y) &= residue_I(x,y), \text{ and } \\ pad_I(x+n\cdot y,y) &= pad_I(x,y). \end{aligned}
```

Note that the div_I^t and rem_I^t from LIA-1 do not fulfill similar useful equalities, due to the disruption around 0 for this pair of operations.

And, when there is no exception, the sign rules are:

```
quot_I(x,y) = -group_I(-x,y),

quot_I(x,y) = -group_I(x,-y),

ratio_I(x,y) = -ratio_I(-x,y),

ratio_I(x,y) = -ratio_I(x,-y),

mod_I(x,y) = -pad_I(x,-y), and

residue_I(x,y) = residue_I(x,-y).
```

Finally, when there is no exception, the integer division and remainder operations come in pairs that fulfill:

```
quot_I(x, y) \cdot y + mod_I(x, y) = x,

ratio_I(x, y) \cdot y + residue_I(x, y) = x, and

group_I(x, y) \cdot y - pad_I(x, y) = x.
```

B.5.1.8 Greatest common divisor and least common positive multiple

The greatest common divisor is useful in reducing a rational number to its lowest terms. The least common multiple is useful in converting two rational numbers to have the same denominator.

Returning 0 for $gcd_I(0,0)$, as is sometimes suggested, would be incorrect, since the greatest common divisor for 0 and 0 should be the supremum (upper limit) of \mathbb{Z}^+ , since elements of \mathbb{Z}^+ divide 0. The supremum of \mathbb{Z}^+ is infinity. Note also that for an n > 0, $gcd_I(n, +\infty)$ should be n, given a reasonable extension of gcd_I to cover infinity arguments.

 gcd_I will overflow only if $bounded_I = \mathbf{true}$, $minint_I = maxint_I - 1$, and both arguments are $minint_I$. The greatest common divisor is then $-minint_I$, which then is not in I.

Least common positive multiple, $lcm_I(x, y)$, overflows for many "large" arguments. E.g., if x and y are relatively prime, then the least common multiple is $|x \cdot y|$, which may be greater than $maxint_I$.

B.5.1.9 Support operations for extended integer range

These operations would typically be used to extend the range of the highest level integer datatype supported by the underlying hardware of an implementation.

The two parts of an integer product, $mul_ov_I(x, y)$ and $mul_wrap_I(x, y)$ together provide the complete integer product. Similarly for addition and subtraction. The use of $wrap_I$ guarantees that **overflow** will not occur.

B.5.2 Basic floating point operations

F must be a subset of \mathcal{R} . Floating point datatypes can have infinity values as well as NaN values, and also may have a $-\mathbf{0}$. These values are not in F. The special values are, however, commonly available in floating point datatypes today, thanks to the wide adoption of IEC 60599.

Note that for some operations the exceptional value **invalid** is produced only for argument values involving -0, $+\infty$, $-\infty$, or **sNaN**. For these operations the signature given in LIA-2 does not contain **invalid**.

The report *Floating-Point C Extensions* [57] discusses possible ways of exploiting the IEC 60559 special values, much of which is now integrated in C. The report identifies some of its suggestions

as controversial and cites Branch Cuts for Complex Elementary Functions, or Much Ado about Nothing's Sign Bit [52] as justification.

The following paragraphs is a short summary of the specifications of IEC 60559 regarding the creation and propagation of signed zeros, infinities, and NaNs. There is also some discussion of the material in references [52, 53, 50] where suggestions on this matter are made.

IEC 60559 specifies that 0 and -0 compare equal. The sign is supposed to indicate the direction of approach to zero. The sign is reliable for a zero generated by underflow in a multiplication or division operation, and should be reliable also for operations that approximate elementary transcendental functions (see the LIA-2 specifications in clause 5.3). It is not reliable for a zero generated by an implied subtraction of two floating point numbers with the same value, for which case the zero is arbitrarily given a + sign. The phrase "implied subtraction" indicates either the addition of two oppositely signed numbers or the subtraction of two like signed numbers.

On occurrence of floating overflow or division of a non-zero number by zero, an implementation conforming to IEC 60559 sets the appropriate status flag (if trapping is not enabled) and then continues execution with a result of $+\infty$ or $-\infty$ if rounding is to nearest. Infinities as such do not indicate that an overflow or division by zero has occurred; infinities can be exact values. IEC 60559 states that the arithmetic of infinities is that associated with mathematical infinities. Thus, an infinity times, plus, minus, or divided by a non-zero finite floating point number yields an infinity for the result; no status flag is set and execution continues. These rules are not necessarily valid for infinities generated by overflow, though they are valid if the infinitary arguments are exact.

NaNs are generated by invalid operations on infinities, 0 divided by 0, and the square root of a negative number (other than $-\mathbf{0}$). Thus NaNs can represent unknown real or complex values, as well as totally undefined values. IEC 60559 requires that the result of any of its basic operations with one or more NaN arguments shall be a NaN. This principle is not extended to the numerical functions by Floating-Point C Extensions [57]. The controversial specifications in Floating-Point C Extensions [57], Branch Cuts for Complex Elementary Functions, or Much Ado about Nothing's Sign Bit [52], and Lecture Notes on the Status of IEEE Standard 754 for Binary Floating-Point Arithmetic [53] are based on an assumption that all NaN operands represent finite non-zero real-valued numbers.

The LIA-2 policy (for clauser 5.2 and 5.3) for dealing with signed zeros, infinities, and NaNs is as follows:

- a) The output is a quiet NaN for any operation for which one (or more) arguments is a quiet NaN, and none of the other arguments is a signalling NaN. There is then no notification.
- b) If a mathematical function h(x) is such that h(0) = 0, the corresponding operation $op_F(x)$ returns x if $x \in \{0, -\mathbf{0}\}$ and h has a positive derivative at 0, and $op_F(x)$ returns $neg_F(x)$ if $x \in \{0, -\mathbf{0}\}$ and h has a negative derivative at 0.
- c) For an argument vector, \overrightarrow{x} , where that argument vector involves $0, -\mathbf{0}, +\infty$, or $-\infty$, the result of the operation $op_F(\overrightarrow{x})$ is

$$\lim_{\overrightarrow{z} \to \overrightarrow{x}} h(\overrightarrow{z})$$

where an approach to zero is from the positive side if $\vec{x} = (..., 0, ...)$, and the approach is from the negative side if $\vec{x} = (..., -0, ...)$. There is no notification if the limit exists, is finite, and is path independent. The returned value is $+\infty$ or $-\infty$ if the limiting value is unbounded, and the approach is towards a point infinitely far from the origin. The returned value is **infinitary** $(+\infty)$ or **infinitary** $(-\infty)$ if the limiting value is unbounded,

and the approach is towards a finite point. The result is $-\mathbf{0}$ if the limit is zero and the approaching values are path independently negative. The result is 0 if the limit is zero and the approaching values are path independently positive. If a path independent limit does not exist the value returned is **invalid**, and a notification occurs, with a continuation value of \mathbf{qNaN} if appropriate.

An exception is made for the arc_F and $arcu_F$ operations, where it is found significantly more useful to return certain non-exceptional values for the origin and for the four double infinity argument cases, than to return an exceptional value, even with non-NaN continuation values.

B.5.2.1 The rounding and floating point result helper functions

The $result_F$ helper function notifies overflow when the result is too large to be approximated by a value in F. The $result_F$ helper function notifies underflow when there is (risk for) denormalisation loss for a tiny result. The $result_F$ helper function also ensures that a properly signed zero is the continuation value when a zero is appropriate for an underflow continuation value. When an overflow or underflow occurs, and recording of indicators is the method for handling (floating point) overflow or underflow, a continuation value must be provided. LIA-2 specifies a continuation value, and if that can be represented in the target datatype, that value should be used as continuation value. If the parameter iec_559_F has the value true, then IEC 60559 in many cases require particular continuation values (consistent with what is specified by LIA-2) to be used.

The continuation values for overflow are defined to be in accordance with IEC 60559. These particular choices for continuation values are useful for interval arithmetic.

B.5.2.2 Floating point maximum and minimum

As for the integer case, the maximum and minimum of empty sequences need be handled, but for floating point datatypes, infinities are usually available.

For floating point datatypes there is also usually a negative zero available, and returning the correct sign on a zero result for the maximum and minimum operations requires more than simple comparisons to implement. The signs of zeroes may need to be inspected using *copysign* or *isnegativezero*.

B.5.2.3 Floating point diminish

As for the integer case, this operation computes the positive difference. Note that $dim_F(+\infty, +\infty) =$ **invalid**(**qNaN**) is consistent with that $sub_F(+\infty, +\infty) =$ **invalid**(**qNaN**) according to IEC 60559.

An implementation of dim_F could be if x >= y then x-y else 0.

B.5.2.4 Floor, round, and ceiling

Since $fmax_F$ always has an integral value for floating point types that conform to LIA-1, no overflow can occur for these operations.

Note that the sign of a zero result is maintained in accordance with IEC 60559:

```
floor_F(x) = neg_F(ceiling_F(neg_F(x)))

rounding_F(x) = neg_F(rounding_F(neg_F(x)))

ceiling_F(x) = neg_F(floor_F(neg_F(x)))
```

Negative zeroes, if available, are handled in such a way as to maintain these identities.

Note that $rounding_rest_F$ always is an exact operation, while $floor_rest_F$ is not always exact for negative arguments, and $ceiling_rest_F$ is not always exact for positive arguments.

B.5.2.5 Remainder after division and round to integer

The remainder after division and unbiased round to integer (IEC 60559 remainder, or IEEE remainder) is always an exact operation (unless there is an implied division by zero), even if the floating point datatype only conforms to LIA-1, but not to the more specific IEC 60559.

Remainder after floating point division and floor to integer cannot be exact for all pairs of arguments from F. For a small negative numerator and a positive denominator, the resulting value loses much absolute accuracy in relation to the original value. Such an operation is therefore not included in LIA-2. Similarly for floating point division and ceiling.

See also the radian normalisation and the argument angular-unit normalisation operations (5.3.8.1, and especially 5.3.9.1).

B.5.2.6 Square root and reciprocal square root

The inverses of squares are double valued, the two possible results having the same magnitude with opposite signs. For a non-zero result, LIA-2 requires that each of the corresponding operations return a positive result.

 \sqrt{x} cannot be exactly halfway between two values in F if $x \in F$. For \sqrt{x} to be exactly halfway between two values in F would require that it had exactly (p+1) digits (last digit non-zero) for its exact representation. The square of such a number would require at least $(2 \cdot p + 1)$ digits with last p+1 digits not all zero, which could not equal the p-digit number x.

The extensions $sqrt_F(+\infty) = +\infty$ and $sqrt_F(-0) = -0$ are mandated by IEC 60559. LIA-2 also requires that these hold for implementations which support infinities and signed zeros. However, it should be noted that while the second is harmless, the first may lead to erroneous results for a $+\infty$ generated by an addition or subtraction with result just barely outside of $[-fmax_F, fmax_F]$ after rounding. Hence its square root would be well within the representable range. The possibility that LIA-2 should require that $sqrt_F(+\infty) = invalid(+\infty)$ was considered, but rejected because of the principle of regarding arguments as exact, even if they are not exact, when there is a non-degenerate neighbourhood around the argument point, for which the mathematical function on \mathcal{R} is defined. In addition $sqrt_F(+\infty) = +\infty$ is already required by IEC 60559.

Note that the requirement that $sqrt_F(x) = invalid(qNaN)$ for x strictly less than zero is mandated by IEC 60559. It follows that NaNs generated in this way represent imaginary values, which would become complex through addition and subtraction, and even imaginary infinities on multiplication by ordinary infinities.

The rec_sqrt_F operation will increase performance for scaling a vector into a unit vector. Such an operation involves division of each component of the vector by the magnitude of the vector or, equivalently and with higher performance, multiplication by the reciprocal of the magnitude.

B.5.2.7 Multiplication to higher precision floating point datatype

This operation is intended for the case that there exist at least two floating point datatypes F and F', ideally such that the product of two numbers of type F is always exactly representable in type F'.

To obtain higher precision for multiplication, in the absence of a suitable level of precision F', a programmer can exploit the paired mul_F and mul_lo_F operations.

B.5.2.8 Support operations for extended floating point precision

These operations would typically be used to extend the precision of the highest level floating point datatype supported by the underlying hardware of an implementation. There is, however, no intent to provide a set of operations suitable for the implementation of a *complete* package for the support of calculations at an arbitrarily high level of precision.

The major motivation for including them in LIA-2 is to provide a capability for accurately evaluating residuals in an iterative algorithm. The residuals give a measure of the error in the current solution. More importantly they can be used to estimate a correction to the current solution. The accuracy of the correction depends on the accuracy of the residuals. The residuals are calculated as a difference in which the number of leading digits cancelled increases as the accuracy of the solution increases. A doubled precision calculation of the residuals is usually adequate to produce a reasonably efficient iteration.

For the basic floating point arithmetic doubled precision operations, the high parts may be calculated by the corresponding floating point operations as specified in LIA-1. Note, however, that in order to implement exact floating point addition and subtraction, rnd_F must round to nearest. If $add_F(x,y)$ rounds to nearest then the high and low parts represent x + y exactly.

When the high parts of an addition or subtraction overflows, the low parts, as specified by LIA-2, return their results as if there was no overflow. $add lo_F$ and $sub lo_F$ can underflow only when subnormals are not supported. In addition, if the high part underflows, then the low part is zero.

The product of two numbers, each with p digits of precision, is always exactly representable in at most $2 \cdot p$ digits. The high and low parts of the product will always represent the true product.

The remainder for division is more useful than a $2 \cdot p$ -digit approximation. The remainder will be exactly representable if the high part differs from the true quotient by less than one ulp. The true quotient can be constructed p digits at a time by division of the successive remainders by the divisor.

The remainder for square root is more useful than a low part for the same reason that the remainder is more useful for division. The remainder for the square root operation will be exactly representable only if the high part is correctly rounded to nearest, as is required by the specification for $sqrt_F$.

See Semantics for Exact Floating Point Operations [62] for more information on exact floating point operations.

See Proposal for Accurate Floating-Point Vector Arithmetic [63] for more information on exact, or high accuracy, floating point summation and dot product. These operations may be the subject of an amendment to LIA-2.

B.5.3 Elementary transcendental floating point operations

The basic floating point operations of LIA-2 and the elementary transcendental floating point operations have been separated into two different clauses of LIA-2, since they use slightly different specification mechanisms. The basic floating point operations need no approximation helper functions. The elementary transcendental floating point operations need approximation helper function in order to express the wider error tolerance for these operations.

B.5.3.1 Maximum error requirements

The $max_error_op_F$ parameters measure the discrepancy between the computed value $op_F(x)$ and the true mathematical value f(x) in ulps of the true value. The magnitude of the error bound is thus available to a program from the computed value $op_F(x)$. Note that for results at an exponent boundary for F, y, the error away from zero is in terms of $ulp_F(y)$, whereas the error toward zero is in terms of $ulp_F(y)/r_F$, which is the ulp of values slightly smaller in magnitude than y.

Within limits, accuracy and performance may be varied to best meet customer needs. Note also that LIA-2 does not prevent a vendor from offering two or more implementations of the various operations.

The operation specifications define the domain and range for the operations. This is done partly by the given signature, and partly by the specification of cases that do not return **invalid**. In addition, the computational domain and range are more limited for the operations than for the corresponding mathematical functions because the arithmetic datatypes are subsets of \mathcal{R} . Further, any (conforming) F is limited in range, and the operations may return an overflow or an underflow. Thus the actual domain of $exp_F(x)$ is approximately given by $x \leq \ln(fmax_F)$. For larger values of x, $exp_F(x)$ will overflow, though for $x = +\infty$ the exact result $+\infty$ will be returned. The actual range extends over positive F, although there are non-negative values, $v \in F$, for which there is no $x \in F$ satisfying $exp_F(x) = v$.

The thresholds for the **overflow** and **underflow** notifications are determined by the parameters defining the arithmetic datatypes. The threshold for an **invalid** notification is determined by the domain of arguments for which the mathematical function being approximated is defined. The **infinitary** notification is the **operation**'s counterpart of a mathematical pole of the mathematical function being approximated by the operation. The threshold for **absolute_precision_underflow** is determined by the parameters $big_angle_r_F$ and $big_angle_u_F$.

LIA-2 imposes a fairly tight bound on the maximum error allowed in the implementation of each operation. The tightest possible bound is given by requiring rounding to nearest, for which the accompanying performance penalty is often unacceptably high for the operations approximating elementary transcendental functions. LIA-2 does not require round to nearest for such operations, but allows for a slightly wider error bound characterised via the $max_error_op_F$ parameters. The $max_error_op_F$ parameters must be documented by the implementation for each such parameter required by LIA-2. A comparison of the values of these parameters with the values of the specified maximum value for each such parameter will give some indication of the "quality" of the routines provided. Further, a comparison of the values of this parameter for two versions of a frequently used operation will give some indication of the accuracy sacrifice made in order to gain performance.

Language bindings are free to modify the error limits provided in the specifications for the operations to meet the expected requirements of their users.

Material on the implementation of high accuracy operations is provided in for example [50, 52, 59].

B.5.3.2 Sign requirements

The requirements imply that the sign of the result or continuation value is to be reliable, except for the sign of an infinite result or continuation value, where except for a signed zero argument, it is often the case that one cannot determine the sign of the infinity. Still for sign symmetric mathematical functions, the approximating operation is also sign symmetric, including infinitary results.

Note that the sign requirements stated generally imply some requirements that are not given explicitly for each operation specification in LIA-2. For example, $sin_F^*(n \cdot 2 \cdot \pi + \pi) = 0$ is a requirement implied by the general sign requirements.

B.5.3.3 Monotonicity requirements

A maximum error of 0.5 ulp implies that an approximation helper function must be a monotonic approximation to the mathematical function. When the maximum error is greater than 0.5 ulp, and the rounding is not directed, this is not automatically the case.

There is no general requirement that the approximation helper functions are *strictly* monotone on the same intervals on which the corresponding exact function is strictly monotone, however, since such a requirement cannot be made due to the fact that all floating point types are discrete, not continuous.

The monotonicity requirements are not extended to the angular unit arguments (for the operations that take such an argument or arguments). The reason for this is that it is thought both hard to implement, and also of no interest to users to have monotonicity on that (those) argument(s), since the angular unit is not normally varied, except when converting between angular units, and even then the unit arguments involved are usually constants.

The monotonicity requirements together with the extra accuracy requirements also imply requirements not explicitly stated. For example $-1 \leq sin_F^*(x) \leq 1$ is such an implied requirement. Therefore, even if some of the extra accuracy requirements are relaxed (see annex A), it may be necessary to reintroduce some of the requirements that were implied.

B.5.3.4 The result helper function

The $result_F^*$ helper function is more suitable than the $result_F$ helper function when the approximation is not guaranteed to be 0.5 ulp nor guaranteed to be directed.

Ideally, however, though not expressed in the LIA-2 specifications, also the operations approximating elementary transcendental functions obey the rounding mode (if the implementation has rounding modes) in the sense that "round towards negative infinity" gives a result that is less than or equal to the true result (and similarly for "round towards positive infinity"). However, and in contrast to the basic arithmetic operations, the error may then be more than 1 ulp. So even if the rounding modes are heeded also for the operations approximating elementary transcendental functions, the interpretation of the rounding modes are still looser than for the basic arithmetic operations. LIA-2 as such does not require even this looser interpretation. It is up to implementations, and the accompanying documentation, to implement this, if desired, and to

document the behaviour of these operations under different rounding modes. To get reliable upper and lower bounds of the true result, that are also close to the true result, is useful for interval arithmetic. Still, using the returned result, within the error bounds specified by LIA-2, together with the relevant $max_error_op_F$ parameter, one can get a (perhaps slightly wider) safe interval containing the true result.

B.5.3.5 Hypotenuse

The $hypot_F$ operation can produce an overflow only if both arguments have magnitudes very close to the overflow threshold. $hypot_F$ only underflows if both arguments are subnormal numbers. Care must be taken in its implementation to either avoid or properly handle overflows and underflows which might occur in squaring the arguments. The function approximated by this operation is mathematically equivalent to complex absolute value, which is needed in the calculation of the argument (also called phase; see arc_F) and modulus (also called absolute value; $hypot_F$) of a complex number. It is important for this application that an implementation satisfy the constraint on the magnitude of the result returned.

LIA-2's $hypot_F$ does not follow the recommendations in $Branch\ Cuts$ for $Complex\ Elementary$ Functions, or $Much\ Ado\ about\ Nothing$'s $Sign\ Bit\ [52]$ and in $Lecture\ Notes\ on\ the\ Status\ of\ IEEE$ $Standard\ 754$ for $Binary\ Floating-Point\ Arithmetic\ [53]$ which recommend that

