

IEEE Standard for Interval Arithmetic

IEEE Computer Society

Sponsored by the
IEEE Microprocessor Standards Committee

IEEE Standard for Interval Arithmetic

Sponsor

Microprocessor Standards Committee
of the
IEEE Computer Society

Approved 11 June 2015

IEEE-SA Standards Board

Abstract: This standard specifies basic interval arithmetic (IA) operations selecting and following one of the commonly used mathematical interval models. This standard supports the IEEE 754TM floating-point formats of practical use in interval computations. Exception conditions are defined, and standard handling of these conditions is specified. Consistency with the interval model is tempered with practical considerations based on input from representatives of vendors, developers and maintainers of existing systems.

The standard provides a layer between the hardware and the programming language levels. It does not mandate that any operations be implemented in hardware. It does not define any realization of the basic operations as functions in a programming language.

Keywords: arithmetic, computing, decoration, enclosure, hull, IEEE 1788TM, interval, operation, verified

The Institute of Electrical and Electronics Engineers, Inc.
3 Park Avenue, New York, NY 10016-5997, USA

Copyright © 2015 by The Institute of Electrical and Electronics Engineers, Inc.
All rights reserved. Published 30 June 2015. Printed in the United States of America.

IEEE is a registered trademark in the U.S. Patent & Trademark Office, owned by The Institute of Electrical and Electronics Engineers, Incorporated.

PDF: ISBN 978-0-7381-9720-3 STD20228
Print: ISBN 978-0-7381-9721-0 STDPD20228

IEEE prohibits discrimination, harassment, and bullying.

For more information, visit <http://www.ieee.org/web/aboutus/whatis/policies/p9-26.html>.

No part of this publication may be reproduced in any form, in an electronic retrieval system or otherwise, without the prior written permission of the publisher.

Important Notices and Disclaimers Concerning IEEE Standards Documents

IEEE documents are made available for use subject to important notices and legal disclaimers. These notices and disclaimers, or a reference to this page, appear in all standards and may be found under the heading “Important Notice” or “Important Notices and Disclaimers Concerning IEEE Standards Documents.”

Notice and Disclaimer of Liability Concerning the Use of IEEE Standards Documents

IEEE Standards documents (standards, recommended practices, and guides), both full-use and trial-use, are developed within IEEE Societies and the Standards Coordinating Committees of the IEEE Standards Association (“IEEE-SA”) Standards Board. IEEE (“the Institute”) develops its standards through a consensus development process, approved by the American National Standards Institute (“ANSI”), which brings together volunteers representing varied viewpoints and interests to achieve the final product. Volunteers are not necessarily members of the Institute and participate without compensation from IEEE. While IEEE administers the process and establishes rules to promote fairness in the consensus development process, IEEE does not independently evaluate, test, or verify the accuracy of any of the information or the soundness of any judgments contained in its standards.

IEEE does not warrant or represent the accuracy or content of the material contained in its standards, and expressly disclaims all warranties (express, implied and statutory) not included in this or any other document relating to the standard, including, but not limited to, the warranties of: merchantability; fitness for a particular purpose; non-infringement; and quality, accuracy, effectiveness, currency, or completeness of material. In addition, IEEE disclaims any and all conditions relating to: results; and workmanlike effort. IEEE standards documents are supplied “AS IS” and “WITH ALL FAULTS.”

Use of an IEEE standard is wholly voluntary. The existence of an IEEE standard does not imply that there are no other ways to produce, test, measure, purchase, market, or provide other goods and services related to the scope of the IEEE standard. Furthermore, the viewpoint expressed at the time a standard is approved and issued is subject to change brought about through developments in the state of the art and comments received from users of the standard.

In publishing and making its standards available, IEEE is not suggesting or rendering professional or other services for, or on behalf of, any person or entity nor is IEEE undertaking to perform any duty owed by any other person or entity to another. Any person utilizing any IEEE Standards document, should rely upon his or her own independent judgment in the exercise of reasonable care in any given circumstances or, as appropriate, seek the advice of a competent professional in determining the appropriateness of a given IEEE standard.

IN NO EVENT SHALL IEEE BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO: PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE PUBLICATION, USE OF, OR RELIANCE UPON ANY STANDARD, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE AND REGARDLESS OF WHETHER SUCH DAMAGE WAS FORESEEABLE.

Translations

The IEEE consensus development process involves the review of documents in English only. In the event that an IEEE standard is translated, only the English version published by IEEE should be considered the approved IEEE standard.

Official statements

A statement, written or oral, that is not processed in accordance with the IEEE-SA Standards Board Operations Manual shall not be considered or inferred to be the official position of IEEE or any of its committees and shall not be considered to be, or be relied upon as, a formal position of IEEE. At lectures, symposia, seminars, or educational courses, an individual presenting information on IEEE standards shall make it clear

that his or her views should be considered the personal views of that individual rather than the formal position of IEEE.

Comments on standards

Comments for revision of IEEE Standards documents are welcome from any interested party, regardless of membership affiliation with IEEE. However, IEEE does not provide consulting information or advice pertaining to IEEE Standards documents. Suggestions for changes in documents should be in the form of a proposed change of text, together with appropriate supporting comments. Since IEEE standards represent a consensus of concerned interests, it is important that any responses to comments and questions also receive the concurrence of a balance of interests. For this reason, IEEE and the members of its societies and Standards Coordinating Committees are not able to provide an instant response to comments or questions except in those cases where the matter has previously been addressed. For the same reason, IEEE does not respond to interpretation requests. Any person who would like to participate in revisions to an IEEE standard is welcome to join the relevant IEEE working group.

Comments on standards should be submitted to the following address:

Secretary, IEEE-SA Standards Board
445 Hoes Lane
Piscataway, NJ 08854 USA

Laws and regulations

Users of IEEE Standards documents should consult all applicable laws and regulations. Compliance with the provisions of any IEEE Standards document does not imply compliance to any applicable regulatory requirements. Implementers of the standard are responsible for observing or referring to the applicable regulatory requirements. IEEE does not, by the publication of its standards, intend to urge action that is not in compliance with applicable laws, and these documents may not be construed as doing so.

Copyrights

IEEE draft and approved standards are copyrighted by IEEE under U.S. and international copyright laws. They are made available by IEEE and are adopted for a wide variety of both public and private uses. These include both use, by reference, in laws and regulations, and use in private self-regulation, standardization, and the promotion of engineering practices and methods. By making these documents available for use and adoption by public authorities and private users, IEEE does not waive any rights in copyright to the documents.

Photocopies

Subject to payment of the appropriate fee, IEEE will grant users a limited, non-exclusive license to photocopy portions of any individual standard for company or organizational internal use or individual, non-commercial use only. To arrange for payment of licensing fees, please contact Copyright Clearance Center, Customer Service, 222 Rosewood Drive, Danvers, MA 01923 USA; +1 978 750 8400. Permission to photocopy portions of any individual standard for educational classroom use can also be obtained through the Copyright Clearance Center.

Updating of IEEE Standards documents

Users of IEEE Standards documents should be aware that these documents may be superseded at any time by the issuance of new editions or may be amended from time to time through the issuance of amendments, corrigenda, or errata. An official IEEE document at any point in time consists of the current edition of the document together with any amendments, corrigenda, or errata then in effect.

Every IEEE standard is subjected to review at least every ten years. When a document is more than ten years old and has not undergone a revision process, it is reasonable to conclude that its contents, although still of some value, do not wholly reflect the present state of the art. Users are cautioned to check to determine that they have the latest edition of any IEEE standard.

In order to determine whether a given document is the current edition and whether it has been amended through the issuance of amendments, corrigenda, or errata, visit the IEEE-SA Website at <http://ieeexplore.ieee.org/xpl/standards.jsp> or contact IEEE at the address listed previously. For more information about the IEEE SA or IEEE's standards development process, visit the IEEE-SA Website at <http://standards.ieee.org>.