```
hhypot_F(+\infty, \mathbf{qNaN}) = +\infty

hhypot_F(-\infty, \mathbf{qNaN}) = +\infty

hhypot_F(\mathbf{qNaN}, +\infty) = +\infty

hhypot_F(\mathbf{qNaN}, -\infty) = +\infty
```

which are based on the claim that a **qNaN** represents an (unknown) real valued number. Such a claim is not always valid, though it may sometimes be.

B.5.3.6 Operations for exponentiations and logarithms

For all of the exponentiation operations, overflow occurs for sufficiently large values of the argument(s).

There is a problem for $power_F(x, y)$ if both x and y are zero:

- Ada raises an 'exception' for the operation that is close in semantics to $power_F$ when both arguments are zero, in accordance with the fact that 0^0 is mathematically undefined.
- The X/OPEN Portability Guide, as well as C99, specifies for pow(0.0, 0.0) a return value of 1, and no notification. Those specifications agree with the recommendations in [50, 52, 53, 56]

The specification in LIA-2 follows Ada, and returns **invalid** for $power_F(0,0)$, because of the risks inherent in returning a result which might be inappropriate for the application at hand. Note however, that $power_{F,I}(0,0)$ is 1, without any notification. The reason is that the limiting value for the corresponding mathematical function, when following either of the only two continuous paths, is 1. This also agrees with the Ada specification for a floating point value raised to a power in an integer datatype, as well as that for other programming languages which distinguish these operations. The C99 (and X/OPEN) specification for the pow can be regarded as a combination of the $power_{F,I}$ and $power_F$ operations. Due to this combination, LIA-2 has a requirement that

 $power_F$ for integral second arguments and $power_{F,I}$ (in the same library) are related by equality for positive first arguments.

Along any path defined by $y = k/\ln(x)$ the mathematical function x^y has the value e^k . It follows that some of the limiting values for x^y depend on the choice of k, and hence are undefined, as indicated in the specification.

The result of the $power_F$ operation is **invalid** for negative values of the base x. The reason is that the floating point exponent y might imply an implicit extraction of an even root of x, which would have a complex value for negative x. This constraint is explicit in Ada, and is widely imposed in existing numerical packages provided by vendors, as well as several other programming languages.

The arguments of $power_F$ are floating point numbers. No special treatment is provided for integer floating point values, which may be approximate. The cases for integer values of the arguments are covered by the operations $power_{F,I}$ and $power_I$. In the example binding for C a specification for pow_F is supplied. pow_F combines $power_F$ and $power_{FZ}$ in a way suitable for C's pow operation.

For implementations of the $power_F$ operation there is an accuracy problem with an algorithm based on the following, mathematically valid, identity:

$$x^y = r_F^{y \cdot \log_{r_F}(x)}$$

The integer part (floor plus one, *not* truncation) of the product $y \log_{r_F}(x)$ defines the exponent of the result and the remaining fractional part defines the reduced argument. If the exponent is large, and one calculates p_F digits of this intermediate result, there will be fewer than p_F digits for the fraction. Thus, in order to obtain a reduced argument accurately rounded to p_F digits, it may be necessary to calculate an approximation to $y \cdot \log_{r_F}(x)$ to a few more than $\log_{r_F}(emax_F) + p_F$ base r_F digits.

In Ada95 the operation closest to *power_F*, r is specified to be computed by successive multiplications, for which the error in the evaluation increases linearly with the size of the exponent. In a strict Ada implementation there is no way that a prescribed error limit of a few ulps can be met for large exponents.

The special exponentiation operations, corresponding to 2^x and 10^x , have specifications which are minor variations on those for $exp_F(x)$. Accuracy and performance can be increased if they are specially coded, rather than evaluated as, e.g., $exp_F(mul_F(x, ln_F(2)))$ or $power_F(2, x)$. Similar comments hold for the base 2 and base 10 logarithms operations.

The $expm1_F$ operation has two advantages: Firstly, $expm1_F(x)$ is much more accurate than $sub_F(exp_F(x), 1)$ when the exponent argument is close to zero. Secondly, the $expm1_F$ operation does not **underflow** for "very" negative exponent arguments, something which may be advantageous if **underflow** handling is slow, and high accuracy for "very" negative arguments is not needed. Note in addition that underflow is avoided for this operation. This can be done only since LIA-2 adds requirements beyond those of LIA-1 regarding minimum precision (see clause 4). If those extra requirements were not done, underflow would not be justifiably removable for this operation. Similar argumentation applies to $ln1p_F$.

Similarly, there are two advantages with the $power1pm1_F$ operation: Firstly, $power1pm1_F(b, x)$ is much more accurate than $sub_F(power_F(add_F(1, b), x), 1)$ when the exponent argument is close to zero. Secondly, the $power1pm1_F$ operation does not **underflow** for "very" negative exponent arguments (when the base is greater than 1), something which may be advantageous if **underflow** handling is slow, and high accuracy for "very" negative arguments is not needed.

The handling of infinites and negative zero as arguments to the exponentiation and logarithm operations, like for all other LIA operations, follow the principles for dealing with these values as explained in section B.5.2. Note in particular that $logbase_F(b, x)$ is specified so as to be consistent with $div_F(ln_F(x), ln_F(b))$ except that $logbase_F(b, x)$ is required to be more accurate.

The $expm1_F$ and $ln1p_F$ operations are required to return the argument when the argument is in a certain interval around 0. Some floating point parameters from LIA-1 had to be made a bit stricter for LIA-2 to guarantee that this interval always is wider than the interval of subnormal numbers (this change is to be integrated with LIA-1 when LIA-1 is revised). This way underflow can always be avoided for these operations, and in the interval specified, they can with high accuracy return the argument unchanged.

Several of the operations have requirements that push the result towards a finite limiting value, so that that the limiting value is actually reached (within a reasonable margin) after rounding, even if the limiting value cannot, or otherwise need not, be reached before rounding. Similar requirements appear also in the other subclauses of clause 5.3.

Note also that even the use of the nearest approximation to e that is representable in F as a base argument to the $power_F$ and $logbase_F$ operations do not produce a duplication of exp_F and ln_F .

B.5.3.7 Introduction to operations for trigonometric elementary functions

The real trigonometric functions sin, cos, tan, cot, sec, and csc are all periodic. The smallest period for sin, cos, sec, and csc is $2 \cdot \pi$ radians (360 degrees). The smallest period for tan and cot is π radians (180 degrees) (and thus also have a period of $2 \cdot \pi$ radians (360 degrees)). The mathematical trigonometric functions are perfectly periodic. Their numerical counterparts are not that perfect, for two reasons.

Firstly, the radian normalisation cannot be exact, even though it can be made very good given very many digits for the approximation(s) of π used in the angle normalisation, returning an offset from the nearest axis, and including guard digits. The unit argument normalisation, however, can be made exact regardless of the (non-zero and, in case $denorm_F = false$ not too small) unit and the original angle, returning only a plain angle in F. LIA-2 requires unit argument angle normalisation to be exact.

Secondly, the length of one revolution is of course constant, but the density of floating point values gets sparser (in absolute spacing rather than relative) the larger the magnitude of the values are. This means that the number of floating point values gets sparser per revolution the larger the magnitude of the angle value is. For this reason the notification **absolute_precision_underflow** is introduced, together with two parameters, one for radians and one for other angular units. This notification is given when the magnitude of the angle value is "too big". Exactly when the representable angle values get too sparse depends upon the application at hand, but LIA-2 gives a default value for the parameters that define the cut-off.

Note that the absolute_precision_underflow notification is unrelated to any argument reduction problems. Argument reduction is (implicitly for radians, explicitly for other angular units) required by LIA-2 to be very accurate. But no matter how accurate the argument reduction is, floating point values are still sparser in absolute terms the larger the values are. The trigonometric operations return a result within about an ulp, and that high accuracy is wasted if the angular argument is not kept at a high accuracy too, both relative and absolute.

LIA-2 includes specifications for high accuracy angle normalisation operations, both for radians and for other angular units. The angle normalisation operations give a result within minus half a cycle to plus half a cycle (as does the angle conversion operations), unless the argument angular value is too big (or there is some other error). These operations should be used to keep the representation of angles at a high accuracy. LIA-2 also includes angle normalisation operations that can be used to maintain an even higher degree of accuracy, giving the offset from the nearest axis (though without any extra guard digits). To use these, one need to keep track of the currently nearest axis, and make appropriate adjustments in the calculations, which unfortunately complicates programs that use these nearest-axis normalisations.

Note that $rad(x) = \arccos(\cos(x))$ if $\sin(x) > 0$ and $rad(x) = -\arccos(\cos(x))$ if $\sin(x) < 0$. The first part of $axis_rad(x)$ indicates which axis is nearest to the angle x. The second part of $axis_rad(x)$ is an angle offset from the axis that is nearest to the angle x. The second part of $axis_rad(x)$ is equal to rad(x) if $\cos(x) \ge 1/\sqrt{2}$ (i.e. if the first part of $axis_rad(x)$ is (1.0)). More generally, the second part of $axis_rad(x)$ is equal to $rad(4 \cdot x)/4$.

rad(x) returns the same angle as the angle value x, but the returned angle value is between $-\pi$ and π . The rad function is defined to be used as the basis for the angle normalisation operations. The $axis_rad$ function is defined to be used as the basis for a numerically more accurate radian angle normalisation operation. The arc function is defined to be used as the basis for the arcus (angle) operations, which are used for conversion from Cartesian to polar co-ordinates.

B.5.3.8 Operations for radian trigonometric elementary functions

The radian trigonometric approximation helper functions (including those for normalisation and conversion from radians) are required to have the same zero points as the approximated mathematical function only if the absolute value of the argument is less than or equal to $big_angle_r_F$. Likewise, the radian trigonometric approximation helper functions are required to have the same sign as the approximated mathematical function only if the absolute value of the argument is less than or equal to $big_angle_r_F$. Indeed, the radian trigonometric approximation helper functions need not be defined at all outside of $[\neg big_angle_r_F, big_angle_r_F]$.

The $big_angle_r_F$ parameter may be adjusted by bindings, or even by some compiler flag, or mode setting within a program. However, this method should only allow the value of this parameter to be set to a value greater than $2 \cdot \pi$, so that at least arguments within the first two (plus and minus) cycles are allowed, and such that $ulp_F(big_angle_r_F) < \pi/1000$, so that at least 2000 evenly distributed points within the 'last' cycle (farthest away from 0) are distinguishable. The latter gives a rather low accuracy at the far ends of the range, especially if p_F is comparatively large, so values this large for $big_angle_r_F$ are not recommendable unless the application is such that high accuracy trigonometric operations are not needed. Note that if $big_angle_r_F$ is allowed to be increased, then, for conformity with LIA-2, the radian angle reduction may need to be made more precise.

For reduction of an argument given in radians, implementations use one or several approximate value(s) of π (or of a multiple of π), valid to, say, n digits. The division implied in the argument reduction cannot be valid to more than n digits, which implies a maximum absolute angle value for which the reduction yields an accurate reduced angle value.

Regarding argument reduction for radians, there is a particular problem when the result of the trigonometric operation is very small (or very big), but the angular argument is not very small. In such cases the argument reduction must be very accurate, using an extra-precise approximation

to π , relative to what is normally used for arguments of similar magnitude, so that significant digits in the result are not lost. Such loss would imply non-conformance to LIA-2 by the error in the final result being greater than that specified by LIA-2. In general, extra care has to be taken when the second part of $axis_rad(x)$ is close to 0.

Note that

- tan and sec have poles at odd multiples of $\pi/2$ radians (90 degrees).
- cot and csc have poles at multiples of π radians (180 degrees).

All four of the corresponding operations with poles may produce **overflow** for arguments sufficiently close to the poles of the functions. The tan_F operation produces no **infinitary** notification. The reason is that the poles of tan(x) are at odd multiples of $\pi/2$, which are not representable in F. The mathematical cotangent function has poles at even multiples of $\pi/2$, of which only the origin is representable in F. For a system which supports signed zeros and infinities, the continuation values are $+\infty$ and $-\infty$ for arguments of 0 and -0 respectively to $cot_F(x)$. Although the mathematical function sec has poles at odd multiples of $\pi/2$, the sec_F operation will not generate any **infinitary** notification because such arguments are not representable in F.

The **infinitary** notification cannot occur for any non-zero argument in radians because π is not representable in F, nor is $\pi/2$. For the angular unit argument trigonometric operations the sign of the infinitary continuation value has been chosen arbitrarily for an **infinitary** which occurs for a non-zero argument. However, sign symmetry, when appropriate, is maintained.

The operations may produce **underflow** for arguments sufficiently close to the zeros of the function. For a subnormal argument x, the sin_F , tan_F , $arcsin_F$, and $arctan_F$ return x for the result, with very high accuracy. Similarly, for a subnormal argument, cos_F and sec_F can return a result of 1.0 with very high accuracy.

The trigonometric inverses are multiple valued. They are rendered single valued by defining a principal value range. This range is closely related to a branch cut in the complex plane for the corresponding complex function. Among the floating point numerical functions this branch cut is "visible" only for the arc_F operation. The arc function has a branch cut along the negative real axis. For x < 0 the function has a discontinuity from $-\pi$ to $+\pi$ as y passes through zero from negative to positive values. Thus for x < 0, systems supporting signed zeros can handle the discontinuity as follows:

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

 $arc_F(x, 0) = down_F(\pi)$

There is a problem for zero argument values for this operation. The values given for the operation $arc_F(x,y)$ for the four combinations of signed zeros for x and y are those given in [52]. The following table of values is given in [52] for the value of $arc_F(x,y)$ with both of the arguments zero:

Zero arguments			
\boldsymbol{x}	y	$arc_F(x,y)$	
0	0	0	
-0	0	π	
-0	-0	$-\pi$	
0	- 0	- 0	

Note that the mathematical arc function is indeterminate (undefined) for (0,0), but close representable approximations the above result are numerically more useful than giving an **invalid** notification for such arguments. LIA-2 therefore specifies results as above.

There is also a problem for argument values of $+\infty$ or $-\infty$ for this operation. The following table of values is given in [52] for the value of $arc_F(x,y)$ with at least one of the arguments infinite:

Infinite arguments			
x	y	$arc_F(x,y)$	
+∞	$\geqslant 0$	0	
$+\infty$	+∞	$\pi/4$	
finite	+∞	$\pi/2$	
$-\infty$	+∞	$3 \cdot \pi/4$	
$-\infty$	$\geqslant 0$	π	
$-\infty$	- 0	$-\pi$	
$-\infty$	< 0	$-\pi$	
$-\infty$	$-\infty$	$-3 \cdot \pi/4$	
finite	$-\infty$	$-\pi/2$	
$+\infty$	$-\infty$	$-\pi/4$	
$+\infty$	< 0	-0	
$+\infty$	-0	-0	

If one of x and y is infinite and the other is finite, the result tabulated is consistent with that obtained by a conventional limiting process. However, the results of $\pi/4$, $-\pi/4$, $3 \cdot \pi/4$, and $-3 \cdot \pi/4$ corresponding to infinite values for both x and y, are of questionable validity, since only the quadrant is known, not the angle within the quadrant. However, these results are numerically more useful than giving an **invalid** notification for such arguments. LIA-2 therefore specifies results as above.

B.5.3.9 Operations for trigonometrics with given angular unit

At present only Ada specifies trigonometric operations with angular unit argument. LIA-2 has adopted angular unit argument operations in order to encourage uniformity among languages which might include such operations in the future. The angular units in T appear to be particularly important and have therefore been given a tighter error bound requirement. An implementation can of course have the same (tighter) error bound for all angular units. Some programming languages provide trigonometric operations with an implicit angular unit argument with value 360.

The trigonometric approximation helper functions with angular unit argument (including those for normalisation and conversion from radians) are required to have the same zero points as the approximated mathematical function. Likewise, the trigonometric approximation helper functions with angular unit argument are required to have the same sign as the approximated mathematical function. However, the trigonometric approximation helper functions with angular unit argument need not be defined at all outside of $[-big_angle_u_F \cdot |u|, big_angle_u_F \cdot |u|]$, where u is the value of the angular unit argument.

The $big_angle_u_F$ parameter may be adjusted by bindings, or even by some compiler flag, or mode setting within a program. However, this method should only allow the value of this parameter to be set to a value greater than or equal to 1, so that at least arguments within

the first two (plus and minus) cycles are allowed, and such that $ulp_F(big_angle_u_F) \leq 1/2000$, so that at least 2000 evenly distributed points within the 'last' cycle (farthest away from 0) are distinguishable. The latter gives a rather low accuracy at the far ends of the range, especially if p_F is comparatively large, so values this large for $big_angle_u_F$ are not recommendable unless the application is such that high accuracy trigonometric operations are not needed.

The $min_angular_unit_F$ parameter is specified for two reasons. Firstly, if the type F has no subnormal values ($denorm_F = \mathbf{false}$), some angle values in F are not representable after normalisation if the angular unit has too small magnitude. This gives the firm limit given in cluase 5.3.9. Secondly, even if F has subnormal values ($denorm_F = \mathbf{true}$), angular units with very small magnitude do not allow the representable angles to be particularly dense, not even if the angular value is within the first cycle. This does in itself not give rise to a particular limit value, but the limit value defined in cluase 5.3.9 is reasonable.

Provided that $|u| \ge min_angular_unit_F$, an angular unit u can be either positive or negative. If it's negative, growing angular values turns the angle "clockwise" rather than counter-clockwise as for radians and other positive angular units. Ada does not permit negative angular units, but since there is no mathematical nor numerical reason to not allow them, LIA-2 allows negative angular unit argument values, avoiding an unjustifiable and arbitrary decision to disallow them. This only very marginally complicates the specifications given in LIA-2 as well as the implementations that follow those specifications.

Note that the angular unit argument need not be integral, even though several common non-radian angular units are integral, 360, 400, etc. Note also that even the use of the nearest approximation to $2 \cdot \pi$ that is representable in F as angular unit argument does not produce a duplication of the radian trigonometric operations. The radian trigonometric operations need to use one or more approximations to π (or an integer fraction of π) that are more accurate than can be represented in F, in order to fulfill the accuracy requirements of LIA-2.

All of the argument angular unit trigonometric, and argument angular unit inverse trigonometric, approximation helper functions, including those for normalisation, angular unit conversion, and arc, are exempted from the monotonicity requirement for the angular unit argument.

If the angular unit argument, u, \mathbf{is} such that $u/4 \in F$, the $tanu_F$ operation has poles at odd multiples of u/4. This is the case for degrees (u = 360). As for $tanu_F$, if the angular unit argument, u, is such that $u/4 \in F$ the $secu_F$ operation has poles (**infinitary**) at odd multiples of u/4.

The same comments hold for the $arcu_F$ operation as for arc_F operation, except that the discontinuity in the mathematical function is from -u/2 to +u/2.

B.5.3.10 Operations for angular-unit conversions

Angular conversion operations are commonly found on 'scientific' calculators and also in Java, though then often only between degrees and radians.

Conversion of an angular value x from angular unit a to angular unit b appears simple: compute $x \cdot b/a$. Basing a numerical conversion of angular values directly on the above mathematical equality (e.g. $div_F(mul_F(x,b),a)$) loses much absolute angular accuracy, however, especially for large angular values. Instead computing $arcu_F(b,cosu_F(a,x),sinu_F(a,x))$ then gives a more accurate result. This might still not be within the accuracy required by LIA-2 for the angular unit conversion operations specified by LIA-2, which here requires a maximum error of 0.5 ulp.

Note that all of the angular conversion operations return an angularly normalised result. This is in order to maintain high absolute accuracy of the angle being represented.

B.5.3.11 Operations for hyperbolic elementary functions

The hyperbolic sine operation, $sinh_F(x)$, will overflow if |x| is in the immediate neighbourhood of $\ln(2 \cdot fmax_F)$, or greater.

The hyperbolic cosine operation, $cosh_F(x)$, will overflow if |x| is in the immediate neighbourhood of $ln(2 \cdot fmax_F)$, or greater.

The hyperbolic cotangent operation, $coth_F(x)$, has a pole at x = 0.

The inverse of cosh is double valued, the two possible results having the same magnitude with opposite signs. The value returned by $arccosh_F$ is always greater than or equal to 1

The inverse hyperbolic tangent operation $arctanh_F(x)$ has poles at x = +1 and at x = -1.

The inverse hyperbolic cotangent operation $arccoth_F(x)$ has poles at x = -1.

Like for the exponentiation and logarithm operations, there are extra accuracy requirements, for certain arguments.

When appropriate, there are also sign symmetry requirements on the approximation helper functions. These sign symmetry requirements for "ordinary" arguments are followed through in the operation specification to cover also signed zeroes and infinites. Similar requirements appear also in the other subclauses of clause 5.3.

For $sinh_F$, $tanh_F$, $arcsinh_F$, and $arctanh_F$, for a specified interval around 0, the operation returns its argument unchanged, and does so with high accuracy. Underflow notifications are also avoided for these cases, since there is no denormalisation loss.

B.5.4 Operations for conversion between numeric datatypes

Clause 5.2 of LIA-1 covers conversions from an integer type to another integer type and to a floating point type, as well as between (LIA-1 conforming) floating point types of the same radix.

LIA-2 extends these conversions to cover conversions to and from non-LIA conforming datatypes, such as conversion to and from strings, and also extends the floating point conversion specifications to handle conversions where the radices may be different.

In ordinary string formats for numerals, the string "Hello world!" is an example of a signalling NaN.

LIA-2 does not specify any string formats, not even for the special values $-\mathbf{0}$, $+\infty$, $-\infty$, and quiet NaN, but possibilities for the special values include the strings used in the text of LIA-2, as well as strings like "+infinity" or "positiva oändligheten", etc, and the strings used may depend on preference settings, as they may also for non-special values. For instance, one may use different notation for the decimal separator character (like period, comma, Arabic comma, ...), use superscript digits for exponents in scientific notation, or use Arabic digits or traditional Thai digits. String formats for numerical values, and if and how they may depend on preference settings, is also an issue for bindings or programming language specifications, not for this part of LIA.

If the value converted is greater than those representable in the target, or less than those representable in the target, even after rounding, then an overflow will result. E.g., if the target is

a character string of at most 3 digits, and the target radix is 10, then an integer source value of 1000 will result in an overflow. As for other operations, if the notification handling is by recording of indicators, a suitable continuation value must be used.

Most language standards contain (partial) format specifications for conversion to and from strings, usually for a decimal representation.

LIA-2 requires, like C99, all floating point conversion operations to be such that the error is at most 0.5 ulp. This is in contrast to IEC 60559, which allows conversion operations to have an error of up to 0.97 ulp.

B.5.5 Numerals as operations in a programming language

B.5.5.1 Numerals for integer datatypes

Negative values (except $minint_I$ if $minint_I = -maxint_I - 1$) can be obtained by using the negation operation (neg_I) .

Integer numerals in radix 10 are normally available in programming languages. Other radices may also be available for integer numerals, and the radix used may be part of determining the target integer datatype. E.g., radix 10 may be for signed integer datatypes, and radix 8 or 16 may be for unsigned integer datatypes.

Syntaxes for numerals for different integer datatypes need not be different, nor need they be the same. This part does not further specify the format for integer numerals. That is an issue for bindings.

Overflow for integer numerals can be detected at "compile time", and warned about. Likewise can notifications about invalid, e.g. for infinitary or NaN numerals that cannot be converted to the target type, be detected at "compile time" and be warned about.

B.5.5.2 Numerals for floating point datatypes

If the numerals used as operations in a program, and numerals read from other sources use the same radix, then "internal" numerals and "external" numerals (strings) denoting the same value in \mathcal{R} and converted to the same target datatype should be converted to the same value. Indeed, the requirement on such conversions to round to nearest implies this. But even if this requirements is relaxed by a binding (see Annex A), external and internal conversions should not differ.

Negative values (including negative 0, -0, if avaliable) can be obtained by using the negation operation (negative 0, -0).

Radices other than 10 may also be available for floating point numerals.

Integer numerals may also be floating point numerals, i.e. their syntaxes need not be different. Nor need syntaxes for numerals for different floating point datatypes be different, nor need they be the same. This part does not specify the syntax for numerals. That is an issue for bindings or programming language specifications.

Overflow or underflow for floating point numerals can be detected at "compile time", and warned about. Likewise can notifications about **infinitary** or **invalid**, e.g. for infinitary or NaN numerals that cannot be converted to the target type, be detected at "compile time" and be warned about.

B.6 Notification

An intermediate overflow on computing approximations to, for example, x^2 or y^2 during the calculation of $hypot_F(x,y) \approx \sqrt{x^2 + y^2}$ does not result in an overflow notification, unless the end result overflows. This is clear from the specification of the $hypot_F$ operation in this part. It is not helpful for the user of an operation to let intermediary overflows or underflows that are not reflected in the end result be propagated. Implementations of LIA-2 operations are required to shield the user from such intermediary overflows for all of the LIA-2 operations. More generally, well-made numerical software should similarly shield users of that software from overflows and underflows that are not reflected in a properly calculated end result. However, such requirements in general are beyond the scope of LIA-2.

If a single argument operation op_F , for the corresponding mathematical function f, is such that f(x) very closely approximates x, when $|x| \leq fminN_F$, then $op_F(x)$ returns x for $|x| \leq fminN_F$, and does not give a notification if there cannot be any denormalisation loss relative to f(x). For details, see the individual operation specifications for $expm1_F$, $ln1p_F$, sin_F , $arcsin_F$, tan_F , $arcsin_F$, tan_F , and $arctan_F$.

Operations specified in LIA-2 return **invalid**(**qNaN**) when passed a signalling **NaN** (**sNaN**) as an argument. Most operations specified in LIA-2 return **qNaN**, without any notification when passed a quiet **NaN** (**qNaN**) as an argument.

The different kinds of notifications occur under the following circumstances:

- a) **invalid**: when an argument is not valid for the operation, and no value in F^* or any special value result makes mathematical sense.
- b) **infinitary**: when the input operand corresponds to a pole of the mathematical function approximated by the operation, or, more generally, when the true result is infinitary, but none of the arguments is infinitary.
- c) **overflow**: when the (rounded) result is outside of the range of the result datatype.
- d) **underflow**: when a sufficiently closely approximating result of the operation has a magnitude that is so small that it might not be sufficiently accurately represented in the result datatype.
- e) absolute_precision_underflow: when the magnitude of the angle argument to a floating point trigonometric operation exceeds the maximum value of the argument for which the density of floating point values is deemed sufficient for the operation to make sense. See clause 5.3.7 and the associated discussion in this rationale (section B.5.3.7).

In order to avoid **absolute_precision_underflow** notifications, and to maintain a high accuracy implementors are encouraged to provide, and programmers encouraged to use, the angle normalisation operations specified in 5.3.8.1 and 5.3.9.1.

The difference between the **infinitary** and **overflow** notifications for floating point operations is that the first corresponds to a true mathematical singularity, and the second corresponds to a well-defined mathematical result that happens to lie outside the range of F.

B.6.1 Continuation values

For handling of notifications, the method that does recording of indicators (LIA-1, clause 6.1.2) is preferred.

100 Rationale

An implementation which supports recording of indicators must supply continuation values to be used when execution is continued following the occurrence of a notification. For systems which support signed zeros, infinities and NaNs, LIA-2 specifies how these values, as well as ordinary values, are used as continuation values. Other implementations which use recording of indicators must supply other suitable continuation values and document the values selected.

B.7 Relationship with language standards

The datatypes involved in implicit conversions need not be accessible to the programmer. For example, trigonometric operations may be evaluated in extended double precision, even though that datatype is not made available to programmers using a particular programming language. These extra datatypes should be made available, however, and the implicit conversions should be expressible as explicit conversions. At least in order to be able to show exactly which expression is going to be evaluated without having to look at the machine code.

B.8 Documentation requirements

To make good use of an implementation that conforms to LIA-2, programmers need to know not only that the implementation conforms, but *how* it conforms. LIA-2 requires implementations to document the binding between the LIA-2 operations and parameters and the total arithmetic environment provided by the implementation.

It is expected that an implementation will meet part of its documentation requirements by incorporation of the relevant language standard. However, there will be aspects of the implementation that the language standard does not specify in the required detail, and the implementation needs to document those details. For example, the language standard may specify the range of allowed parameter values, but the implementation must document the actual value. The combination of the language standard and the implementation documentation together should meet all the requirements in clause 8.

ECHORN.COM. Click to view the full politic agost 2.200 h.

Rationale

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-2. It contains suggestions for binding the "abstract" operations specified in LIA-2 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-2 operation name, or parameter name, bound to an identifier (or expression) by the binding.

Portability of programs can be improved if two conforming LIA-2 systems using the same programming language agree in the manner with which they adhere to LIA-2. For instance, LIA-2 requires that the parameter $big_angle_r_F$ 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. Name consistency cannot, however, be fully maintained between different programming languages, due to already existing differences in naming conventions, and LIA does not require wholesale naming changes, nor expression syntax changes.

The following clauses are suggestions rather than requirements because the areas covered are the responsibility of the various programming 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,
BASIC,
C,
C++,
Fortran,
Haskell,
Java,
Common bisp,
ISLisp,
Modula-2,
Pascal and Extended Pascal,
PL/I, and
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 \star .)

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 examplified 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 examplified 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 [11].

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 LIA-2 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. Ada allows in general that formal parameter names are used in calls, though one normally does not write out the formal parameter names in calls. However, in some cases a formal parameter name is needed in the call to make the overloaded name resolve to the appropriate definition, rather than some other definition.

The Ada packages which contain some of the operations listed below are not detailed in this abbreviated example binding. For such details, see ISO/IEC 8652:1995, *Information Technology* –

Programming Languages – Ada [11]. A full binding would include information regarding packages also for the operations that are not included in the Ada standard.

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

```
INT, Max(x, y)
max_I(x,y)
min_I(x,y)
                           INT, Min(x, y)
                                                                       *
max\_seq_I(xs)
                           Max(xs)
                           Min(xs)
min\_seq_I(xs)
                                      dim_I(x,y)
                           Dim(x, y)
power_I(x,y)
                           x ** y
shift2_I(x,y)
                           Shift2(x, y)
shift 10_I(x,y)
                           Shift10(x, y)
                           Sqrt(x)
sqrt_I(x)
divides_I(x,y)
                           x \neq 0 and then y \mod x = 0
even_I(x)
                           x \mod 2 = 0
odd_I(x)
                           x \mod 2 \neq 0
                           Quotient(x, y)
quot_I(x,y)
mod_I(x,y)
                           x \mod y
ratio_I(x, y)
                           Ratio(x, y)
residue_I(x,y)
                           Residue(x, y)
                           Group(x, y)
group_I(x,y)
pad_I(x,y)
                           Pad(x, y)
                           Gcd(x, y)
gcd_I(x,y)
                           Lcm(x, y)
lcm_I(x,y)
                           Gcd(xs)
gcd\_seq_I(xs)
                           Lcm(xs)
lcm\_seq_I(xs)
                           Add_Wrap(x, y)
add\_wrap_I(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 array (Integer range \ll) of INT.