Errata

Errata, if any, for all IEEE standards can be accessed on the IEEE-SA Website at the following URL: <http://standards.ieee.org/findstds/errata/index.html>. Users are encouraged to check this URL for errata periodically.

Patents

Attention is called to the possibility that implementation of this standard may require use of subject matter covered by patent rights. By publication of this standard, no position is taken by the IEEE with respect to the existence or validity of any patent rights in connection therewith. If a patent holder or patent applicant has filed a statement of assurance via an Accepted Letter of Assurance, then the statement is listed on the IEEE-SA Website at <http://standards.ieee.org/about/sasb/patcom/patents.html>. Letters of Assurance may indicate whether the Submitter is willing or unwilling to grant licenses under patent rights without compensation or under reasonable rates, with reasonable terms and conditions that are demonstrably free of any unfair discrimination to applicants desiring to obtain such licenses.

Essential Patent Claims may exist for which a Letter of Assurance has not been received. The IEEE is not responsible for identifying Essential Patent Claims for which a license may be required, for conducting inquiries into the legal validity or scope of Patent Claims, or determining whether any licensing terms or conditions provided in connection with submission of a Letter of Assurance, if any, or in any licensing agreements are reasonable or non-discriminatory. Users of this standard are expressly advised that determination of the validity of any patent rights, and the risk of infringement of such rights, is entirely their own responsibility. Further information may be obtained from the IEEE Standards Association.

Participants

At the time this IEEE standard was completed, the Interval Standard Working Group had the following membership:

Nathalie Revol, *Chair*
R. Baker Kearfott, *Vice Chair and Acting Chair*
William Edmonson, *Secretary*
J. Wolff von Gutenberg, *Web Master*
Guillaume Melquiond, *Archivist*
George Corliss, *Voting Tabulator*
John Pryce, *Senior Technical Editor*
Christian Keil, *Deputy Technical Editor*
Michel Hack, Vincent Lefèvre, Ian McIntosh, Dmitry Nadezhin,
Ned Nedialkov and J. Wolff von Gutenberg, *Assistant Technical Editors*

Alexandru Amaricai	Malgorzata Jankowska	Evgenija Popova
Ayman Bakr	Michel Kieffer	Tarek Raissi
Ahmed Belhani	Walter Krämer	Nacim Ramdani
Gerd Bohlender	Vladik Kreinovich	Andreas Rauh
Gilles Chabert	Ulrich Kulisch	Gaby Dos Reis
Rudnei Dias da Cunha	Dorina Lanza	Michael Schulte
Hend Dawood	David Lester	Kyarash Shahriari
Bo Einarsson	Dominique Lohez	Stefan Siegel
Alan Eliassen	Wolfram Luther	Iwona Skalna
Hossam A. H. Fahmy	Amin Maher	Mark Stadtherr
Richard Fateman	Svetoslav Markov	James Stine
Scott Ferson	Günter Mayer	Pipop Thienprapasith
Haitham Gad	Jean-Pierre Merlet	Warwick Tucker
Kathy Gerber	Jean-Michel Muller	Alfredo Vaccaro
Alexandre Goldsztejn	Humberto Munoz	Maarten van Emden
Frederic Goualard	Jose Antonio Munoz	Erik-Jan van Kampen
Michael Groszkiewicz	Kaori Nagatou	Van Snyder
Mohamed Guerfel	Mitsuhiro Nakao	Josep Vehi
Robert Hanek	Markus Neher	Julio Villalba-Moreno
Behnam Hashemi	Marco Nehmeier	G. William Walster
Nathan Hayes	Diep Nguyen	Yan Wang
Oliver Heimlich	Michael Nooner	Lee Winter
Timothy Hickey	Shinichi Oishi	Pierre-Alain Yvars
Werner Hofschuster	Sylvain Pion	Sergei Zhilin
Chenyi Hu	Antony Popov	Mohamed Zidan
Trevor Jackson, III		Dan Zuras

The working group wishes to record with regret the loss of two members. Antony Popov of Sofia University, Bulgaria, died suddenly in 2012 at the young age of 49. Walter Krämer, of the Bergische Universität Wuppertal, Germany, died in October 2014 at the age of 62. Despite illness, Walter had remained an active participant until June 2014.

The following members of the individual balloting committee voted on this guide. Balloters may have voted for approval, disapproval, or abstention.

Bakul Banerjee	Oliver Heimlich	JeanMichel Muller
Juan Carreon	Werner Hoelzl	Dmitry Nadezhin
Keith Chow	Chenyi Hu	Marco Nehmeier
George Corliss	Piotr Karocki	John Pryce
Hossam Fahmy	Ralph Kearfott	Nathalie Revol
Andrew Fieldsend	Vladik Kreinovich	Eugene Stoudenmire
Alexander Gelman	Vincent Lefevre	Gerald Stueve
Frederic Goualard	Vincent Lipsio	J. Wolff Von Gudenberg
Randall Groves	William Lumpkins	Forrest Wright
Michel Hack	Guillaume Melquiond	Oren Yuen
Peter Harrod	James Moore	

When the IEEE-SA Standards Board approved this standard on 11 June 2015, it had the following membership:

John Kulick, *Chair*
Jon Walter Rosdahl, *Vice Chair*
Richard H. Hulett, *Past Chair*
Konstantinos Karachalios, *Secretary*

Masayuki Ariyoshi	Joseph L. Koepfinger*	Stephen J. Shellhammer
Ted Burse	David J. Law	Adrian P. Stephens
Stephen Dukes	Hung Ling	Yatin Trivedi
Jean-Philippe Faure	Andrew Myles	Phillip Winston
J. Travis Griffith	T. W. Olsen	Don Wright
Gary Hoffman	Glenn Parsons	Yu Yuan
Michael Janezic	Ronald C. Petersen	Daidi Zhong
	Annette D. Reilly	

*Member Emeritus

Don Messina
IEEE-SA Content Production and Management
Jonathan Goldberg
IEEE-SA Operational Program Management

Introduction

This introduction is not part of IEEE Std 1788-2015, IEEE Standard for Interval Arithmetic.

This introduction explains some of the alternative interpretations, and sometimes competing objectives, that influenced the design of this standard. Implementers should study it for a fuller understanding of the design choices made in this standard among these interpretations and objectives. For more information on interval computations, including history, applications and software, see^a e.g. [B1, B14] and the references therein, and also the interval computations web site [B5].

Mathematical context

Interval computation is a collaboration between human programmer and machine infrastructure which, correctly done, produces mathematically proven numerical results about continuous problems—for instance, rigorous bounds on the global minimum of a function or the solution of a differential equation. It is part of the discipline of “constructive real analysis.” In the long term, the results of such computations might become sufficiently trusted to be accepted as contributing to legal decisions. The machine infrastructure acts as a body of theorems on which the correctness of an interval algorithm relies, so it must be made as reliable as is practical. In its logical chain are many links—hardware, underlying floating-point system, etc.—over which this standard has no control. The standard aims to strengthen one specific link, by defining interval objects and operations that are theoretically well-founded and practical to implement.

This document uses the standard notation $[a, b]$ for “the interval between numbers a and b ,” with various detailed meanings depending on the underlying theory. The “classical” interval arithmetic (IA) of R.A. Moore [B8] uses only bounded, closed, nonempty intervals in the real numbers \mathbb{R} —that is, $[a, b] = \{x \in \mathbb{R} \mid a \leq x \leq b\}$ where $a, b \in \mathbb{R}$ with $a \leq b$. So, for instance, division by an interval containing 0 is not defined in it. It was agreed early on that this standard should strictly extend classical IA in virtue of allowing an interval to be unbounded or empty.