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

```
\begin{array}{llll} max_F(x,y) & FLT \operatorname{Max}(x,\ y) & \star \\ min_F(x,y) & FLT \operatorname{Min}(x,\ y) & \star \\ mmax_F(x,y) & \operatorname{MMax}(x,\ y) & \dagger \\ mmin_F(x,y) & \operatorname{MMin}(x,\ y) & \dagger \\ max\_seq_F(xs) & \operatorname{Max}(xs) & \dagger \\ min\_seq_F(xs) & \operatorname{Min}(xs) & \dagger \\ \end{array}
```

C.1 Ada 105

```
mmax\_seq_F(xs)
                             MMax(xs)
mmin\_seq_F(xs)
                             MMin(xs)
                                                                           †
                             Dim(x, y)
dim_F(x,y)
                             FLT'Floor(x)
floor_F(x)
floor\_rest_F(x)
                             x - FLT'Floor(x)
rounding_F(x)
                             FLT'Unbiased_Rounding(x)
                             x - FLT, Unbiased_Rounding(x)
rounding\_rest_F(x)
                             FLT'Ceiling(x)
ceiling_F(x)
ceiling\_rest_F(x)
                             x - FLT'Ceiling(x)
                             FLT'Remainder(x, y)
residue_F(x,y)
sqrt_F(x)
                             Sqrt(x)
                             Rec\_Sqrt(x)
rec\_sqrt_F(x)
mul_{F \to F'}(x, y)
                             Prod(x, y)
add\_lo_F(x,y)
                             Add_Low(x, y)
sub\_lo_F(x,y)
                             Sub\_Low(x, y)
mul\_lo_F(x,y)
                             Mul_Low(x, y)
                             Div_Rest(x, y)
div\_rest_F(x,y)
sqrt\_rest_F(x)
                             Sqrt_Rest(x)
```

where x and y are expressions of type FLT, and where xs is an expression of type array (Integer range <>) of FLT.

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

$max_error_hypot_F$	$\operatorname{Err}_{Hypotenuse}(x)$	†
$max_error_exp_F$ $max_error_power_F$	Err_Exp(x) Err_Power(x)	† †
big_angle_r _F max_error_rad _F max_error_sin _F max_error_tan _F	Big_Radian_Angle(x) Err_Rad(x) Err_Sin(x) Err_Tan(x)	† † † †
$min_angular_unit_F$ $big_angle_u_F$ $max_error_sinu_F(u)$ $max_error_tanu_F(u)$	<pre>Smallest_Angular_Unit(x) Big_Angle(x) Err_Sin_Cycle(u) Err_Tan_Cycle(u)</pre>	† † † †
$max_error_sinh_F$ $max_error_tanh_F$	<pre>Err_Sinh(x) Err_Tanh(x)</pre>	† †
$max_error_convert_F$ $max_error_convert_{D'}$ $max_error_convert_{D'}$	<pre>Err_Convert(x) Err_Convert_To_String Err_Convert_To_String</pre>	† † †

where x and u are expressions of type FLT, and F' and D' are non-special value sets for string formats. 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. (This is in Ada normally done as 'type attributes', but new such cannot be defined outside of the Ada standard itself.)

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

```
hypot_F(x,y)
                            Hypotenuse(x, y)
                                                                           †
power_{F,I}(b,z)
                             b ** z
                             Exp(x)
exp_F(x)
expm1_F(x)
                             ExpM1(x)
exp2_F(x)
                             Exp2(x)
exp10_F(x)
                             Exp10(x)
power_F(b, y)
                             b ** y
power1pm1_F(b, y)
                             Power1PM1(b, y)
ln_F(x)
                            Log(x)
ln1p_F(x)
                            Log1P(x)
log 2_F(x)
                            Log2(x)
log 10_F(x)
                            Log10(x)
                                         (note parameter order)
logbase_F(b, x)
                            Log(x, b)
                            Log1P1P(x, b)
logbase1p1p_F(b, x)
                             Rad(x)
rad_F(x)
                             Rad(x, h, v)
axis\_rad_F(x)
                                              (note out parameters)
                             Sin(x)
sin_F(x)
cos_F(x)
                             Cos(x)
tan_F(x)
                             Tan(x)
                             Cot(v)
cot_F(x)
                             Sec(x)
sec_F(x)
                             Csc(x)
csc_F(x)
                             CosSin(x, c, s)
cossin_F(x)
                                                (note out parameters)
                             ArcSin(x)
arcsin_F(x)
arccos_F(x)
                             ArcCos(x)
arctan_F(x)
                             ArcTan(x)
arccot_F(x)
                             ArcCotS(x)
arccotc_F(x)
                             ArcCot(x)
arcsec_F(x)
                             ArcSec(x)
arccsc_F(x)
                             ArcCsc(x)
arc_F(x,y)
                             ArcTan(y, x) or ArcCot(x, y)
                                                                           ★(invalid at origin)
cycle_F(u,x)
                             Cycle(x, u) (note parameter order)
axis\_cycle_F(u,x)
                             Cycle(x, u, h, v)
sinu_F(u,x)
                            Sin(x, u)
                                         (note parameter order)
cosu_F(u,x)
                             Cos(x, u)
tanu_F(u,x)
                            Tan(x, u)
```

C.1 Ada 107

```
cotu_F(u,x)
                            Cot(x, u)
                                                                          †
secu_F(u,x)
                            Sec(x, u)
cscu_F(u,x)
                            Csc(x, u)
cossinu_F(u,x)
                            CosSin(x, u, c, s)
                            ArcSin(x, u)
arcsinu_F(u,x)
arccosu_F(u,x)
                            ArcCos(x, u)
                            ArcTan(x, Cycle=>u)
arctanu_F(u,x)
arccotu_F(u,x)
                            ArcCotS(x, u)
                            ArcCot(x, Cycle=>u)
arccotcu_F(u,x)
arcsecu_F(u,x)
                            ArcSec(x, u)
arccscu_F(u, x)
                            ArcCsc(x, u)
                                                                          ★(invalid at origin)
arcu_F(u, x, y)
                            ArcTan(y, x, u) or ArcCot(x, y, u)
                                        FUIL POF OF ISOINE CAOSE
rad\_to\_cycle_F(x, w)
                            Rad_to_Cycle(x, w)
cycle\_to\_rad_F(u,x)
                            Cycle_to_Rad(u, x)
cycle\_to\_cycle_F(u, x, w)
                            Cycle_to_Cycle(u, x, w)
sinh_F(x)
                            SinH(x)
                            CosH(x)
cosh_F(x)
tanh_F(x)
                            TanH(x)
coth_F(x)
                            CotH(x)
sech_F(x)
                            SecH(x)
csch_F(x)
                            CscH(x)
arcsinh_F(x)
                            ArcSinH(x)
arccosh_F(x)
                            ArcCosH(x)
                            ArcTanH(x)
arctanh_F(x)
arccoth_F(x)
                            ArcCotH(x)
arcsech_F(x)
                            ArcSecH(x)
arccsch_F(x)
                            ArcCscH(x)
```

where b, x, y, u, and w are expressions of type FLT, z is an expression of type INT, and c, s, h, and v are variables of type FLT.

Ada95 specifies (in other words) that $power_{F,I}$ must be computed by repeated multiplication (in an unspecified order). That computation method cannot, whatever the order of multiplications, guarantee fulfillment of the LIA-2 accuracy requirements, and cannot fulfill the required relationship between $power_{F,I}$ and $power_{F}$. Further, Ada95 specifies that angular units must be positive, and implicitly has a value for $min_angular_unit_F$ of $fmin_F$. LIA-2 allows also negative angular units, but has a larger value for $min_angular_unit_F$. A real Ada binding for LIA-2 must state how these conflicts are resolved (see clause 2).

Arithmetic value conversions in Ada are always explicit and usually use the destination datatype name as the name of the conversion function, except when converting to/from string formats.

```
convert_{I \rightarrow I'}(x) INT2(x) \star convert_{I'' \rightarrow I}(s) Get(s, n, w); \star convert_{I'' \rightarrow I}(f) Get(f?, n, w?); \star convert_{I \rightarrow I''}(x) Put(s, x, base?); \star convert_{I \rightarrow I''}(x) Put(h?, x, w?, base?);
```

```
floor_{F \to I}(y)
                                INT(FLT'Floor(y))
rounding_{F \to I}(y)
                                INT(FLT', Unbiased_Rounding(y))
                                INT(FLT', Ceiling(y))
ceiling_{F \to I}(y)
convert_{I \to F}(x)
                                FLT(x)
convert_{F \to F'}(y)
                                FLT2(y)
convert_{F''\to F}(s)
                                Get(s, n, w?);
                                Get(f?, n, w?);
convert_{F''\to F}(f)
                                Put(s, y, Aft=>a?, Exp=>e?);
convert_{F \to F''}(y)
                                Put(h?, y, Fore=>i?, Aft=>a?, Exp=>e?);
convert_{F \to F''}(y)
                                FLT(z)
convert_{D \to F}(z)
convert_{D' \to F}(s)
                                Get(s, n, w?);
convert_{D' \to F}(f)
                                Get(f?, n, w?);
convert_{F\to D}(y)
                                FXD(y)
                               Put(s, y, Aft=>a?, Exp=>0);
convert_{F\to D'}(y)
                                Put(h?, y, Fore=>i?, Aft\rightleftharpoons>n?, Exp=>0);
convert_{F \to D'}(y)
```

Ada provides non-negative numerals for all its integer and floating point types. The default base is 10, but all bases from 2 to 16 can be used. There is no differentiation between the numerals for different floating point types, nor between numerals for different integer types, but integer numerals (without a point) cannot be used for floating point types, and 'real' numerals (with a point) cannot be used for integer types. Integer numerals can have an exponent part though. The details are not repeated in this example binding, see ISO/IEC 8652:1995, clause 2.4 Numeric Literals, clause 3.5.4 Integer Types, and clause 3.5.6 Real Types.

The Ada standard does not specify any numerals for infinities and NaNs. Suggestion:

+∞	FLT'Infinity	†
qNaN	FLT'NaN	†
sNaN	FLT'NaNSignalling	†

as well as string formats for reading and writing these values as character strings.

Ada has a notion of 'exception' that implies a non-returnable, but catchable, change of control flow. Ada uses its exception mechanism as its default means of notification. **underflow** does not cause any notification in Ada, and the continuation value to the **underflow** is used directly,

C.1 Ada 109

since an Ada exception is inappropriate for an **underflow** notification. On **underflow** the continuation value (specified in LIA-2) is used directly without recording the **underflow** itself. Ada uses the exception Constraint_Error for **infinitary** and **overflow** notifications, and the exceptions Numerics.Argument_Error, IO_Exceptions.Data_Error, and IO_Exceptions.End_Error for **invalid** notifications. Since Ada exceptions are non-returnable changes of control flow, no continuation value is provided for these 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.2 BASIC

The programming language BASIC is defined by ANSI X3.113-1987 (R1998) [40], endorsed by ISO/IEC 10279:1991, Information technology – Programming languages – Full BASIC [16].

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.

BASIC has no user accessible datatype corresponding to the LIA datatype Boolean.

BASIC has one primitive computational datatype, numeric. The model presented by the BASIC language is that of a real number with decimal radix and a specified (minimum) number of significant decimal digits. Numeric data is not declared directly, but any special characteristics are inferred from how they are used and from any OPTIONS that are in force.

The BASIC statement OPTION ARITHMETIC MATIVE ties the numeric type more closely to the underlying implementation. The precision and type of NATIVE numeric data is implementation dependent.

For the trigonometric operations, if OPTION ANGLE DEGREES is in effect, the argument or result is given in degrees. If OPTION ANGLE RADIANS (default) is in effect, the argument or result is given in radians.

Since the BASIC numeric datatype does not match LIA-1 integer datatypes, this binding example does not include any of the LIA-2 operations for integer datatypes.

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

```
max_F(x,y)
                              MAX(x, y)
min_F(x,y)
                              MIN(x, y)
mmax_F(x,y)
                              MMAX(x, y)
mmin_F(x,y)
                              MMIN(x, y)
max\_seq_F(xs)
                              MAXS(xs)
                              MINS(xs)
min\_seq_F(xs)
mmax\_seq_F(xs)
                              \texttt{MMAXS}(xs)
mmin\_seq_F(xs)
                              \mathtt{MMINS}(xs)
dim_F(x,y)
                              MONUS(x, y)
floor_F(x)
                              INT(x)
```

```
floor\_rest_F(x)
                              x - INT(x)
rounding_F(x)
                              ROUND(x)
rounding\_rest_F(x)
                              x - ROUND(x)
ceiling_F(x)
                              \mathtt{CEIL}(x)
ceiling\_rest_F(x)
                              x - CEIL(x)
residue_F(x,y)
                              RESIDUE(x, y)
sqrt_F(x)
                              SQR(x)
rec\_sqrt_F(x)
                              REC\_SQRT(x)
add\_lo_F(x,y)
                              ADD_LOW(x, y)
sub\_lo_F(x,y)
                              SUB\_LOW(x, y)
mul\_lo_F(x,y)
                              MUL\_LOW(x, y)
div\_rest_F(x,y)
                              DIV_REST(x, y)
sqrt\_rest_F(x)
                              SQRT_REST(x)
```

where x and y are expressions of type numeric, and where xs is an expression of type array of numeric.

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

$max_error_hypot_F$	ERR_HYPOTENUSE	†
$max_error_exp_F$ $max_error_power_F$	ERR_EXP ERR_POWER	† †
$big_angle_r_F$	BIG_RADIAN_ANGLE	†
$max_error_rad_F$	ERR_RAD (**)	†
$max_error_sin_F$	ERR_SIN CONTROL OF THE PROPERTY OF THE PROPERT	†
$max_error_tan_F$	ERR_TAN	†
$min_angular_unit_F$	MIN_ANGLE_UNIT	†
$big_angle_u_F$	BIG_ANGLE	†
$max_error_sinu_F(u)$	ERR_SIN_CYCLE(u)	†
$max_error_tanu_F(u)$	$ERR_TAN_CYCLE(u)$	†
$max_error_sinh_{\mathbf{F}}$	ERR_SINH	†
$max_error_tanh_F$	ERR_TANH	†
$max_error_convert_F$	ERR_CONVERT	†
$max_error_convert_{F'}$	ERR_CONVERT_TO_STRING	†
$max_error_convert_{D'}$	ERR_CONVERT_TO_STRING	†

where u is an expression of type numeric.

The LIA-2 floating point operations are listed below, along with the syntax used to invoke them. BASIC has a degree mode and a radian mode for the trigonometric operations.