Beyond this, various extensions of classical IA were considered. One choice that distinguishes between theories is: Are arithmetic operations purely algebraic, or do they involve topology? An example of the latter is containment set (cset) theory [B13], which extends functions over the reals to functions over the extended reals, e.g., $\sin(+\infty)$ is the set of all possible limits of $\sin x$ as $x \rightarrow +\infty$, which is $[-1, 1]$. The complications of this were deemed to outweigh the advantages, and it was agreed that operations should be purely algebraic.

Another choice is: Is an interval a set—a subset of the number line—or is it something different? The most widely used forms of IA are *set-based* and define an interval to be a set of real numbers [B10]. They have established software to find validated solutions of linear and nonlinear algebraic equations, optimization problems, differential equations, etc.

However, *Kaucher* IA and the nearly equivalent *modal* IA have significant applications. In the former, an interval is formally a pair (a, b) of real numbers, which for $a \leq b$ is “proper” and identified with the normal interval $\{x \in \mathbb{R} \mid a \leq x \leq b\}$, and for $a > b$ is “improper.” In the latter, an interval is a pair (X, Q) , where X is a normal interval and Q is a quantifier, either \exists or \forall . At the time of writing, it finds commercial use in the graphics rendering industry. Both forms are referred to as *Kaucher* IA henceforth.

In view of their significance, it was decided to support both set-based and *Kaucher* IA. Because of their different mathematical bases, this led to the concept of *flavors* (see Clause 7). A flavor is a version of IA that extends classical IA in a precisely defined sense, such that when only classical intervals and restricted operations are used (avoiding, e.g., division by an interval containing zero), all flavors produce the same results at the mathematical level and also—up to roundoff—in finite precision.

Currently, the standard includes only the set-based flavor. Among other possible flavors are *Kaucher*/modal intervals; containment-sets; and the interval system of Siegfried Rump [B15], which handles the relation between floating-point numbers and intervals, including overflow, in an elegant way, as well as being able to support open and half-open intervals. All of these extend classical IA in the defined sense.

^aThe numbers in brackets correspond to those of the bibliography in Annex A.

Clause 1 through Clause 9 contain a common set of definitions and requirements that apply to all flavors. Clause 10 through Clause 14 contain the set-based flavor, in which:

- Intervals are subsets of the set \mathbb{R} of real numbers. At the mathematical level (Level 1 in the structure defined in Clause 5) they are precisely all topologically closed and connected subsets of \mathbb{R} . The finite-precision level (Level 2) uses the notion of an interval type, which is a finite set of Level 1 intervals.
- The interval version of an elementary function such as $\sin x$ is essentially the natural extension to sets of the corresponding pointwise function on real numbers.

Fuzzy sets, like intervals, are a way to handle uncertain knowledge, and the two topics are related. However, to consider this relation was beyond the scope of this project.

Specification levels

The floating-point standard IEEE Std 754TM-2008 describes itself as layered into four Specification Levels. To manage complexity, the present standard uses a corresponding structure. It deals mainly with Level 1, of mathematical *interval theory*, and Level 2, the finite set of *interval datums* in terms of which finite-precision interval computation is defined. It has some concern with Level 3, of *representations* of intervals as data structures; and with Level 4, of interchange *encoding* in *bit strings*.

There is another important player: the programming language. It was a recognized omission of the first (1985) version of IEEE Std 754-2008 that it specified individual operations but not how they should be used in expressions. Optimizing compilers have, since well before that standard, used clever transformations so that it is impossible to know the precisions used and the roundings performed while evaluating an expression, or whether the compiler has even “optimized away” $(1.0 + x) - 1.0$ to become simply x . The 2008 revision specifies this by placing requirements on how operations should be used in expressions, though as of this writing, few programming languages have adopted that.

The lack of any restrictions is also a problem for intervals. Thus the standard makes requirements and recommendations on language implementations, thereby defining the notion of a standard-conforming implementation of intervals within a language.

The language does not constitute a fifth level in some linear sequence; from the user’s viewpoint, most current languages sit above datum level 2, alongside theory level 1, as a practical means to implement interval algorithms by manipulating Level 2 entities (though most languages have influence on Levels 3 and 4 also). This standard extends them to provide an instantiation of Level 2 entities.

The Fundamental Theorem

Moore’s [B8] Fundamental Theorem of Interval Arithmetic (FTIA) is central to interval computation. Roughly, it says as follows. Let f be an *explicit arithmetic expression*—that is, it is built from finitely many elementary functions (arithmetic operations) such as $+$, $-$, \times , \div , \sin , \exp , \dots , with no non-arithmetic operations such as intersection, so that it defines a real function $f(x_1, \dots, x_n)$. Then evaluating f “in interval mode” over any interval input box (x_1, \dots, x_n) is guaranteed to enclose (i.e., give a set that contains) the range of f over those inputs. Typically there are sub-cases, where extra conditions lead to stronger conclusions such as f being continuous on the input box.

A version of the FTIA holds in all variants of interval theory, but with varying hypotheses and conclusions. In the context of this standard, an expression should be evaluated entirely in one flavor, and inferences made strictly from that flavor’s FTIA; otherwise, a user might believe an FTIA holds in a case where it does not, with possibly serious effects in applications. As stated, the FTIA is about the mathematical level. Moore’s achievements were to see that “outward rounding” makes the FTIA hold also in finite precision and to follow through the consequences. An advantage of the level structure used by the standard is that the mapping between Levels 1 and 2 defines a framework where it is easily proved that

Each flavor’s finite-precision FTIA holds in any conforming implementation.

Generally, during program execution it can only be decided *after* evaluating an expression whether the conditions for any sub-case of the FTIA hold. The purpose of each flavor’s *decoration system* is to make such decisions computable, see 6.4 and Clause 7.

For the set-based flavor, see [B12] for background on its decoration system, 6.4 for a statement of its FTIA, and Annex B for a statement and proof of its FTDIA.

Operations

There are several interpretations of *evaluation outside an operation's domain* and *operations as relations rather than functions*. This includes classical alternative meanings of division by an interval containing zero, or square root of an interval containing negative values. To illustrate the different interpretations, consider $y = \sqrt{x}$ where $x = [-1, 4]$.

- a) In *optimization*, when computing lower bounds on the objective function, it is generally appropriate to return the result $y = [0, 2]$, and ignore the fact that $\sqrt{\cdot}$ has been applied to negative elements of x .
- b) In applications (such as solving differential equations) where one must check whether the hypotheses of a *fixed point theorem* are satisfied:
 - 1) one might need to be sure that the function is defined and continuous on the input and, hence, report an illegal argument when, as in the above case, this fails; or
 - 2) one might need the result $y = [0, 2]$, but must flag the fact that $\sqrt{\cdot}$ has been evaluated at points where it is undefined or not continuous.
- c) In *constraint propagation*, the equation is often to be interpreted as: find an interval enclosing all y such that $y^2 = x$ for some $x \in [-1, 4]$. In this case the answer is $[-2, 2]$.

The standard provides means to meet these diverse needs, while aiming to preserve clarity and efficiency. A language might achieve this by binding one of the above three interpretations—usually some variant of b)—to its built-in operations, and providing the others as library procedures.

In the context of flavors, a key idea is that of *common operation instances*: those elementary interval calculations that at the mathematical level are required to give the same result in all flavors. For example $[1, 2]/[3, 4] = [1/4, 2/3]$ is common, while division by an interval containing zero is not common.

Decorations

Many interval algorithms are only valid if certain mathematical conditions are satisfied: for instance, one might need to know that a function f , defined by an expression, is everywhere continuous on a box in \mathbb{R}^n defined by n input intervals x_1, \dots, x_n . The IEEE 754 use of global flags to record events such as division by zero was considered inadequate in an era of massively parallel processing. In this standard, such events are recorded locally by *decorations*.

A *decorated interval* is an ordinary interval tagged with a few bits that encode the decoration. A small number of decorations is provided, designed for efficient propagation of such property information. For instance, if evaluation outputs an interval y with the **dac** decoration, then f is *defined and continuous* on its input box, with range contained in y . This allows a rigorous check, for instance, that the conditions for applying a fixed-point theorem hold. Depending on need, a programmer may use bare (undecorated) intervals, or decorated intervals, or (if provided) the optional *compressed* decorated interval arithmetic that offers less decoration capability in exchange for faster execution.

The Basic standard

To make the standard more accessible and speed up production of implementations, a subset of the set-based standard called the Basic Standard for Interval Arithmetic (BSIA) has been written. It includes just one finite-precision interval type—intervals whose endpoints are IEEE 754 **binary64** numbers—and those operations that in the editors' view are most commonly used. A minimal implementation of the BSIA is not a conforming implementation of the full standard since some required operations of the latter are omitted or provided in a restricted form. However a program that runs using such an implementation should run, and give identical output within roundoff, using an implementation of the full standard. At the time of writing a project is underway to issue the BSIA as a separate IEEE standard.

Contents

Part 1. General Requirements	2
1. Overview	2
1.1. Scope	2
1.2. Purpose	2
1.3. Inclusions	2
1.4. Exclusions	2
1.5. Word usage	3
1.6. The meaning of conformance	3
1.7. Programming environment considerations	3
1.8. Language considerations	4
2. Normative references	4
3. Notation, abbreviations, and special terms	5
3.1. Frequently used notation and abbreviations	5
3.2. Special terms	5
4. Conformance	10
4.1. Conformance overview	10
4.2. Set-based interval arithmetic	11
4.2.1. IEEE 754 conformance	11
4.2.2. Compressed decorated interval arithmetic	11
4.3. Conformance claim	11
4.4. Implementation conformance questionnaire	12
5. Structure of the standard in levels	13
6. Functions and expressions	14
6.1. Function definitions	14
6.2. Expression definitions	15
6.3. Function libraries	17
6.4. The FTIA	18
6.5. Related issues	19
7. Flavors	19
7.1. Flavors overview	19
7.2. Flavor basic properties	20
7.3. Common evaluations	21
7.4. Primary versions and Level 1 interval versions	21
7.4.1. Arithmetic operations	21
7.4.2. Nonarithmetic operations required in all flavors	22
7.4.3. Flavor-defined nonarithmetic operations	22
7.5. The relation of Level 1 to Level 2	22
7.5.1. Types	23
7.5.2. Hull	23
7.5.3. Level 2 operations	23
7.5.4. Measures of accuracy	24
8. Decoration system	25
8.1. Decorations overview	25
8.2. Decoration definition and propagation	26
8.3. Recognizing common evaluation	26
9. Operations and related items required in all flavors	27
9.1. Arithmetic operations	27
9.2. Cancellative addition and subtraction	29
9.3. Set operations	29
9.4. Numeric functions of intervals	29
9.5. Boolean functions of intervals	29
9.6. Operations on/with decorations	29
9.7. All-flavor interval and number literals	30

9.7.1.	Overview	30
9.7.2.	All-flavor number literals	31
9.7.3.	Unit in last place	31
9.7.4.	All-flavor bare interval literals	31
9.7.5.	All-flavor decorated interval literals	31
9.7.6.	Grammar for all-flavor literals	32
9.8.	Constructors	32
Part 2.	Set-Based Intervals	34
10.	Level 1 description	34
10.1.	Non-interval Level 1 entities	34
10.2.	Intervals	34
10.3.	Hull	35
10.4.	Functions and expressions	35
10.5.	Required operations	36
10.5.1.	Interval literals	36
10.5.2.	Interval constants	36
10.5.3.	Forward-mode elementary functions	36
10.5.4.	Reverse-mode elementary functions	36
10.5.5.	Two-output division	37
10.5.6.	Cancellative addition and subtraction	38
10.5.7.	Set operations	38
10.5.8.	Constructors	38
10.5.9.	Numeric functions of intervals	39
10.5.10.	Boolean functions of intervals	39
10.6.	Recommended operations	40
10.6.1.	Forward-mode elementary functions	40
10.6.2.	Slope functions	41
10.6.3.	Boolean functions of intervals	41
10.6.4.	Extended interval comparisons	41
11.	The decoration system at Level 1	44
11.1.	Decorations and decorated intervals overview	44
11.2.	Definitions and basic properties	44
11.3.	The ill-formed interval	45
11.4.	Permitted combinations	45
11.5.	Operations on/with decorations	45
11.5.1.	Initializing	45
11.5.2.	Disassembling and assembling	46
11.5.3.	Comparisons	46
11.6.	Decorations and arithmetic operations	46
11.7.	Decoration of non-arithmetic operations	47
11.7.1.	Interval-valued operations	47
11.7.2.	Non-interval-valued operations	47
11.8.	User-supplied functions	47
11.9.	Notes on the <code>com</code> decoration	48
11.10.	Compressed arithmetic with a threshold (optional)	49
11.10.1.	Motivation	49
11.10.2.	Compressed interval types	49
11.10.3.	Operations	50
12.	Level 2 description	51
12.1.	Level 2 introduction	51
12.1.1.	Types and formats	51
12.1.2.	Operations	51
12.1.3.	Exception behavior	52

12.2.	Naming conventions for operations	52
12.3.	Tagging, and the meaning of equality at Level 2	52
12.4.	Number formats	53
12.5.	Bare and decorated interval types	54
12.5.1.	Definition	54
12.5.2.	Inf-sup and mid-rad types	55
12.6.	754-conformance	55
12.6.1.	Definition	55
12.6.2.	754-conforming mixed-type operations	55
12.7.	Multi-precision interval types	55
12.8.	Explicit and implicit types, and Level 2 hull operation	56
12.8.1.	Hull in one dimension	56
12.8.2.	Hull in several dimensions	56
12.9.	Level 2 interval extensions	56
12.10.	Accuracy of operations	57
12.10.1.	Measures of accuracy	57
12.10.2.	Accuracy requirements	58
12.10.3.	Documentation requirements	58
12.11.	Interval and number literals	58
12.11.1.	Overview	58
12.11.2.	Number literals	58
12.11.3.	Bare intervals	59
12.11.4.	Decorated intervals	59
12.11.5.	Grammar for portable literals	59
12.12.	Required operations on bare and decorated intervals	60
12.12.1.	Interval constants	60
12.12.2.	Forward-mode elementary functions	61
12.12.3.	Two-output division	61
12.12.4.	Reverse-mode elementary functions	61
12.12.5.	Cancellative addition and subtraction	61
12.12.6.	Set operations	62
12.12.7.	Constructors	62
12.12.8.	Numeric functions of intervals	63
12.12.9.	Boolean functions of intervals	64
12.12.10.	Interval type conversion	65
12.12.11.	Operations on/with decorations	65
12.12.12.	Reduction operations	65
12.13.	Recommended operations	66
12.13.1.	Forward-mode elementary functions	66
12.13.2.	Slope functions	66
12.13.3.	Boolean functions of intervals	66
12.13.4.	Extended interval comparisons	66
12.13.5.	Exact reduction operations	66
13.	Input and output (I/O) of intervals	67
13.1.	Overview	67
13.2.	Input	67
13.3.	Output	67
13.4.	Exact text representation	68
13.4.1.	Conversion of IEEE 754 numbers to strings	69
13.4.2.	Exact representations of comparable types	70
14.	Levels 3 and 4 description	70
14.1.	Overview	70
14.2.	Representation	70
14.3.	Operations and representation	71

14.4. Interchange representations and encodings	71
Annex A. Bibliography (informative)	74
Annex B. The fundamental theorem of decorated interval arithmetic for the set-based flavor (informative)	76
B1. Preliminaries	76
B2. The theorem	78

IEEE Standard for Interval Arithmetic

IMPORTANT NOTICE: IEEE Standards documents are not intended to ensure safety, security, health, or environmental protection, or ensure against interference with or from other devices or networks. Implementers of IEEE Standards documents are responsible for determining and complying with all appropriate safety, security, environmental, health, and interference protection practices and all applicable laws and regulations.