$hypot_F(x,y)$	HYPOT(x, y)	†
$exp_F(x)$	EXP(x)	*
$expm1_F(x)$	EXPM1(x)	†

 $C.2\ BASIC$

```
exp2_F(x)
                            EXP2(x)
                                                                          †
                            EXP10(x)
exp10_F(x)
power_F(b, y)
                            POWER(b, y)
power1pm1_F(b, y)
                            POWER1PM1(b, y)
ln_F(x)
                            LOG(x)
ln1p_F(x)
                            LOG1P(x)
log 2_F(x)
                            LOG2(x)
log 10_F(x)
                            LOG10(x)
logbase_F(b, x)
                            LOGBASE(b, x)
logbase1p1p_F(b, x)
                            LOGBASE1P1P(b, x)
                                             (when in radian mode)
rad_F(x)
                            NORMANGLE(x)
sin_F(x)
                            SIN(x)
                                      (when in radian mode)
cos_F(x)
                            COS(x)
                                      (when in radian mode)
                                      (when in radian mode)
tan_F(x)
                            TAN(x)
cot_F(x)
                            COT(x)
                                      (when in radian mode)
                                      (when in radian mode)
sec_F(x)
                            SEC(x)
                                      (when in radian mode)
                            CSC(x)
csc_F(x)
                                       (when in radian mode)
arcsin_F(x)
                            ASIN(x)
                                       (when in radian mode)
arccos_F(x)
                            ACOS(x)
                                      (when in radian mode)
arctan_F(x)
                            ATN(x)
                                       (when in radian mode)
arccot_F(x)
                            ACOT(x)
arccotc_F(x)
                            ACOTC(x)
                                         (when in radian mode)
                            ASEC(x) when in radian mode)
arcsec_F(x)
arccsc_F(x)
                            ACSC(x)
                                       (when in radian mode)
arc_F(x,y)
                            ANGLE (x, y)
                                            (when in radian mode)
                                                                          ★(invalid at origin)
                            NORMANGLEU(u, x)
cycle_F(u,x)
                            SINU(u, x)
sinu_F(u,x)
cosu_F(u,x)
                            COSU(u, x)
tanu_F(u,x)
                            TANU(u, x)
cotu_F(u,x)
                            COTU(u, x)
secu_F(u,x)
                            SECU(u, x)
                            CSCU(u, x)
cscu_F(u,x)
arcsinu_F(u,x)
                            ASINU(u, x)
                            ACOSU(u, x)
arccosu_F(u,x)
arctanu_F(u,x)
                            ATNU(u, x)
                            ACOTU(u, x)
arccotu_F(u,x)
arccotcu_F(u, x)
                            ACOTCU(u, x)
arcsecu_F(u,x)
                            ASECU(u, x)
                            ACSCU(u, x)
arccscu_F(u,x)
arcu_F(u, x, y)
                            ANGLEU(u, x, y)
```

```
cycle_F(360,x)
                                  NORMANGLE(x)
                                                    (when in degree mode)
     sinu_F(360,x)
                                  SIN(x)
                                             (when in degree mode)
     cosu_F(360,x)
                                  COS(x)
                                             (when in degree mode)
     tanu_F(360,x)
                                             (when in degree mode)
                                  TAN(x)
     cotu_F(360,x)
                                  COT(x)
                                             (when in degree mode)
     secu_F(360,x)
                                  SEC(x)
                                             (when in degree mode)
     cscu_F(360,x)
                                             (when in degree mode)
                                  CSC(x)
     arcsinu_F(360,x)
                                              (when in degree mode)
                                  ASIN(x)
                                              (when in degree mode)
     arccosu_F(360,x)
                                  ACOS(x)
     arctanu_F(360,x)
                                  ATN(x)
                                             (when in degree mode)
     arccotu_F(360,x)
                                  ACOT(x)
                                              (when in degree mode)
                                               (when in degree mode)
     arccotcu_F(360,x)
                                  ACOTC(x)
                                              (when in degree mode)
     arcsecu_F(360, x)
                                  ASEC(x)
     arccscu_F(360, x)
                                  ACSC(x)
                                              (when in degree mode)
                                                  (when in degree mode)
                                                                                 ★(invalid at origin)
     arcu_F(360, x, y)
                                  ANGLE(x, y)
     rad\_to\_cycle_F(x, 360)
                                  DEG(x)
     cycle\_to\_rad_F(360, x)
                                  RAD(x)
     rad\_to\_cycle_F(x, w)
                                  RAD_TO_CYCLE(x, w)
     cycle\_to\_rad_F(u, x)
                                  CYCLE_TO_RAD(u, x)
     cycle\_to\_cycle_F(u, x, w)
                                  CYCLE_TO_CYCLE(u
     sinh_F(x)
                                  SINH(x)
     cosh_F(x)
                                  COSH(x)
     tanh_F(x)
                                  TANH(x)
                                  COTH (x)
     coth_F(x)
     sech_F(x)
                                  SECH(x)
                                   CSCH(x)
     csch_F(x)
                                  ARCSINH(x)
     arcsinh_F(x)
     arccosh_F(x)
                                  ARCCOSH(x)
     arctanh_F(x)
                                  ARCTANH(x)
     arccoth_F(x)
                                  ARCCOTH(x)
     arcsech_F(x)
                                  ARCSECH(x)
     arccsch_F(x)
                                  ARCCSCH(x)
where b, x, y, d, and w are expressions of type numeric.
```

Arithmetic value conversions in BASIC are always tied to reading and writing text.

```
convert_{F''\to F}(stdin)
                                      \mathtt{READ} \ x
convert_{F \to F''}(y)
                                      PRINT y
convert_{D' \to F}(stdin)
                                      READ x
```

where x is a variable of type numeric, and y is an expression of type numeric.

BASIC provides non-negative numerals for numeric in base 10.

BASIC does not specify any numerals for infinities and NaNs. Suggestion:

C.2 BASIC 113

$+\infty$	INFINITY	†
${ m qNaN}$	NAN	†
sNaN	NANSIGNALLING	†

as well as string formats for reading and writing these values as character strings.

BASIC has a notion of 'exception' that implies a non-returnable change of control flow. BASIC uses its exception mechanism as its default means of notification. **underflow** does not cause any notification in BASIC, and the continuation value to the **underflow** is used directly, since an BASIC exception is inappropriate for an **underflow** notification. BASIC uses the exception numbers 1001 to 1008 for **overflow**, exception numbers 1502 and 1503 for *handled* **underflow**, the exception number 3001 to 3004 for **infinitary**, the exception numbers -3000, 3002, and 3004 to 3011 for **invalid**, and the exception numbers -3050 and -3051 for **absolute_precision_underflow**. Since BASIC exceptions are non-returnable changes of control flow, no continuation value is provided for these notifications, except that *unhandled* underflow uses the continuation value specified without any BASIC exception.

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.3 C

The programming language C is defined by ISO/IEC 9899:1999, Information technology – Programming languages – C [17]. This edition of the C standard is often referred to as C99, which is also used below.

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 LIA datatype **Boolean** is implemented by the C datatype int (1 = true and 0 = false) or the new C99 _Bool datatype.

C defines numerous integer datatypes. They may be aliases of each other in an implementation defined way. The description here is not complete. See the C99 standard. Some of the integer datatypes have a predetermined bit width, and the signed ones use 2's complement for representation of negative values: intn_t and uintn_t, where n is the bit width expressed as a decimal numeral. Some bit widths are required. There are also minimum width, fastest minimum width, and special purpose integer datatypes (like size_t). Also provided are the more well-known integer datatypes char, short int, int, long int, long long int (new in C99), unsigned char, unsigned short int, unsigned int, unsigned long int, and unsigned long long int (new in C99). Finally there are the integer datatypes intmax_t and uintmax_t may even be unbounded with a negative integer infinity as INTMAX_MIN and a positive integer infinity as INTMAX_MAX and UINTMAX_MAX. INT is used below to designate one of the integer datatypes.

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

For some of the operations below, the C standard defines 'type generic macros'. These are fixed by the C standard, and new ones cannot be defined in program libraries.

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

```
imaxt(x, y)
max_I(x,y)
                             (x < y ? y :
max_I(x,y)
                                               x)
                                                                            *
                             imint(x, y)
min_I(x,y)
                             (x < y ? x : y)
min_I(x,y)
                             imax_arrt(xs, nr_of_items)
max\_seq_I(xs)
min\_seq_I(xs)
                             imin_arrt(xs, nr_of_items)
dim_I(x,y)
                             idimt(x, y)
power_I(x,y)
                             ipowert(x, y)
                             shift2t(x, y)
shift2_I(x,y)
shift 10_I(x,y)
                             shift10t(x, y)
sqrt_I(x)
                             isqrtt(x)
divides_I(x,y)
                             does\_dividet(x, y)
divides_I(x,y)
                             x = 0 & y & x = 0
even_I(x)
                             x \% 2 == 0
                             x \% 2 != 0
odd_I(x)
quot_I(x,y)
                             quott(x, y)
                             modt(x, y)
mod_I(x,y)
ratio_I(x,y)
                             ratiot(x, y)
                             iremaindert(x)
residue_I(x,y)
group_I(x,y)
                             groupt(x, y)
pad_I(x,y)
                             padt(x,
gcd_I(x,y)
                             gcdt(x, y)
                             \mathbf{Tcm}t(x, y)
lcm_I(x,y)
                             gcd_arrt(xs, nr_of_items)
gcd\_seq_I(xs)
                             lcm_arrt(xs, nr_of_items)
lcm\_seq_I(xs)
add\_wrap_I(x, y)
                             add\_wrapt(x, y)
add\_ov_I(x,y)
                             add\_overt(x, y)
sub\_wrap_I(x,y)
                             sub\_wrapt(x, y)
sub\_ov_I(x,y)
                             sub\_overt(x, y)
mul\_wrap_I(x,y)
                             mul\_wrapt(x, y)
mul\_ov_I(x,y)
                             mul\_overt(x, y)
```

where x and y are expressions of the same integer type and where xs is an expression of type array of an integer type. t is a string (part of the operation name in C), that is the empty string for int, is 1 for long int, is u for unsigned int, and is ul for unsigned long int.

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

C.3 C 115

```
mmax_F(x,y)
                            fmaxt(x, y) or fmax(x, y)
                                                                         ★(C99)
mmin_F(x,y)
                            fmint(x, y) or fmin(x, y)
                                                                         ★(C99)
max\_seq_F(xs)
                            nmax_arrt(xs, nr_of_items)
                            nmin_arrt(xs, nr_of_items)
                                                                         †
min\_seq_F(xs)
mmax\_seq_F(xs)
                            fmax_arrt(xs, nr_of_items)
                            fmin_arrt(xs, nr_of_items)
mmin\_seq_F(xs)
dim_F(x,y)
                            fdimt(x, y) or fdim(x, y)
                                                            (dev. for special values) \star(C99)
floor_F(x)
                            floort(x) or floor(x)
floor\_rest_F(x)
                            x - floort(x)
                            nearbyintt(x) (when in round to nearest mode) \star(C99)
rounding_F(x)
rounding\_rest_F(x)
                            x - nearbyint t(x) (when in round to nearest mode) \star (C99)
ceiling_F(x)
                            ceilt(x) or ceil(x)
ceiling\_rest_F(x)
                            x - ceilt(x)
                            remainder (x, y) or remainder (x, y)
residue_F(x,y)
sqrt_F(x)
                            sqrtt(x) or sqrt(x)
rec\_sqrt_F(x)
                            rec_sqrtt(x)
mul_{F \to F'}(x, y)
                            dprodt(x, y)
add\_lo_F(x,y)
                            add_lowt(x, y)
sub\_lo_F(x,y)
                            sub\_lowt(x, y)
mul\_lo_F(x,y)
                            mullowt(x, y)
div\_rest_F(x,y)
                            div_restt(x, y)
sqrt\_rest_F(x)
                            sqrt_restt(x)
```

where x and y are expressions of the same floating point type, and where xs is an expression of type array of a floating point type. t is a string (part of the operation name), that is the empty string for double, is f for float, and is l for long double (the same applies to the parameters and operations below).

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

```
max\_error\_hypot_F
                                  \mathtt{err\_hypot} t
max\_error\_exp_F
                                  err_expt
max_error_powers
                                  err_powert
big\_angle\_r_{F}
                                  big_radian_anglet
max\_error\_rad_F
                                  \mathtt{err\_rad}t
max\_error\_sin_F
                                  \mathtt{err\_sin} t
max_{error\_tan_F}
                                  err_tant
min\_angular\_unit_F
                                  smallest\_angle\_unitt
                                  {\tt big\_angle} t
big\_angle\_u_F
                                  err_sin_cyclet(u)
max\_error\_sinu_F(u)
max\_error\_tanu_F(u)
                                  err_tan_cyclet(u)
max\_error\_sinh_F
                                  \mathtt{err\_sinh} t
max\_error\_tanh_F
                                  \mathtt{err\_tanh} t
```

where u is an expression of a floating point type. No conversion error parameters are needed, since C99 requires all floating point datatype conversion (even to and from strings) to always have an error that is ≤ 0.5 ulp.

C has a pow operation that does not conform to LIA-2, but may be specified in LIA-2 terms:

```
pow_F(x,y) = power_{F\mathcal{Z}}(x,y) if y \in F \cap \mathcal{Z}

= pow_F(x,0) if y = -\mathbf{0}

= 1 if x \in \{-1,1\} and y \in \{-\infty, +\infty\}

= 1 if x is a quiet NaN and y = 0

= power_F(x,y) otherwise
```

C99 has a hypot operation that does not conform to LIA-2, but may be specified in LIA-2 terms:

```
hhypot_F(x,y) = +\infty if x is a quiet NaN and y \in \{-\infty, +\infty\} = +\infty if x \in \{-\infty, +\infty\} and y is a quiet NaN otherwise
```

The LIA-2 elementary floating point operations are listed below, together with the non-LIA-2 pow_F and $hhypot_F$, along with the syntax used to invoke them:

```
hypot_F(x,y)
                            hypotenuset(x, y)
                            hypott(x, y) or hypot(x, y)
hhypot_F(x,y)
                                                                          ★ Not LIA-2!
power_{F,I}(b,z)
                            powerit(b, z)
exp_F(x)
                            \exp t(x) or \exp(x)
expm1_F(x)
                            expm1t(x) or expm1(x)
                                                                          *(C99)
                            \exp 2t(x) or \exp 2(x)
                                                                          ★(C99)
exp2_F(x)
                            exp10t(x)
exp10_F(x)
                            powert(b, y)
power_F(b, y)
                            powt(b \lor y) or pow(b, y)
                                                                          ★ Not LIA-2!
pow_F(b,y)
power1pm1_F(b, y)
                            power1pm1t(b, y)
                            \log t(x) or \log(x)
ln_F(x)
                                                                          ★(C99)
ln1p_F(x)
                            log1pt(x) or log1p(x)
                                                                          ★(C99)
log 2_F(x)
                            log2t(x) or log2(x)
log 10_F(x)
                            log10t(x) or log10(x)
logbase_F(b,x)
                            logbaset(b, x)
logbase1p1p
                            logbase1p1pt(b, x)
                            radiant(x)
                            axis\_radt(x, \&h, \&v) (note out parameters)
sin_F(x)
                            sint(x)
                                          sin(x)
                                      or
cos_F(x)
                            cost(x)
                                          cos(x)
                                      or
tan_F(x)
                            tant(x)
                                          tan(x)
cot_F(x)
                            \cot t(x)
sec_F(x)
                            sect(x)
csc_F(x)
                            csct(x)
cossin_F(x)
                            cossint(x, \&c, \&s)
```

C.3 C 117

```
arcsin_F(x)
                             asint(x)
                                        or
                                             asin(x)
arccos_F(x)
                             acost(x)
                                             acos(x)
                                        or
arctan_F(x)
                             atant(x)
                                        or
                                             atan(x)
                             acott(x)
arccot_F(x)
arccotc_F(x)
                             acotct(x)
                             asect(x)
arcsec_F(x)
arccsc_F(x)
                             acsct(x)
arc_F(x,y)
                             atan2t(y, x)
                                             or atan2(y, x)
cycle_F(u,x)
                             cyclet(u, x)
axis\_cycle_F(u, x)
                             axis\_cyclet(u, x, &h, &v)
sinu_F(u,x)
                             sinut(u, x)
                             cosut(u, x)
cosu_F(u,x)
tanu_F(u,x)
                             tanut(u, x)
cotu_F(u,x)
                             cotut(u, x)
                             secut(u, x)
secu_F(u,x)
cscu_F(u,x)
                             cscut(u, x)
                             cossinut(u, x, &c, &s)
cossinu_F(u,x)
arcsinu_F(u,x)
                             asinut(u, x)
arccosu_F(u,x)
                             acosut(u, x)
                             atanut(u, x)
arctanu_F(u,x)
                             acotut(u, x)
arccotu_F(u,x)
                             acotcut(u, x)
arccctcu_F(u,x)
arcsecu_F(u,x)
                             asecut(u, x)
arccscu_F(u,x)
                             acscut(u, x)
                             atan2ut(u, y, x)
arcu_F(u, x, y)
rad\_to\_cycle_F(x, w)
                             radian_to_cyclet(x, w)
cycle\_to\_rad_F(u,x)
                             cycle_to_radiant(u, x)
cycle\_to\_cycle_F(u, x, w)
                             cycle_to_cyclet(u, x, w)
sinh_F(x)
                             sinht(x)
                                             sinh(x)
                                                                            ★(C99)
                                        or
cosh_F(x)
                                                                            ★(C99)
                             cosht(x)
                                        or
                                             cosh(x)
tanh_F(x)
                             tanht(x)
                                            tanh(x)
                                                                            \star(C99)
                                        or
coth_F(x)
                             cotht(x)
sech_F(x)
                             secht(x)
csch_{F}(x)
                             \operatorname{csch} t(x)
arcsinh_F(x)
                             asinht(x)
                                              asinh(x)
                                                                            ★(C99)
                                          or
arccosh_F(x)
                             acosht(x)
                                              acosh(x)
                                                                            ★(C99)
arctanh_F(x)
                             atanht(x)
                                                                            ★(C99)
                                          or
                                              atanh(x)
arccoth_F(x)
                             acotht(x)
arcsech_F(x)
                             asecht(x)
arccsch_F(x)
                             \operatorname{acsch} t(x)
```

where b, x, y, u, and w are expressions of type FLT, h, v, c, and s are lvalue expressions of type FLT, and z is an expression of type INT.

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 string formats. The rules for when implicit conversions are applied is not repeated here, but work as if a cast had been applied.

When converting to/from string formats, format strings are used. The format string is used as a pattern for the string format generated or parsed. The description of format strings here is not complete. Please see the C99 standard for a full description. In the format strings % is used to indicate the start of a format pattern. After the %, optionally a string field width (w below) may be given as a positive decimal integer numeral. For the floating and fixed point format patterns, there may then optionally be a '.' followed by a positive integer numeral (d below) indicating the number of fractional digits in the string. The C operations below use HYPHEN-MINUS rather than MINUS (which would have been typographically better), and only digits that are in ASCII, independently of so-called locale. For generating or parsing other kinds of digits, say Arabic digits or Thai digits, another API must be used, that is not standardised in C. For the floating and fixed point formats, $+\infty$ may be represented as either inf or infinity, $-\infty$ may be represented as either $-\inf$ or $-\inf$

For the integer formats then follows an internal type indicator, of which some are new to C99. Not all C99 integer types have internal type indicators. However, for t below, hh indicates char, h indicates short int, the empty string indicates int, 1 (the letter l) indicates long int, 11 (the letters ll) indicates long long int, and j indicates intmax_t or uintmax_t. (For system purposes there are also special type names like size_t, and z indicates size_t and t indicates ptrdiff_t as type format letters.) Finally, there is a radix (for the string side) and signedness (both sides) format letter (r below): d for signed decimal; o, u, x, X for octal, decimal, hexadecimal with small letters, and hexadecimal with capital letters, all unsigned. E.g., %jd indicates decimal numeral string for unsigned char, with a two character field width, and %lu indicates decimal numeral string for unsigned long int.

For the floating point formats instead follows another internal type indicator. Not all C99 floating point types have standard internal type indicators for the format strings. However, for u below the empty string indicates double and L indicates long double. Finally, there is a radix (for the string side) format letter: e or E for decimal, a or A for hexadecimal. E.g., 15.8LA indicates hexadecimal floating point numeral string for long double, with capital letters for the letter components, a field width of 15 characters, and 8 hexadecimal fractional digits.

For the fixed point formats also follows the internal type indicator as for the floating point formats. But for the final part of the pattern, there is another radix (for the string side) format letter (p below), only two are standardised, both for the decimal radix: f or F. E.g., %Lf indicates decimal fixed point numeral string for long double, with a small letter for the letter component. (There is also a combined floating/fixed point string format: g.)

```
\begin{array}{llll} convert_{I \to I'}(x) & (INT2)x & \star \\ convert_{I'' \to I}(s) & \operatorname{sscanf}(s, \text{"}\%wtr\text{"}, \&i) & \star \\ convert_{I'' \to I}(f) & \operatorname{fscanf}(f, \text{"}\%wtr\text{"}, \&i) & \star \\ convert_{I \to I''}(x) & \operatorname{sprintf}(s, \text{"}\%wtr\text{"}, x) & \star \\ convert_{I \to I''}(x) & \operatorname{fprintf}(h, \text{"}\%wtr\text{"}, x) & \star \\ floor_{F \to I}(y) & (INT) \operatorname{floort}(y) & \star \\ \end{array}
```

C.3 C

```
floor_{F \to I}(y)
                               (INT) nearbyintt(y) (when in round towards -\infty mode) \star(C99)
rounding_{F \to I}(y)
                               (INT)nearbyintt(y)
                                                        (when in round to nearest mode) \star(C99)
ceiling_{F \to I}(y)
                               (INT) nearbyintt(y) (when in round towards +\infty mode) \star(C99)
ceiling_{F \to I}(y)
                               (INT) ceilt(y)
convert_{I \to F}(x)
                               (FLT)x
convert_{F \to F'}(y)
                               (FLT2)y
                               sscanf(s, "%w.duv", &r)
convert_{F''\to F}(s)
                               fscanf(f, "%w.duv", &r)
convert_{F''\to F}(f)
                               sprintf(s, "%w.duv", y)
convert_{F \to F''}(y)
                               fprintf(h, "%w.duv", y)
convert_{F \to F''}(y)
                               sscanf(s, "%wup", \&g)
convert_{D' \to F}(s)
                               fscanf(f, "%wup", &q)
convert_{D'\to F}(f)
                               sprintf(s, "%w.dup", y)
convert_{F \to D'}(y)
convert_{F \to D'}(y)
                               fprintf(h, "%w.dup", y)
```

where s is an expression of type char*, f is an expression of type FILE*, i is an lvalue expression of type int, g is an lvalue expression of type double, x is an expression of type INT, y is an expression of type FLT, INT2 is the integer datatype that corresponds to I', and FLT2 is the floating point datatype that corresponds to F'.

C99 provides non-negative numerals for all its integer and floating point types. The default base is 10, but base 8 (for integers) and 16 (both integer and float) can be used too. 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, 1 for long double. Numerals for floating point types must have a '.' or an exponent in them. The details are not repeated in this example binding, see ISO/IEC 9899:1999, clause 6.4.4.1 Integer constants, and clause 6.4.4.2 Floating constants.

C99 specifies numerals (as macros) for infinities and NaNs for float:

$+\infty$	c)	INFINITY	*	
qNaN		NAN	*	
sNaN		NANSIGNALLING	†	

as well as string formats for reading and writing these values as character strings.

C99 has two ways of handling arithmetic errors. One, for backwards compatibility, is by assigning to error. The other is by recording of indicators, the method preferred by LIA-2, which can be used for floating point errors. For C99, the **absolute_precision_underflow** notification is ignored. The behaviour when integer operations initiate a notification is, however, not defined by C99.

C.4 C++

The programming language C++ is defined by ISO/IEC 14882:1998, Programming languages – C++ [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.

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

The LIA 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:

```
For 180 IEC 1096
                             max(x, y)
max_I(x,y)
min_I(x,y)
                             min(x, y)
max\_seq_I(xs)
                             xs.max()
min\_seq_I(xs)
                             xs.min()
                             dim(x, y)
dim_I(x,y)
power_I(x, y)
                             power(x, y)
                             shift2(x, y)
shift 2_I(x,y)
                             shift10(x, y)
shift 10_I(x,y)
                             sqrt(x)
sqrt_I(x)
                             does_divide(x, y)
divides_I(x,y)
                             x != 0 &  y % x == 0
divides_I(x,y)
                             x \% == 0
even_I(x)
                             2 != 0
odd_I(x)
quot_I(x,y)
                             quot(x, y)
mod_I(x,y)
                             mod(x, y)
ratio_I(x,y)
                             ratio(x, y)
residue_I(x,y)
                             remainder (x, y)
                             group(x, y)
group_I(x, y)
                             pad(x, y)
pad_I(x,y)
gcd_{\mathbf{I}}(x,y)
                             gcd(x, y)
lcm_I(x,y)
                             lcm(x, y)
gcd\_seq_I(xs)
                             xs.gcd()
lcm\_seq_I(xs)
                             xs.lcm()
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)
```

C.4 C++121

```
mul\_ov_I(x,y) mul\_over(x, y)
```

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:

```
max_F(x,y)
                             nmax(x, y)
min_F(x,y)
                             nmin(x, y)
mmax_F(x,y)
                             \max(x, y)
mmin_F(x,y)
                             min(x, y)
max\_seq_F(xs)
                             xs.nmax()
min\_seq_F(xs)
                             xs.nmin()
mmax\_seq_F(xs)
                             xs.max()
mmin\_seq_F(xs)
                             xs.min()
                                                DF of 15011EC 10967
dim_F(x,y)
                             dim(x, y)
floor_F(x)
                             floor(x)
floor\_rest_F(x)
                             x - floor(x)
rounding_F(x)
                             round(x)
rounding\_rest_F(x)
                             x - \text{round}(x)
ceiling_F(x)
                             ceil(x)
                             x - ceil(x)
ceiling\_rest_F(x)
residue_F(x,y)
                             remainder (x, y)
sqrt_F(x)
                             sqrt(x)
                             rec\_sqrt(x)
rec\_sqrt_F(x)
mul_{F \to F'}(x, y)
                             dprod(x, y)
add\_lo_F(x,y)
                             add_low(x, y)
sub\_lo_F(x,y)
                             sub_low(x, y)
mul\_lo_F(x,y)
                             mullow(x, y)
div\_rest_F(x,y)
                             div_rest(x, y)
sqrt\_rest_F(x)
                            \neg sqrt\_rest(x)
```

where x and y are expressions of the same floating point type, and where xs is an expression of type valarray of a floating point type. The C++ standard does not make clear how to handle NaN arguments, in particular for max and min.

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

```
\begin{array}{lll} max\_error\_hypot_F & \text{numeric\_limits} < FLT > :: err\_hypotenuse() & \dagger \\ max\_error\_exp_F & \text{numeric\_limits} < FLT > :: err\_exp() & \dagger \\ max\_error\_power_F & \text{numeric\_limits} < FLT > :: err\_power() & \dagger \\ big\_angle\_r_F & \text{numeric\_limits} < FLT > :: big\_radian\_angle() & \dagger \\ max\_error\_rad_F & \text{numeric\_limits} < FLT > :: err\_rad() & \dagger \\ max\_error\_sin_F & \text{numeric\_limits} < FLT > :: err\_sin() & \dagger \\ max\_error\_tan_F & \text{numeric\_limits} < FLT > :: err\_tan() & \dagger \\ \end{array}
```

```
numeric_limits<FLT>::smallest_angle_unit()
     min\_angular\_unit_F
                                 numeric_limits<FLT>::big_angle()
     big\_angle\_u_F
                                 numeric_limits<FLT>::err_sin_cycle(u)
     max\_error\_sinu_F(u)
     max\_error\_tanu_F(u)
                                 numeric_limits<FLT>::err_tan_cycle(u)
     max\_error\_sinh_F
                                 numeric_limits<FLT>::err_sinh()
                                 numeric_limits<FLT>::err_tanh()
     max\_error\_tanh_F
                                 numeric_limits<FLT>::err_convert()
     max\_error\_convert_F
                                 numeric_limits<FLT>::err_convert_to_string()
     max\_error\_convert_{F'}
                                 numeric_limits<FLT>::err_convert_to_string()
     max\_error\_convert_{D'}
where u is an expression of a floating point type.
                                              th the of Isolific 10961.2,
   The LIA-2 elementary floating point operations are listed below, along with the syntax used
to invoke them:
                                 hypotenuse(x, y)
     hypot_F(x,y)
     power_{F,I}(b,z)
                                 power(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)
                                  power(b, y)
                                 pow(b, y)
                                                                               ★ Not LIA-2! (See C.)
     pow_F(b, y)
     power1pm1_F(b, y)
                                  power1pm1(b
     ln_F(x)
                                  log(x)
     ln1p_F(x)
                                  log1p(x)
                                  log2(x)
     log 2_F(x)
                                  log10(x)
     log 10_F(x)
     logbase_F(b, x)
                                  logbase(b, x)
     logbase1p1p_F(b, x)
                                  logbase1p1p(b, x)
     rad_F(x)
                                  rad(x)
                                  axis\_rad(x, \&h, \&v) (note out parameters) †
                                  sin(x)
     cos_F(x
                                  cos(x)
     tan_F(x)
                                  tan(x)
     cot_{E}(x)
                                  cot(x)
     sec_F(x)
                                  sec(x)
     csc_F(x)
                                  csc(x)
     cossin_F(x)
                                  cossin(x, \&c, \&s)
```

C.4 C++ 123

asin(x)

acos(x)

atan(x)

acot(x)

 $arcsin_F(x)$ $arccos_F(x)$

 $arctan_F(x)$

 $arccot_F(x)$

```
arccotc_F(x)
                            acotc(x)
                                                                           †
                            asec(x)
arcsec_F(x)
arccsc_F(x)
                            acsc(x)
                            atan2(y, x)
arc_F(x,y)
cycle_F(u,x)
                            cycle(u, x)
axis\_cycle_F(u,x)
                            axis\_cycle(u, x, &h, &v)
sinu_F(u,x)
                            sinu(u, x)
cosu_F(u,x)
                            cosu(u, x)
tanu_F(u,x)
                            tanu(u, x)
                            cotu(u, x)
cotu_F(u,x)
secu_F(u,x)
                            secu(u, x)
                            cscu(u, x)
cscu_F(u,x)
                            cossinu(u, x, &c, &s)
cossinu_F(x)
arcsinu_F(u,x)
                            asinu(u, x)
arccosu_F(u,x)
                            acosu(u, x)
                            atanu(u, x)
arctanu_F(u,x)
                            acotu(u, x)
arccotu_F(u,x)
                            acotcu(u, x)
arccotcu_F(u,x)
                            asecu(u, x)
arcsecu_F(u,x)
arccscu_F(u,x)
                            acscu(u, x)
                            atan2u(u, y, x)
arcu_F(u, x, y)
rad\_to\_cycle_F(x, w)
                            radian_to_cycle(x, w)
cycle\_to\_rad_F(u,x)
                            cycle_{to}_{radian}(u, x)
cycle\_to\_cycle_F(u, x, w)
                             cycle_to\cycle(u, x, w)
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)
```

where b, x, y, u, and w are expressions of type FLT, h, v, c, and s are lvalue expressions of type FLT, and z is an expression of type INT.

Arithmetic value conversions in C++ can be explicit or implicit. The rules for when implicit conversions are applied are not repeated here. C++ also deals with stream input/output in other ways, see clause 22.2.2 of ISO/IEC 14882:1998, 'Locale and facets'. The explicit arithmetic value

conversions are usually expressed as 'casts', except when converting to/from string formats.

When converting to/from string formats, format strings are used. The format string is used as a pattern for the string format generated or parsed. The description of format strings here is not complete. Please see the C++ standard for a full description. In the format strings % is used to indicate the start of a format pattern. After the %, optionally a string field width (w below) may be given as a positive decimal integer numeral. For the floating and fixed point format patterns, there may then optionally be a '.' followed by a positive integer numeral (d below) indicating the number of fractional digits in the string. The C++ operations below use HYPHEN-MINUS rather than MINUS (which would have been typographically better), and only digits that are in ASCII, independently of so-called locale. For generating or parsing other kinds of digits, say Arabic digits or Thai digits, another API must be used, that is not standardised in C++. For the floating and fixed point formats, $+\infty$ may be represented as either inf or infinity, $-\infty$ may be represented as either -inf or -infinity, and a NaN may be represented as NaN; all independently of so-called locale. For language dependent representations of these values another API must be used, that is not standardised in C.

For the integer formats then follows an internal type indicator. For t below, the empty string indicates int, 1 (the letter l) indicates long int. Finally, there is a radix (for the string side) and signedness (both sides) format letter (r below): d for signed decimal; o, u, x, X for octal, decimal, hexadecimal with small letters, and hexadecimal with capital letters, all unsigned. E.g., %d indicates decimal numeral string for int and %lu indicates decimal numeral string for unsigned long int.

For the floating point formats instead follows another internal type indicator. For u below the empty string indicates double and L indicates long double. Finally, there is a radix (for the string side) format letter: e or E for decimal. E.g. 15.8LE indicates hexadecimal floating point numeral string for long double, with a capital letter for the letter component, a field width of 15 characters, and 8 hexadecimal fractional digits.

For the fixed point formats also follows the internal type indicator as for the floating point formats. But for the final part of the pattern, there is another radix (for the string side) format letter (p below), only two are standardised, both for the decimal radix: f or F. E.g., %Lf indicates decimal fixed point numeral string for long double, with a small letter for the letter component. (There is also a combined floating/fixed point string format: g.)

```
static\_cast < INT2 > (x)
convert_{I \to I'}(x)
                                                                                   *
convert_{I''\to I}(s)
                                sscanf(s, "%wtr", \&i)
convert_{I''\to I}(f)
                                fscanf(f, "%wtr", &i)
                                sprintf(s, "%wtr", x)
convert_{I \to I}(x)
                                fprintf(h, "%wtr", x)
convert_{I \to I''}(x)
floor_{P \to I}(y)
                                static_cast<INT>(floor(y))
                                static_cast<INT>(round(y))
rounding_{F\to I}(y)
                                static_cast<INT>(ceil(y))
ceiling_{F\to I}(y)
                                static\_cast < FLT > (x)
convert_{I \to F}(x)
convert_{F \to F'}(y)
                                (FLT2)y
convert_{F''\to F}(s)
                                sscanf(s, "%w.duv", &r)
                                fscanf(f, "%w.duv", &r)
convert_{F''\to F}(f)
```

C.4 C++ 125

```
convert_{F 	o F''}(y) sprintf(s, "\%w.duv", y) \star convert_{F 	o F''}(y) fprintf(h, "\%w.duv", y) \star convert_{D' 	o F}(s) sscanf(s, "\%wup", &g) \star convert_{D' 	o F}(f) fscanf(f, "\%wup", &g) \star convert_{F 	o D'}(y) sprintf(s, "\%w.dup", y) \star convert_{F 	o D'}(y) fprintf(h, "\%w.dup", y) \star
```

where s is an expression of type char*, f is an expression of type FILE*, i is an expression of type int, g is an expression of type double, x is an expression of type INT, y is an expression of type FLT, INT2 is the integer datatype that corresponds to I', and FLT2 is the floating point datatype that corresponds to F'.