This IEEE document is made available for use subject to important notices and legal disclaimers. These notices and disclaimers appear in all publications containing this document and may be found under the heading Important Notice or Important Notices and Disclaimers Concerning IEEE Documents. They can also be obtained on request from IEEE or viewed at <http://standards.ieee.org/IPR/disclaimers.html>.

PART 1

General Requirements

1. Overview

1.1 Scope

This standard specifies basic interval arithmetic (IA) operations selecting and following one of the commonly used mathematical interval models. This standard supports the IEEE 754TM floating-point formats of practical use in interval computations. Exception conditions are defined, and standard handling of these conditions is specified. Consistency with the interval model is tempered with practical considerations based on input from representatives of vendors, developers and maintainers of existing systems.

The standard provides a layer between the hardware and the programming language levels. It does not mandate that any operations be implemented in hardware. It does not define any realization of the basic operations as functions in a programming language.

1.2 Purpose

The aim of the standard is to improve the availability of reliable computing in modern hardware and software environments by defining the basic building blocks needed for performing interval arithmetic. There are presently many systems for interval arithmetic in use; lack of a standard inhibits development, portability, and ability to verify correctness of codes.

1.3 Inclusions

This standard specifies

- Types for interval data based on underlying numeric formats, with a special class of type derived from IEEE 754 floating-point formats.
- Constructors for intervals from numeric and character sequence data.
- Addition, subtraction, multiplication, division, fused multiply add, square root; other interval-valued operations for intervals.
- Midpoint, radius and other numeric functions of intervals.
- Interval comparison relations and other boolean functions of intervals.
- Elementary interval functions of intervals.
- Conversions between different interval types.
- Conversions between interval types and external representations as text strings.
- Interval-related exceptional conditions and their handling.

1.4 Exclusions

This standard does not specify

- Which numeric formats supported by the underlying system shall have an associated interval type.
- How (for implementations supporting IEEE 754 arithmetic) operations act on the IEEE 754 status flags.
- How an implementation represents intervals at the level of programming language data types or bit patterns.

1.5 Word usage

In this standard three words are used to differentiate between different levels of requirements and optionality, as follows:

- **may** indicates a course of action permissible within the limits of the standard with no implied preference (“may” means “is permitted to”);
- **shall** indicates mandatory requirements strictly to be followed to conform to the standard and from which no deviation is permitted (“shall” means “is required to”);
- **should** indicates that among several possibilities, one is recommended as particularly suitable, without mentioning or excluding others; or that a certain course of action is preferred but not necessarily required; or that (in the negative form) a certain course of action is deprecated but not prohibited (“should” means “is recommended to”).

Further:

- **optional** indicates features that may be omitted, but if provided shall be provided exactly as specified;
- **can** is used for statements of possibility and capability (“can” means “is able to”);
- **might** indicates the possibility of a situation that could occur, with no implication of the likelihood of that situation (“might” means “could possibly”);
- **comprise** indicates members of a set or list are exactly those objects having some property. An unqualified **consist of** merely asserts all members of a set have some property, e.g., “a binary floating-point format consists of numbers with a terminating binary representation.” “Comprises” means “consists exactly of.”
- **Example** introduces text that is informative (is not a requirement of this standard). Such text is set in slanted sanserif font within brackets. *[Example. Like this.]*

1.6 The meaning of conformance

Clause 4 lists the requirements on a conforming implementation in summary form, with references to where these are stated in detail.

1.7 Programming environment considerations

This standard does not define all aspects of a conforming programming environment. Such behavior should be defined by a programming language definition supporting this standard, if available; otherwise by a particular implementation. Some programming language specifications might permit some behaviors to be defined by the implementation.

Language-defined behavior should be defined by a programming language standard supporting this standard. Standards for languages intended to reproduce results exactly on all platforms are expected to specify behavior more tightly than do standards for languages intended to maximize performance on every platform.

Because this standard requires facilities that are not currently available in common programming languages, the standards for such languages might not be able to conform fully to this standard if they are no longer being revised. If the language can be extended by a function library or class or package to provide a conforming environment, then that extension should define all the language-defined behaviors that would normally be defined by a language standard.

Implementation-defined behavior is defined by a specific implementation of a specific programming environment conforming to this standard. Implementations define behaviors not specified by this standard nor by any relevant programming language standard or programming language extension. Conformance to this standard is a property of a specific implementation of a specific programming environment, rather than of a language specification. However, a language standard could also be said to conform to this standard if it were constructed so that every conforming implementation of that language also conformed automatically to this standard.

1.8 Language considerations

All relevant languages are based on the concepts of data and transformations. In procedural languages, data are held in variables, which are transformed by assignment. In functional languages, input data are supplied as arguments; the transformed form is returned as results. Dataflow languages vary considerably, but use some form of the data and transformation approach.

Similarly, all relevant languages are based on the concept of mapping the pseudo-mathematical notation that is the program code to approximate real arithmetic, almost exclusively using some form of floating-point. The unit of mapping and transformation might be individual operations and built-in functions, expressions, statements, complete procedures, or other. This standard is applicable to all of these.

In this standard, the units of transformation are called operations and written as named functions; in a specific implementation they might be represented by operators (e.g., using an infix notation), or by families of type-specific functions, or by operators or functions whose names might differ from those used here.

The least requirement on a conforming language standard, compiler or interpreter is that it shall:

- a) define bindings so that the programmer can specify Level 2 data (in the sense of the levels defined in Clause 5) as described in this standard;
- b) define bindings so that the programmer can specify the operations on such data as described in this standard;
- c) define any properties of such data and operations that this standard requires to be defined;
- d) honor the rules of interval transformations on such data and operations as described in this standard; such units of transformation that the language standard, compiler or interpreter uses.

Specifically, if the data before and after the unit of transformation are regarded as sets of mathematical intervals, the transformed form of all combinations of the elements (the real values) represented by the prior set shall be a member of the posterior set.

If a conforming language standard supports reproducible interval arithmetic, it shall also:

- e) Use the data bindings as specified in point a) above for reproducible operations;
- f) Define bindings to the reproducible operations as described in this standard;
- g) Define any modes and constraints that the programmer needs to specify or obey in order to obtain reproducible results.

If a conforming language standard supports both non-reproducible and reproducible interval arithmetic, it shall also:

- h) Permit a reproducible transformation unit to be used as a component in a non-reproducible program, possibly via a suitable wrapping interface.

2. Normative references

The following referenced documents are indispensable for the application of this document (i.e., they must be understood and used, so each referenced document is cited in text and its relationship to this document is explained). For dated references, only the edition cited applies. For undated references, the latest edition of the referenced document (including any amendments or corrigenda) applies.

IEEE Std 754TM-2008, IEEE Standard for Floating-Point Arithmetic¹².

¹IEEE publications are available from The Institute of Electrical and Electronics Engineers at <http://standards.ieee.org>.

²The IEEE standards or products referred to in this clause are trademarks of The Institute of Electrical and Electronics Engineers, Inc.

3. Notation, abbreviations, and special terms

3.1 Frequently used notation and abbreviations

IA	Interval arithmetic.
\mathbb{R}	the set of real numbers.
$\overline{\mathbb{R}}$	the set of extended real numbers, $\mathbb{R} \cup \{-\infty, +\infty\}$.
$\overline{\overline{\mathbb{R}}}$	the set of closed real intervals, including unbounded intervals and the empty set.
\mathbb{F}	generic notation for a number format.
\mathbb{T}	generic notation for an interval type.
\emptyset , Empty	the empty set.
Entire	the whole real line.
NaN	Not an Interval.
NaN	Not a Number.
qNaN, sNaN	quiet and signaling NaN.
x, y, \dots [resp. f, g, \dots]	typeface/notation for a numeric value [resp. numeric function].
$\mathbf{x}, \mathbf{y}, \dots$ [resp. $\mathbf{f}, \mathbf{g}, \dots$]	typeface/notation for an interval value [resp. interval function].
f, g, \dots	typeface/notation for an expression, producing a function by evaluation.
$\text{Dom}(f)$	the domain of a point-function f .
$\text{Rge}(f \mid \mathbf{s})$	the range of a point-function f over a set \mathbf{s} ; the same as the image of \mathbf{s} under f .
Val	map from member of concrete set to represented value in abstract set (e.g., from number format \mathbb{F} to $\overline{\mathbb{R}}$).