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, 1 for long double. Numerals for floating point types must have a '.' or an exponent in them. The details are not repeated in this example binding, see ISO/IEC 14882:1998, clause 2.9.1 Integer literals, and clause 2.9.4 Floating literals.

C++ specifies numerals for infinities and NaNs:

```
\begin{array}{lll} +\infty & & \text{numeric\_limits} < FLT > :: infinity() & \star \\ \text{qNaN} & & \text{numeric\_limits} < FLT > :: quiet\_NaN() & \star \\ \text{sNaN} & & \text{numeric\_limits} < FLT > :: signaling\_NaN() & \star \\ \end{array}
```

as well as string formats for reading and writing these values as character strings.

C++ has completely undefined behaviour on arithmetic notification. An implementation wishing to conform to LIA-2 should provide a means for recording of indicators, similar to C99.

C.5 Fortran

The programming language Fortran is defined by ISO/IEC 1539-1:1997, Information technology – Programming languages – Fortran – Part 1: Base language [22]. It is complemented with ISO/IEC TR 15580:1998, Information technology – Programming languages – Fortran – Floating-point exception handling [23].

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 operations marked " (\star) " are not part of the base standard, but included in the *Floating-point exception handling* Technical Report [23].

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 datatype 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, each parameterised with a *kind* parameter.

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

```
max_I(x,y)
                              MAX(x, y)
                              MIN(x, y)
min_I(x,y)
                              \mathtt{MAX}(x_1, x_2, \ldots, x_n)
max\_seq_I(xs)
                              \texttt{MAXVAL}(xs)
max\_seq_I(xs)
min\_seq_I(xs)
                              MIN(x_1, x_2, \ldots, x_n)
min\_seq_I(xs)
                              MINVAL(xs)
dim_I(x,y)
                              DIM(x, y)
power_I(x, y)
                              x ** y
shift 2_I(x,y)
                              SHIFT2(x, y)
shift 10_I(x,y)
                              SHIFT10(x, y)
sqrt_I(x)
                              ISQRT(x)
divides_I(x,y)
                              DIVIDES(x, y)
                              \texttt{MODULO(}x\texttt{,2)} \ == \ 0
even_I(x)
                              MODULO(x,2) \neq 0
odd_I(x)
quot_I(x,y)
                              QUOTIENT(x, y)
                              MODULO(x, y)
mod_I(x,y)
                              RATIO(x, y)
ratio_I(x,y)
                              RESIDUE(x, y)
residue_I(x,y)
                              GROUP(x, y)
group_I(x,y)
pad_I(x,y)
                              PAD(x, y)
                              GCD(x, y)
gcd_I(x,y)
lcm_I(x,y)
                              LCM(x, y)
                              GCDVAL ( & S)
gcd\_seq_I(xs)
                              LCMVAL(xs)
lcm\_seq_I(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\_OVER(x, y)
mul\_ov_I(x,y)
```

where x and y are expressions of type INTEGER(kind) and where xs is an expression of type array of INTEGER(kind).

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

$max_F(x,y)$	MAX(x, y)	*
$min_F(x,y)$	MIN(x, y)	*
$mmax_F(x,y)$	MMAX(x, y)	†
$mmin_F(x,y)$	MMIN(x, y)	†
$max_seq_F(xs)$	$\mathtt{MAX}(x_1, x_2, \ldots, x_n)$	*
$max_seq_F(xs)$	$ exttt{MAXVAL}(xs)$	*
$min_seq_F(xs)$	$\texttt{MIN}(x_1, x_2, \ldots, x_n)$	*
$min_seq_F(xs)$	$\mathtt{MINVAL}(xs)$	*

C.5 Fortran

```
mmax\_seq_F(xs)
                              \texttt{MMAX}(x_1, x_2, \ldots, x_n)
                              MMAXVAL(xs)
mmax\_seq_F(xs)
                              MMIN(x_1, x_2, \ldots, x_n)
mmin\_seq_F(xs)
                              MMINVAL(xs)
mmin\_seq_F(xs)
dim_F(x,y)
                              DIM(x, y)
floor_F(x)
                               IEEE_RINT(x) (if in round towards -\infty mode) (*)
floor\_rest_F(x)
                               x - IEEE_RINT(x) (if in round towards -\infty mode) (\star)
rounding_F(x)
                              IEEE_RINT(x) (if in round to nearest mode) (\star)
rounding\_rest_F(x)
                              x - IEEE_RINT(x) (if in round to nearest mode) (\star)
ceiling_F(x)
                              IEEE_RINT(x) (if in round towards +\infty mode) (*)
                              x - \text{IEEE\_RINT}(x) (if in round towards +\infty mode) (**)
ceiling\_rest_F(x)
residue_F(x,y)
                              IEEE_REM(x, y)
                              SQRT(x)
sqrt_F(x)
                              REC\_SQRT(x)
rec\_sqrt_F(x)
mul_{F \to F'}(x, y)
                              DPROD(x, y)
add\_lo_F(x,y)
                               ADD_LOW(x, y)
sub\_lo_F(x,y)
                              SUB_LOW(x, y)
                              MUL\_LOW(x, y)
mul\_lo_F(x,y)
div\_rest_F(x,y)
                              DIV_REST(x, y)
sqrt\_rest_F(x)
                              SQRT_REST(x)
```

where x and y are expressions of type REAL(kind), and where xs is an expression of type array of REAL(kind).

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

$max_error_hypot_F$	$ERR_HYPOTENUSE(x)$	†
	io io	
$max_error_exp_F$	$ERR_{\cdot}EXP(x)$	†
$max_error_power_F$	ERR_POWER(x)	†
	J.	
$big_angle_r_F$	$BIG_RADIAN_ANGLE(x)$	†
$max_error_rad_F$	$ERR_RAD(x)$	†
$max_error_sin_F$	$ERR_SIN(x)$	†
$max_error_tan_F$	$ERR_TAN(x)$	†
2N		
$min_angular_unit_F$	$MIN_ANGLE_UNIT(x)$	†
$big_angle_u_F$	$BIG_ANGLE(x)$	†
$max_2error_sinu_F(u)$	$ERR_SIN_CYCLE(u)$	†
$max_error_tanu_F(u)$	ERR_TAN_CYCLE(u)	†
$max_error_sinh_F$	$ERR_SINH(x)$	†
$max_error_tanh_F$	$ERR_TANH(x)$	†
$max_error_convert_F$	$ERR_CONVERT(x)$	†
$max_error_convert_{F'}$	ERR_CONVERT_TO_STRING	†
$max_error_convert_{D'}$	ERR_CONVERT_TO_STRING	†

where x and u are expressions of type REAL(kind). 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:

$hypot_F(x,y)$	HYPOT(x, y)	†
$power_{F,I}(b,z)$	b ** z	*
$exp_F(x)$	EXP(x)	^ +
$expm1_F(x)$	EXPM1(x)	+
$exp2_F(x)$	EXP2(x)	+ .
$exp2_F(x)$ $exp10_F(x)$	EXP10(x)	+ ~
= ', ',	EAF10(1)	900
$power_F(b, y)$	b ** y	7.7
$power1pm1_F(b,y)$	POWER1PM1 (b, y)	
$ln_F(x)$	LOG(x)	*
$ln1p_F(x)$	LOG1P(x)	†
$log \mathcal{Z}_F(x)$	LOG2(x)	†
$log 10_F(x)$	LOG10(x)	*
$logbase_F(b,x)$	LOGBASE(b, x)	†
$logbase1p1p_F(b,x)$	LOGBASE1P1P (b, x)	†
009000017177(0,0)	200211221111 (0, 10, 10, 10, 10, 10, 10, 10, 10, 10, 1	'
$rad_F(x)$	b ** z EXP(x) EXPM1(x) EXP2(x) EXP10(x) b ** y POWER1PM1(b, y) LOG(x) LOG1P(x) LOG2(x) LOG10(x) LOGBASE(b, x) LOGBASE1P1P(b, x) RAD(x) SIN(x) COS(x) TAN(x) COT(x) SEC(x)	†
$sin_F(x)$	SIN(x)	*
$cos_F(x)$	COS(x)	*
$tan_F(x)$	TAN(x)	*
$cot_F(x)$	COT.(a)	†
$sec_F(x)$	SEC(x)	† † †
$csc_F(x)$	CSC(x)	†
$arcsin_F(x)$ $arccos_F(x)$ $arctan_F(x)$ $arccot_F(x)$		
$arcsin_F(x)$	ASIN(x)	*
$arccos_F(x)$	ACOS(x)	*
$arctan_F(x)$	ATAN(x)	*
$arccot_F(x)$	ACOT(x)	†
ar ccore F(x)	ACOTC(x)	†
$arcsec_F(x)$	ASEC(x)	† † † †
$arccse_F(x)$	ACSC(x)	†
$arc_F(x,y)$	ATAN2(y, x)	*
$cycle_F(u,x)$	CYCLE(u, x)	†
$sinu_F(u,x)$	SINU(u, x)	†
$cosu_F(u,x)$	COSU(u, x)	†
$tanu_F(u,x)$	TANU(u, x)	†
$cotu_F(u,x)$	COTU(u, x)	†
$secu_F(u,x)$	SECU(u, x)	† † † †
1 () /	. •	'

C.5 Fortran

$cscu_F(u,x)$	CSCU(u, x)	†
$arcsinu_F(u,x)$	ASINU(u, x)	†
$arccosu_F(u,x)$	ACOSU(u, x)	
$arctanu_F(u,x)$		† +
- ()	ATANU(u, x)	<u> </u>
$arccotu_F(u,x)$	ACOTU(u, x)	† †
$arccotcu_F(u,x)$	ACOTCU(u, x)	Ţ
$arcsecu_F(u,x)$	ASECU(u, x)	† †
$arccscu_F(u,x)$	ACSCU(u, x)	Ţ
$arcu_F(u, x, y)$	ATAN2U(u, y, x)	†
$cycle_F(360,x)$	DEGREES(x)	† ^
$sinu_F(360,x)$	SIND(x)	†~00
$cosu_F(360,x)$	COSD(x)	مرابل
$tanu_F(360,x)$	TAND(x)	1
$cotu_F(360,x)$	COTD(x)	†
$secu_F(360, x)$	SECD(x)	†
$cscu_F(360,x)$	COSD(x) TAND(x) COTD(x) SECD(x) CSCD(x) ASIND(x) ACOSD(x) ATAND(x) ACOTD(x) ACOTCD(x) ASECD(x)	†
$arcsinu_F(360, x)$	ASIND(x)	†
$arccosu_F(360, x)$	ACOSD(x)	†
$arctanu_F(360,x)$	ATAND(x)	†
$arccotu_F(360,x)$	ACOTD(x)	†
$arccotcu_F(360,x)$	ACOTCD(x)	†
$arcsecu_F(360,x)$	ASECD(x)	† † † †
$arccscu_F(360,x)$	ACSCD(x)	†
$arcu_F(360, x, y)$	ATAN2D(y)	†
		'
$rad_to_cycle_F(x, w)$	RAD_TO_CYCLE(x , w)	†
$cycle_to_rad_F(u,x)$	$CYCLE_{TO_{RAD}(u, x)}$	†
$cycle_to_cycle_F(u, x, w)$	CYCLE_TO_CYCLE(u, x, w)	†
		'
$sinh_F(x)$	SINH(x)	*
$cosh_F(x)$	COSH(x)	*
$sinh_F(x)$ $cosh_F(x)$ $tanh_F(x)$ $coth_F(x)$ $sech_F(x)$	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)	†

where b, x, y, u, and w are expressions of type REAL(kind), and z is an expression of type INTEGER(kindi).

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

```
convert_{I \to I'}(x)
                               INT(x, kindi2)
                       lbl\_a
                               FORMAT (Bn)
                                                                                 ★(binary)
convert_{I''\to I}(f)
                               READ (UNIT=#f,FMT=lbl_a) r
                               WRITE (UNIT=#h, FMT=lbl_a) x
convert_{I \to I''}(x)
                       lbl\_b
                               FORMAT (On)
                                                                                 \star(octal)
convert_{I''\to I}(f)
                               READ (UNIT=#f,FMT=lbl_b) r
                               WRITE (UNIT=#h, FMT=lbl_b) x
convert_{I \to I''}(x)
                       lbl\_c
                               FORMAT (In)
                                                                                 \star (decimal)
convert_{I''\to I}(f)
                               READ (UNIT=#f,FMT=lbl_c) r
convert_{I \to I''}(x)
                               WRITE (UNIT=#h, FMT=lbl_c) x
                                                                                 ★(hexadecimal)
                       lbl\_d
                               FORMAT (Zn)
convert_{I''\to I}(f)
                               READ (UNIT=#f,FMT=lbl_d) r
convert_{I \to I''}(x)
                               WRITE (UNIT=#h, FMT=lbl_d) x
                               FLOOR(y, kindi?)
floor_{F \to I}(y)
                               ROUND(y, kindi?)
rounding_{F \to I}(y)
                               CEILING(y, kindi?)
ceiling_{F \to I}(y)
                               REAL(x, kind) or sometimes DBLE(x)
convert_{I \to F}(x)
convert_{F \to F'}(y)
                               REAL(y, kind2) or sometimes DBLE(y)
                       lbl\_e
                               FORMAT (Fw.d)
                               FORMAT (\mathbf{D}w.d)
                       lbl_{-}f
                               FORMATO (Ew.d)
                       lbl\_g
                       lbl\_h
                               FORMAT (Ew.dEe)
                              FORMAT (ENw.d)
                       lbl\_i
                       lbl_{-}j
                             FORMAT (ENw.dEe)
                       lblk
                               FORMAT (ESw.d)
                       lbl\_l
                               FORMAT (ESw.dEe)
                               READ (UNIT=#f, FMT=lbl_x) t
                               WRITE (UNIT=#h, FMT=lbl_x) y
                               READ (UNIT=#f, FMT=lbl_x) t
```

where x is an expression of type INTEGER(kindi), y is an expression of type REAL(kind), f is an input file with unit number #f, and h is an output file with unit number #h. w, d, and e are literal digit (0-9) sequences, giving total, decimals, and exponent widths. lbl_x is one of lbl_e to lbl_z ; all of the lbl_x are labels for formats.

Fortran provides base 10 non-negative numerals for all of its integer and floating point types. Numerals for floating point types must have a '.' in them. The details are not repeated in this example binding, see ISO/IEC 1539-1:1997, clause 4.3.1.1 Integer type, and clause 4.3.1.2 Real type.

Fortran does not specify numerals for infinities and NaNs. Suggestion:

 $C.5 \ Fortran$

$+\infty$	INFINITY	†
${ m qNaN}$	NAN	†
sNaN	NANSIGNALLING	†

as well as string formats for reading and writing these values as character strings.

Fortran implementations can provide recording of indicators for floating point arithmetic notifications, the LIA-2 preferred method. See ISO/IEC TR 15580:1998, *Information technology – Programming languages – Fortran – Floating-point exception handling* [23]. **absolute_precision_underflow** notifications are however ignored.

C.6 Haskell

The programming language Haskell is defined by Report on the programming language Haskell 98 [65], together with Standard libraries for the Haskell 98 programming language [66].

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 Haskell datatype Bool corresponds to the LIA datatype Boolean.

Every implementation of Haskell has at least two integer datatypes, Integer, which is unbounded, 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.

Haskell has a type class system that allows for overloading, and allowing static type checking of dynamic overloading. But in contrast to object oriented programming languages, type classes are not types. E.g. + has the type (Num a) => a -> a, where Num is a type class and a is a type variable.

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

```
max_I(x,y)
                              \max x y
min_I(x,y)
                                               x 'min' y
                               \min x y
                                         or
max\_seq_I(xs)
                               \verb"maximum" xs
min\_seq_I(xs)
                               \min \max \ xs
dim_I(x,y)
                               \dim x y or x 'dim' y
power_I(x, y)
                               x \hat{y} or (\hat{y})
shift2_{I}(x,y)
                               shift2 x y or x 'shift2' y
shift 10_I(x,y)
                               shift10 x y or x 'shift10' y
sqrt_I(x)
                               isqrt x
divides_I(x,y)
                               divides x \ y or x 'divides' y
even_I(x)
                               even x
odd_I(x)
                               odd x
quot_I(x,y)
                               \operatorname{div} x y \quad \operatorname{or}
                                              x 'div' y
mod_I(x,y)
                               mod x y or x 'mod'
```

```
ratio_I(x,y)
                             ratio x \ y or x 'ratio' y
residue_I(x,y)
                             residue x \ y or x 'residue' y
group_I(x,y)
                             grp x y or
pad_I(x,y)
                             pad x y or
                                          x 'pad' y
gcd_I(x,y)
                             \gcd x y \text{ or }
                                            x 'gcd' y
lcm_I(x,y)
                             lcm x y or
                                            x 'lcm' y
gcd\_seq_I(xs)
                             gcd\_seq xs
lcm\_seq_I(xs)
                             lcm\_seq xs
add\_wrap_I(x,y)
                             x +:
add\_ov_I(x,y)
                             x + :+ y
sub\_wrap_I(x,y)
sub\_ov_I(x,y)
                             x -:+ y
mul\_wrap_I(x,y)
mul\_ov_I(x,y)
                             x *: + y
```

where x and y are expressions of type INT and where xs is an expression of type [INT].