3.2 Special terms

The labels (AF) and (S) mark special terms that apply to all flavors, or the set-based flavor only, respectively.

IEEE 754 format: (AF) A floating-point format that together with its associated operations conforms to IEEE Std 754-2008. A **basic IEEE 754 format** is one of the five formats `binary32`, `binary64`, `binary128`, `decimal64`, `decimal128`.

IEEE 754 conforming implementation: (S) An implementation of this standard, or a part thereof, that is built on an IEEE 754 system, uses only IEEE 754 conforming types, and satisfies the extra requirements for IEEE 754 conformance.

NOTE—Details in 12.6..

IEEE 754 conforming type: (S) An inf-sup interval type derived from an IEEE 754 format, with its relevant operations; or the associated decorated type with its operations.

IEEE 754 system: (AF) A programming environment, made up of hardware or software or both, that provides floating-point arithmetic conforming to IEEE Std 754-2008.

accuracy mode: (AF) A way to describe the quality of an interval version of a function. See **tightness**.

NOTE—Details in (AF) 7.5.4 and (S) 12.10.1.

arithmetic operation: (AF) A **version** of a **point operation**.

NOTE—Details in 6.1.

A **basic arithmetic operation** of IEEE 754 is one of the six functions `+`, `-`, `×`, `÷`, fused multiply-add `fma` and square root `sqrt`.

arity: (AF) The number of arguments of a function.

NOTE—Details in 6.3.

bare interval: (AF) Same as **interval**; used to emphasize it is not decorated.

bounds: (S) If x is an interval then $\underline{x} = \inf x$ and $\bar{x} = \sup x$, its **lower bound** and **upper bound** respectively, are extended-real numbers that for a nonempty interval satisfy $\underline{x} \leq \bar{x}$, $\underline{x} < +\infty$ and $\bar{x} > -\infty$. By convention the empty set has lower bound $+\infty$ and upper bound $-\infty$. Note $\pm\infty$ can be bounds of an interval but never members of it.

(AF) In any flavor, bounds have the above meaning for common intervals. A flavor may define a notion of bound for non-common intervals. E.g., a Kaucher flavor might define it such that the improper interval $[2, 1]$ has lower bound 2.

box: (AF) A **box** or **interval vector** is an n -dimensional interval, i.e. a tuple (x_1, \dots, x_n) where the x_i are intervals. When an interval may be regarded as a set of points, it may be identified with the cartesian³ product $x_1 \times \dots \times x_n$. (S)
NOTE—Details in 10.2.

common evaluation: (AF) Those evaluations of an operation (over common intervals) that have the same value in all flavors.
NOTE—Details in 7.3 and Clause 9.

common interval: (AF) Those intervals that “modulo the embedding map” belong to all flavors.
NOTE—Details in 7.2.

comparable: (AF) A set of interval types is comparable if for any two of them, one is wider than the other, that is if, regarded as sets of mathematical intervals, the types are linearly ordered by set inclusion.

compressed: (S) Compressed interval types provide a limited form of decorated interval computation, in a storage-efficient way aimed at giving faster execution. They are for certain applications where a decoration value denotes some exceptional condition, and one needs to record either an interval value or a decoration value but never both at once.
NOTE—Details in 11.10.

conforming part: (AF) A subset (possibly the whole) of an implementation’s types and formats, with associated operations, that forms a conforming implementation in its own right. In a multi-flavor implementation, it might be formed from a subset of the flavors.
NOTE—Details in 4.1.

constant: (AF) A **real constant** at Level 1 is defined to be a scalar point function with no arguments; this includes the constant NaN.
NOTE—Details in 6.1.

contain: (AF) A partial order relation between intervals that shall be defined in each flavor, and for common intervals shall coincide with normal set containment. Synonym **enclose**.
NOTE—Details in 7.2.

datum: (AF) One of the entities manipulated by finite precision (Level 2) operations of this standard. It can be a **boolean**, a **decoration**, a (bare) **interval**, a **decorated interval**, a **number** or a **string** datum. Number datums are organized into formats; bare and decorated interval datums are organized into types. An \mathbb{F} -datum or \mathbb{T} -datum means a member of the number format \mathbb{F} , or of the bare or decorated interval type \mathbb{T} .
NOTE—Details in Clause 5, 7.5.1, (s) 12.1.

decoration: (AF) One of the flavor-defined sets of values used to record events, called exceptional conditions, in interval operations: e.g., evaluating a function at points where it is undefined.
NOTE—Details in Clause 8.

(S) One of the five values **com**, **dac**, **def**, **trv** and **ill**.
NOTE—Details in Clause 11.

decorated interval: (AF) A pair (interval, decoration).
NOTE—Details in 8, see also 7.5.1.

³The ultimate accolade to a mathematician is to be uncapitalized: cartesian, hamiltonian, noetherian, ...

domain: (AF) For a function with arguments taken from some set, the **domain** comprises those points in the set at which the function has a value. The domain of a point operation is part of its definition. E.g., the (point) operation of division x/y , in this standard, has arguments (x, y) in \mathbb{R}^2 , and its domain is $\{(x, y) \in \mathbb{R}^2 \mid y \neq 0\}$. See also **natural domain**.

elementary function: (AF) Synonymous with **arithmetic operation**.

enclose: (AF) Synonymous with **contain**.

exception, exceptional condition: (AF) An **exception** is an event that occurs, and may be signaled, when an operation on some particular operands has no outcome suitable for every reasonable application. A signaled exception is handled in a language- or implementation-defined way.

NOTE—Details in 7.5.

An **exceptional condition** is one of the events handled by the decoration system; it is not an exception.

NOTE—Details in Clause 8.

(S) In the set-based standard, exceptions may occur in interval constructors and in the **intervalPart** operation.

NOTE—Details in 12.12.11 and in 12.1.3, 13.4 respectively.

expression: (AF) A symbolic form used to define a function. As used in this standard, it is a mathematical construct extracted from a section of run-time dataflow of a program, from which statements about the behavior of this function may be deduced by using the FTIA or FTDIA of a given flavor. It may be represented by a computational graph, a code list or a normal algebraic expression. An **arithmetic expression** is one whose operations are all arithmetic operations.

NOTE—Details in Clause 6.

explicit type: (S) An interval type that has a uniquely defined interval hull operation.

floating-point format: (AF) A number format like those of IEEE 754, whose numbers have the form $x = s \times d_0.d_1 \dots d_p \times b^e$ where integer $b \geq 2$ is the fixed radix, integer $p \geq 0$ is the fixed precision, $s = \pm 1$ is the sign, $d_0.d_1 \dots d_p$ (a radix- b fraction) is the significand or mantissa, and e is an integer in a fixed exponent range $emin \leq e \leq emax$.

fma: (AF) Fused multiply-add operation, that computes $x \times y + z$. One of the basic arithmetic operations.

NOTE—Details in 9.1.

format: (Or **number format**.) (AF) One of the sets into which number datums are organized at Level 2, usually regarded as a finite set \mathbb{F} of extended-reals with possibly also -0 , $+0$.