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

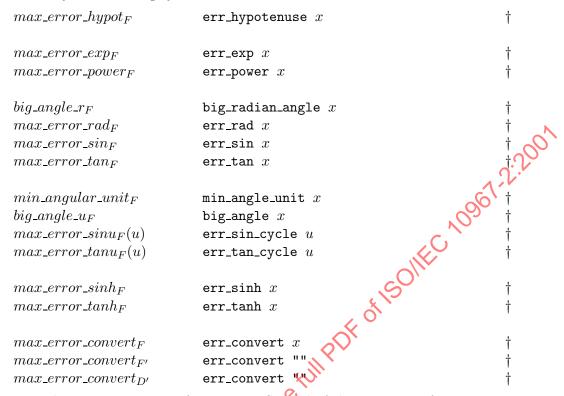
```
max_F(x,y)
min_F(x,y)
                               \min x y or
mmax_F(x,y)
                               mmax x y or
mmin_F(x,y)
                               mmin x y or
                                                    'mmin' y
max\_seq_F(xs)
                               {\tt maximum}\ xs
min\_seq_F(xs)
                               {\tt minimum}\ xs
mmax\_seq_F(xs)
                               mmaximum xs
mmin\_seq_F(xs)
                               mminimum xs
                                \operatorname{\mathtt{dim}} x \ y \quad \text{or} \quad x \text{ 'dim'} \ y
dim_F(x,y)
floor_F(x)
                               fromInteger (floor x)
floor\_rest_F(x)
                               x - fromInteger (floor x)
rounding_F(x)
                               fromInteger (round x)
rounding\_rest_F(x)
                               x - fromInteger (round x)
ceiling_F(x)
                               fromInteger (ceiling x)
ceiling_rest_(x
                               x - fromInteger (ceiling x)
residue_F(x, y)
                               residue x \ y or x 'residue' y
sqrt_F(x)
                               \operatorname{sqrt} x
                               {\tt rec\_sqrt}\ x
rec_sqrt_F(x)
mul_{F\to F'}(x,y)
                               prod x y
add\_lo_F(x,y)
                               x +: - y
sub\_lo_F(x,y)
mul\_lo_F(x,y)
                               x *: - y
div\_rest_F(x,y)
                               x /: * y
sqrt\_rest_F(x)
                               sqrt\_rest x
```

where x and y are expressions of type FLT, and where xs is an expression of type [FLT].

C.6 Haskell

The binding for the floor, round, and ceiling operations here take advantage of the unbounded Integer type in Haskell, and that Integer is the default integer type.

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



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:

```
hypotenuse x y
hypot_F(x,y)
power_{F,I}(b,z)
exp_F(x)
                              expM1 x
exp2_F(x)
                             exp2 x
exp10_F(x)
                             exp10 x
power_{F}(b, y)
                             b ** y \text{ or } (**) b y
power1pm1_F(b, y)
                             power1PM1 b y or b 'power1PM1' y
ln_F(x)
                             \log x
ln1p_F(x)
                             log1P x
log 2_F(x)
                             log2 x
log 10_F(x)
                             log10 x
                             logBase b x or b 'logBase' x
logbase_F(b,x)
logbase1p1p_F(b, x)
                             logBase1P1P \ b \ x
```

$rad_F(x)$	radians x	†
$axis_rad_F(x)$	$axis_radians x$	†
		'
$sin_F(x)$	$\sin x$	*
$cos_F(x)$	$\cos x$	*
$tan_F(x)$		
$cot_F(x)$	cot x	†
$sec_F(x)$	sec x	+
$csc_F(x)$	csc x	+
$cossin_F(x)$	cosSin x	+
$cossin_F(x)$		' ,
$arcsin_F(x)$	tan x cot x sec x csc x cosSin x asin x acos x atan x acot x acotc x asec x atan2 y x cycle u x axis_cycle u x tanU u x cosU u x secU u x secU u x cosSinU u x	$Q_{O\star}$
$arccos_F(x)$	$a\cos x$	- D
$arctan_F(x)$	atan x	. L.
$arccot_F(x)$	acot x	+
$arccotc_F(x)$	acotc x	† †
$arcsec_F(x)$	asec x	+
$arccsc_F(x)$	acsc x	+
$arc_F(x,y)$	atan2 y x	1
$arc_F(x,y)$	ataliz y x	^
$cycle_F(u,x)$	cycle u x	+
$axis_cycle_F(u,x)$	axis_cycle u x	†
		'
$sinu_F(u,x)$	sinU u x	†
$cosu_F(u,x)$	cosU u x	†
$tanu_F(u,x)$	tanU u	†
$cotu_F(u,x)$	cotU w x	†
$secu_F(u,x)$	seculu x	÷
$cscu_F(u,x)$	cscU u x	÷
$cossinu_F(u,x)$	cosSinU u x	† † † † †
aic aic		'
$arcsinu_F(u,x)$	$\verb"asinU" u x$	†
$arccosu_F(u,x)$	acosU u x	†
$arctanu_F(u,x)$	$\mathtt{atanU} \ u \ x$	†
$arccotu_F(u, \mathbf{x})$	${\tt acotU}\ u\ x$	†
$arccotcu_{I\!\!F}(u,x)$	${\tt acotcU}\ u\ x$	†
$arcsecu_{F}(u,x)^{'}$	asecU u x	†
$arcescu_F(u,x)$	$\operatorname{acscU} \ u \ x$	†
$areu_F(u,x,y)$	atan2U u y x	†
		'
$rad_to_cycle_F(x, w)$	${\tt rad_to_cycle}\ w\ x$	†
$cycle_to_rad_F(u,x)$	${\tt cycle_to_rad}\ u\ x$	†
$cycle_to_cycle_F(u, x, w)$	cycle_to_cycle $u \ x \ w$	†
	•	
$sinh_F(x)$	$\sinh x$	*
$cosh_F(x)$	$\cosh x$	*
$tanh_F(x)$	tanh x	*
- 🗸 /		

C.6 Haskell

$coth_F(x)$	$\coth x$	†
$sech_F(x)$	$\operatorname{sech} x$	†
$csch_F(x)$	$\operatorname{csch} x$	†
$arcsinh_F(x)$	$\mathtt{asinh}\ x$	*
$arccosh_F(x)$	$\operatorname{acosh}\ x$	*
$arctanh_F(x)$	$\operatorname{atanh}\ x$	*
$arccoth_F(x)$	$\mathbf{a}\mathbf{coth}\ x$	†
$arcsech_F(x)$	$\operatorname{asech} x$	†
$arccsch_F(x)$	$\operatorname{acsch} x$	†

where b, x, y, u, and w are expressions of type FLT, and z is an expression of type INT.

Arithmetic value conversions in Haskell are always explicit. They are done with the overloaded from Integral and from Fractional operations.

$convert_{I \to I'}(x)$	fromIntegral x	*
$convert_{I''\to I}(x)$	read s	*
$convert_{I \to I''}(x)$	show x	*
$floor_{F \to I}(y)$	floor y	*
$rounding_{F \to I}(y)$	round y	*
$ceiling_{F \to I}(y)$	ceiling y	*
$convert_{I \to F}(x)$	fromIntegral x	*
$convert_{F \to F'}(y)$	fromFractional	*
$convert_{F''\to F}(s)$	read s	*
$convert_{F \to F''}(y)$	show y	*
$convert_{D' \to F}(s)$	readie	*
$convert_{F \to D'}(y)$	show y	*

where x is an expression of type INT, y is an expression of type FLT. show does not allow for format control.

Haskell provides non-negative numerals for all its integer and floating point types in base 10. There is no differentiation between the numerals for different floating point datatypes, nor between numerals for different integer datatypes, and integer numerals can be used for floating point values. Integer numerals stand for a value in Integer (the unbounded integer type) and an implicit fromInteger operation is applied to it. Fractional numerals stand for a value in Rational (the unbounded type of rational numbers) and an implicit fromRational operation is applied to it.

Haskell does not specify any numerals for infinities and NaNs. Suggestion:

$+\infty$	infinity	†
qNaN	${\tt quietNaN}$	†
sNaN	${ t sigalling NaN}$	†

as well as string formats for reading and writing these values as character strings.

Haskell has the notion of error, which results in a change of 'control flow', which cannot be returned from, nor caught. An error results in the termination of the program. (There are suggestions to improve this.) **infinitary** for integer types and **invalid** (in general) are considered to be error. No notification results for **underflow**, and the continuation value (specified by LIA-2) is used directly, since recording of indicators is not available and error is inappropriate for **underflow**. The handling of integer **overflow** is implementation dependent. The handling of floating point **overflow** and **infinitary** should be to return a suitable infinity (specified by LIA-2), if possible, without any notification, since recording of indicators is not available.

C.7 Java

The programming language Java is defined by *The Java Language Specification* [64], plus a number of class libraries (exactly which vary depending on the Java 'edition' and version).

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 LIA-2 operations that are provided in Java 2 (marked "*" below) are in the final class java.lang.Math.

The Java datatype boolean corresponds to the LIA datatype Boolean.

Every implementation of Java has the integer datatypes int and long. The notation *INT* will be used to refer to either one of them.

Java has two floating point datatypes, float and double, which must conform to IEC 60559. The notation FLT will be used to refer to either one of them.

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

```
min_I(x,y)
                             \min(x, y)
                             \max(x, y)
max_I(x,y)
min\_seq_I(xs)
                             min_arr(xs)
max\_seq_I(xs)
                             \mathtt{max\_arr}(xs)
dim_I(x,y)
                              dim(x, y)
power_I(x, y)
                              power(x, y)
shift2_I(x,y)
                              shift2(x, y)
shift 10_I(x,y)
                              shift10(x, y)
sqrt_I(x)
                              sqrt(x)
divides_I(x,y)
                              divides(x, y)
divides_I(x,y)
                              x != 0 \&\& y \% x == 0
                              x \% 2 == 0
even_I(x)
                              x \% 2 != 0
odd_I(x)
                              quot(x, y)
quot_I(x,y)
mod_I(x,y)
                              mod(x, y)
ratio_I(x, y)
                              ratio(x, y)
residue_I(x,y)
                              residue(x, y)
```

C.7 Java

```
group_I(x,y)
                             group(x, y)
pad_I(x,y)
                             pad(x, y)
                                                                            †
gcd_I(x,y)
                             gcd(x, y)
lcm_I(x,y)
                             lcm(x, y)
                             gcd_arr(xs)
gcd\_seq_I(xs)
lcm\_seq_I(xs)
                             lcm_arr(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\_over(x, y)
sub\_ov_I(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 INT[].

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

```
FULL POF OF ISOILE.
min_F(x,y)
                             min(x, y)
max_F(x,y)
                             \max(x, y)
mmax_F(x,y)
                             mmax(x, y)
                             mmin(x, y)
mmin_F(x,y)
                             min_arr(xs)
min\_seq_F(xs)
                             max_arr(xs)
max\_seq_F(xs)
mmax\_seq_F(xs)
                             mmax(xs)
                             mmin(xs)
mmin\_seq_F(xs)
                            dim(x, y)
dim_F(x,y)
                             floor(x)
floor_F(x)
                                                                           ★(only for double)
                             x floor(x)
floor\_rest_F(x)
                                                                           ★(only for double)
rounding_F(x)
                             rint(x)
                                                                           ★(only for double)
rounding\_rest_F(x)
                             x - rint(x)
                                                                           ★(only for double)
ceiling_F(x)
                             ceil(x)
                                                                           ★(only for double)
ceiling\_rest_F(x)
                             x - ceil(x)
                                                                           ★(only for double)
residue_F(x,y)
                             IEEEremainder(x, y)
                                                                           ⋆ (only for double)
sqrt_F(x)
                             sqrt(x)
rec\_sqrt_F(x)
                             rec\_sqrt(x)
                             dprod(x, y)
mul_{F \to F'}(x, y)
add lo_F(x,y)
                             add_low(x, y)
sub\_lo_F(x,y)
                             sub_low(x, y)
mul\_lo_F(x,y)
                             mul_low(x, y)
div\_rest_F(x,y)
                             div_rest(x, y)
                             sqrt_rest(x)
sqrt\_rest_F(x)
```

where x and y are expressions of type FLT, and where xs is an expression of type FLT[].

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

$max_error_hypot_F$	err_hypotenuse(x)	†
$max_error_exp_F$ $max_error_power_F(b, x)$	<pre>err_exp(x) err_power(x)</pre>	† †
$big_angle_r_F$ $max_error_rad_F$ $max_error_sin_F$ $max_error_tan_F$	$\begin{array}{l} \operatorname{big_radian_angle}(x) \\ \operatorname{err_rad}(x) \\ \operatorname{err_sin}(x) \\ \operatorname{err_tan}(x) \end{array}$	† † †
$min_angular_unit_F$ $big_angle_u_F$ $max_error_sinu_F(u)$ $max_error_tanu_F(u)$	<pre>smallest_angular_unit(x) big_angle(x) err_sin_cycle(u) err_tan_cycle(u)</pre>	† †001
$max_error_sinh_F$ $max_error_tanh_F$	err_sinh(x) err_tanh(x)	† †
$max_error_convert_F$ $max_error_convert_{D'}$ $max_error_convert_{D'}$	err_convert(x) err_convert_to_string err_convert_to_string	† † †

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. These are defined only for double not for float.

$hypot_F(x)$	(x,y)	hypotenuse (x, y)	†
$power_{F,I}$	(b,z)	poweri(b, z)	†
$exp_F(x)$	N	exp(x)	*
$expm1_F($	x)	expm1(x)	†
$exp2_F(x)$.0	exp2(x)	†
$exp10_F(x)$	c)	exp10(x)	†
$power_F(l)$	_ (4	power(b, y)	†
$pow_F(b, y)$		pow(b, y)	\star Not LIA-2!
power1pr	$mI_F(b,y)$	power1pm1(b, y)	†
),		
$ln_F(x)$		log(x)	*
$ln1p_F(x)$)	log1p(x)	†
$log 2_F(x)$		log2(x)	†
$log 10_F(x)$)	log10(x)	†
$logbase_{F}$	(b,x)	log(b, x)	†
logbase1p	$p1p_F(b,x)$	log1p1p(b, x)	†
$rad_F(x)$		radian(x)	†
$axis_rad$	$f_F(x)$	<pre>axis_rad(x)</pre>	†

C.7 Java 139

$sin_F(x)$	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)	†
$arccotc_F(x)$	acotc(x)	† ~
$arcsec_F(x)$	asec(x)	100
$arccsc_F(x)$	acsc(x)	d.h
$arc_F(x,y)$	atan2(y, x)	*
- ,	90	
$cycle_F(u,x)$	cycle(u, x)	†
$axis_cycle_F(u,x)$	$axis_cycle(u, x)$	†
. (J
$sinu_F(u,x)$	sinu(u, x)	Ţ
$cosu_F(u,x)$	$\cos(u, x)$	1
$tanu_F(u,x)$	tanu(u, x)	 +
$cotu_F(u,x)$	cotu(u, x)	 +
$secu_F(u,x)$	secu(u, x)	 +
$cscu_F(u,x)$	cscu(u, x)	1
$arcsinu_F(u,x)$	acos(x) atan(x) acot(x) acot(x) asec(x) acsc(x) atan2(y, x) cycle(u, x) axis_cycle(u, x) sinu(u, x) cosu(u, x) tanu(u, x) cotu(u, x) secu(u, x) asinu(u, x) asinu(u, x) asinu(u, x) atanu(u, x) atanu(u, x) atanu(u, x)	†
$arccosu_F(u,x)$	$acosu(u \wedge x)$	†
$arctanu_F(u,x)$	atanu(u, x)	†
$arccotu_F(u,x)$	acotu(u, x)	†
$arccotcu_F(u,x)$	acotcu(u, x)	† † † †
$arcsecu_F(u,x)$	$\mathbf{Q}_{asecu}(u, x)$	†
$arccscu_F(u,x)$	acscu(u, x)	†
$arcu_F(u, x, y)$	atan2u(u, y, x)	†
CO.		
$rad_to_cycle_F(x, 360)$	toDegrees(x)	*
$cycle_to_rad_F(360, x)$	toRadians(x)	*
$rad_to_cycle_F(x, w)$	radian_to_cycle(x, w)	†
$cycle to rad_F(u, x)$	cycle_to_radian(u, x)	†
$cycle_to_cycle_F(u, x, w)$	$cycle_{to}_{cycle}(u, x, w)$	†
$sinh_F(x)$	sinh(x)	†
$cosh_F(x)$	$\cosh(x)$	†
$tanh_F(x)$	tanh(x)	† † † †
$coth_F(x)$	$\coth(x)$	+
$sech_F(x)$	$\operatorname{sech}(x)$	†
$csch_F(x)$	$\operatorname{csch}(x)$	†
	\w/	'