If \mathbb{F} is a format, an **\mathbb{F} -datum** means a member of \mathbb{F} , and an **\mathbb{F} -number** means a non-NaN member of \mathbb{F} . The **value** map Val maps \mathbb{F} -numbers to the extended-reals they denote: $\text{Val}(-0) = \text{Val}(+0) = 0$ and $\text{Val}(x) = x$ for other \mathbb{F} -numbers.

(S) A format is **provided** if the implementation provides a representation of it, and **supported** if also it has properties specified in 12.4.

FTIA, FTDIA: (AF) For a real function f defined by an arithmetic expression, and in a given flavor:

- The Fundamental Theorem of Interval Arithmetic (FTIA) gives relations between the result **y** of evaluating the expression in interval mode over an input box and the behavior of f on the box. The basic property is that **y** contains, in a flavor-defined sense, the range of f over the box; various extra conditions can give extra conclusions, such as that f is continuous on the box.
- The Fundamental Theorem of Decorated Interval Arithmetic (FTDIA) describes how evaluating the expression in decorated interval mode enables the various conditions of the FTIA to be verified, making the corresponding conclusions computable.

NOTE—Details in 6.4.

function: (AF) Has the usual mathematical meaning of a (possibly partial, i.e., not everywhere defined) function.

NOTE—Details in 6.1. Synonymous with **map**, **mapping**.

hull: (AF) (Or **interval hull**.) When not qualified by the name of an interval type, the hull of a subset s of \mathbb{R} is the Level 1 hull, namely the tightest interval containing s in a flavor’s sense of “contain.”

NOTE—Details in 7.5.2.

(s) When \mathbb{T} is an explicit type, the \mathbb{T} -hull of s is the unique tightest \mathbb{T} -interval containing s .

When \mathbb{T} is an implicit type, the \mathbb{T} -hull of s is a minimal \mathbb{T} -interval containing s as specified in the definition of the type.

implementation: (AF) When used without qualification, means a realization of an interval arithmetic conforming to the specification of this standard.

NOTE—Details in Clause 4.

implicit type: (s) An interval type that does not have a uniquely defined interval hull operation: the hull must be specified as part of the definition of the type.

inf-sup: (AF) Describes a representation of an interval based on its lower and upper bounds.

interval: (AF) At Level 1 the entities of the abstract data type forming one part of the interval model of an interval flavor. At Level 2 a \mathbb{T} -interval is a member of the bare interval type \mathbb{T} , or a non-NaI member of the decorated interval type \mathbb{T} .

NOTE—Details in 7.2, 7.5.1.

(s) A (bare) interval x , see 10.2, is a closed connected subset of \mathbb{R} . The set of all intervals is denoted $\overline{\mathbb{IR}}$. It comprises the empty set \emptyset and the nonempty intervals $[\underline{x}, \overline{x}] = \{x \in \mathbb{R} \mid \underline{x} \leq x \leq \overline{x}\}$. See also **bounds**.

interval extension: (s) At Level 1, an interval extension of a point function f is a function \mathbf{f} from intervals to intervals such that $f(x)$ belongs to $\mathbf{f}(x)$ whenever x belongs to x and $f(x)$ is defined. It is the **natural** (or tightest) interval extension if $\mathbf{f}(x)$ is the interval hull of the range of f over x , for all x .

NOTE—Details in 10.4. For Level 2 interval extension, see 12.9.

A **decorated interval extension** of f is a function from decorated intervals to decorated intervals, whose interval part is an interval extension of f , and whose decoration part propagates decorations as specified in 11.6.

interval vector: (AF) See **box**.

library: (AF) In a given flavor, the set of Level 1 operations (Level 1 library) or those provided by an implementation (Level 2 library). Further classification may be made into the **point library**, **bare interval library** and **decorated interval library**.

NOTE—Details in 6.3, 9, see also 6.2.

literal: (AF) A text string that denotes a Level 1 value, typically within input to a constructor and (s) in I/O. The standard uses **integer**, **number**, **decoration**, **bare interval** and **decorated interval** literals, which denote values of the indicated kinds.

NOTE—Details in 9.7, Clause 13.

map, mapping: (AF) See **function**.

mathematical interval of constructor: (AF) The arguments of an interval constructor, if valid, define a mathematical interval x . The actual interval returned by the constructor is the tightest interval of the destination type that contains x , see 9.8 and 7.5.3.

mid-rad: (AF) Describes a representation of an interval based on its midpoint and radius.

NaI, NaN: (AF) At Level 1 NaN is the function $\mathbb{R}^0 \rightarrow \mathbb{R}$ with empty domain. NaI is the natural interval extension of NaN, the map from subsets of \mathbb{R}^0 to subsets of \mathbb{R} whose value is always the empty set.

NOTE—Details in 6.1, including the meaning of \mathbb{R}^0 .

(s) At Level 2 NaN is the Not a Number datum, which is a member of every supported number format. NaI is the Not an Interval datum, which is a member of every decorated interval type.

NOTE—Details in 11.3.

(AF) Other flavors may provide NaN and NaI.

NOTE—Details in 7.5.1.

narrower: (AF) See **wider**.

natural domain: (AF) For an arithmetic expression $f(z_1, \dots, z_n)$, the natural domain is the set of $x = (x_1, \dots, x_n) \in \mathbb{R}^n$, where the expression defines a value for the associated point function $f(x)$, see Clause 6.

no value: (AF) A mathematical (Level 1) operation evaluated at a point outside its domain is said to have no value. Used instead of “undefined,” which can be ambiguous. E.g., in this standard, real number division x/y has no value when $y = 0$.

NOTE—Details in 6.1.

non-arithmetic operation: (AF) An operation on intervals that is not an interval version of a point function; includes intersection and convex hull of two intervals.

NOTE—Details in 6.1.

number: (AF) Any member of the set $\mathbb{R} \cup \{-\infty, +\infty\}$ of extended reals: a **finite number** if it belongs to \mathbb{R} , else an **infinite number**.

NOTE—Details in 10.1.

number format: (AF) See **format**.

octet: (AF) Bit string of length 8, equivalently 8-bit byte.

An **octet-encoding** is a mapping from the conceptual bit string encodings of floating-point datums and of decorations into an octet sequence.

NOTE—Details in 14.4.

operation: (AF) In a given flavor, synonymous with **supported function**.

point function: (AF) A mathematical (Level 1) function of real variables.

NOTE—Details in 6.1.

primary version: (AF) A Level 1 version of a generic function, from which all other versions are derived.

NOTE—Details in 6.1.

provided operation: (AF) In a given implementation of a flavor, an operation for which the implementation provides at least one Level 2 version.

NOTE—Details in 6.1.

range: (AF) The range $\text{Rge}(f | s)$ of a function f over a set s is the set of values $f(x)$ at those points x of s where f is defined.

NOTE—Details in 6.1.

string, text: (AF) A **text string**, or just **string**, is a finite character sequence belonging to some alphabet, see 10.1. The term **text** is also used to mean strings generally; e.g., an operation having “numeric or text input” means each input is a number or a string.

supported function: (AF) In a given flavor, a function for which the flavor defines the primary version.

NOTE—Details in 6.1.

tightest: (AF) Smallest in the partial order of set containment, or in the “contains” relation of a flavor. The tightest set (unique, if it exists) with a given property is contained in every other set with that property.

Also denotes one of the accuracy modes, see **tightness**.

tightness: (AF) The strongest accuracy mode (in (s) these are *tightest*, *accurate*, *valid*) that a given operation, in a given type, achieves for some input box, or uniformly over some set of inputs.

NOTE—Details in 7.5.4, 12.10.

type: (AF) (Or **interval type**.) One of the sets into which bare and decorated interval datums are organized at Level 2, usually regarded as a finite set \mathbb{T} of Level 1 intervals, (a **bare type**) in the bare case; and of Level 1 decorated intervals together with the value NaI, (a **decorated type**) in the decorated case, see 12.5. Each bare type has a corresponding decorated type, and vice versa.

(S) A type is **provided** if the implementation provides a representation of it, and **supported** if also it has the properties specified in 12.5.

(AF) If \mathbb{T} is a type, a **\mathbb{T} -datum** means a member of \mathbb{T} . A **\mathbb{T} -interval** for bare interval types means the same as a \mathbb{T} -datum, and for decorated types means a non-NaI member of \mathbb{T} .

version: (AF) A version of an operation is any of the actual operations denoted by a generic operation name. NOTE—Details in 6.1.

wider, narrower: (AF) An interval type \mathbb{T}' is wider than a type \mathbb{T} , and \mathbb{T} is narrower than \mathbb{T}' , if \mathbb{T} is a subset of \mathbb{T}' when they are regarded as sets of Level 1 intervals, ignoring the type tags and possible decorations, see 12.5.1. Wider means having more precision (IEEE Std 754-2008).

4. Conformance

4.1 Conformance overview

This standard defines conformance for programming environments for interval arithmetic. A programming environment is a collection of processes for developing and executing computer programs. An *implementation*, in the following, is a programming environment that implements the Level 2 behavior of interval arithmetic as specified in this standard. The implementation may include a programming language as well as extensions in the form of libraries, classes, or packages that are necessary to satisfy the requirements for conformance. Requirements are given for the whole implementation; whether they are satisfied by the language itself or by an extension is irrelevant for conformance, see 1.7.

Conformance requirements in this standard follow the guidelines developed by OASIS, see [B11]. In particular the OASIS concept of profiles is used to structure the requirements for a conforming implementation of a flavor of interval arithmetic.

A conforming implementation shall provide at least one of the following profiles, defined in Clause 7, called *standard flavors* of interval arithmetic in the context of this standard. A *profile* in this sense is a specific subset of functionality and requirements of this standard that may be provided by an implementation, but if it is provided shall be implemented exactly as specified in the corresponding section. A *sub-profile* is a nested profile. If the main profile is provided, an implementation may also support the sub-profile. If it does, it shall do so as specified.

An implementation may also provide additional flavors that shall conform to the core specification summarized in 7.1 and detailed in the following subclauses of Clause 7 as well as Clause 8, Clause 9. The decoration system of each additional flavor shall support FTIA sub-cases as described in 6.4.

NOTE—Any person proposing a new flavor should submit it as a revision of this standard, see the item “Comments on standards” in the front matter of this document.⁴

In any flavor of interval arithmetic, a *part* of an implementation is specified by a subset of the implementation’s interval types (and where necessary to remove ambiguity, of its number formats). It comprises those types and formats, and associated operations of the implementation. It is a *conforming part* if it forms a conforming implementation of that flavor in its own right. Note that this implies the full set of decorations of the flavor shall be supported.

⁴Notes in text, tables, and figures of this standard are given for information only and do not contain requirements needed to implement the standard.

4.2 Set-based interval arithmetic

An implementation of the set-based flavor, Clause 10 through Clause 14, shall satisfy these requirements:

- provide the decorations specified in 11.2;
- provide at least one supported bare interval type, see 12.5;
- if multi-precision interval types are supported, define them as a parameterized sequence of interval types, see 12.7;
- provide implementations of the required operations in 12.12; required and recommended accuracies for these operations are in 12.10;
- if any of the recommended operations referenced in 12.13 are provided, provide them in a way that satisfies the requirements specified in the same clause; and
- provide input and output functions to convert intervals from and to strings as well as a public representation as specified in 13.2, 13.3, and 13.4; the string conversions shall satisfy containment in the general case and satisfy accuracy requirements for IEEE 754 conforming types as described in 13.2, 13.3, 13.4.

In all these, the implementation shall follow the representation rules defined in 14.1, essentially stating that booleans, strings, and decorations are represented as given in the references listed in the preceding bullets and that there is a one-to-one correspondence between the abstract number and interval formats in said references and the concrete formats used in the implementation.

As 14.2 states, each Level 2 datum shall be represented by at least one Level 3 object, and each Level 3 object shall represent at most one Level 2 datum.

4.2.1 IEEE 754 conformance

In addition to basic conformance of the set-based flavor, an implementation may claim IEEE 754 conformance for all or part of the set-based flavor if the requirements of 12.6 are satisfied.

4.2.2 Compressed decorated interval arithmetic

An implementation may support the compressed arithmetic sub-profile of the set-based flavor. If compressed arithmetic is supported, it shall be as described in 11.10, in particular the implementation shall:

- provide an inquiry⁵ function to distinguish between intervals and decorations;
- provide a constructor for compressed intervals for each threshold value;
- provide a conversion function from compressed intervals to decorated intervals of the parent type;
- follow a *worst case semantic* for all arithmetic operations on compressed intervals; and
- provide compressed arithmetic implementations of the required operations in 12.12.

4.3 Conformance claim

An implementation may claim its conformance to this standard in the following way.

[*Name of implementation and version*] is conforming to IEEE Std 1788-2015 Standard for Interval Arithmetic. It is conforming to the set-based flavor with IEEE 754 conformance for [*list of IEEE 754 conforming supported interval types*] and [*with / without*] compressed arithmetic. Additionally it provides [*list of non-standard flavors*].

Part of the conformance claim shall be the completion of the questionnaire in 4.4.

⁵In UK English “enquiry” is correct usage.

4.4 Implementation conformance questionnaire

a) Implementation-defined behavior

- 1) What status flags or other means to signal the occurrence of certain decoration values in computations does the implementation provide if any? (see 8.1).

Does the implementation provide the **set-based flavor**? If so answer the following set of questions.

b) Documentation of behavior

- 1) If the implementation supports implicit interval types, how is the interval hull operation realized? The answer may be given via an appropriate algorithm, see 12.8.
- 2) What accuracy is achieved (i.e., tightest, accurate, or valid) for each of the implementation's interval operations? (see 12.10).
- 3) Under what conditions is a constructor unable to determine whether a Level 1 value exists that corresponds to the supplied inputs? (see 12.12.7).
- 4) How are cases for rounding a Level 1 value to an \mathbb{F} -number handled that are not covered by the rules given in 12.12.8? This includes: how is the distance to an infinity calculated when rounding a number? how are ties broken in rounding numbers if multiple numbers qualify as the rounded result?
- 5) How are interval datums converted to their exact text representations? (see 13.4).

c) Implementation-defined behavior

- 1) Does the implementation include the interval overlapping function? (see 10.6.4). If so, how is it made available to the user?
- 2) Does the implementation store additional information in a NaI? What functions are provided for the user to set and read this information? (see 11.3).
- 3) What means if any does the implementation provide for an exception to be signaled when a NaI is produced? (see 11.3).
- 4) What interval types are supported besides the required ones? (see 12.1.1).
- 5) What mechanisms of exception handling are used in exception handlers provided by the implementation? (see 12.1.3). What additional exception handling is provided by the implementation?
- 6) What is the tie-breaking method used in rounding of supported number formats \mathbb{F} that are not IEEE 754 conforming? (see 12.12.8).
- 7) Does the implementation include different versions of the same operation for a given type and how are these provided to the user? (see 12.12).
- 8) What combinations of formats are supported in interval constructors? (see 12.12.7).
- 9) What is the tightness of the result of constructor calls in cases where the standard does not specify it? (see 12.12.7).
- 10) What methods are used to read or write strings from or to character streams? (see 13.1). Does the implementation employ variations in locales (such as specific character case matching)? This includes the syntax used in the strings for reading and writing.
- 11) What is the tightness of the string to interval conversion for non IEEE 754 conforming interval types and the tightness for the interval to string conversion for all interval types? (see 13.2, 13.3).
- 12) What is the result of Level 3 operations for invalid inputs? (see 14.3).
- 13) What are the interchange representations of the fields of the standard Level 3 representation listed in 14.4?

- 14) What decorations does the implementation provide and what is their mathematical definition? (see 8.2). How are these decorations mapped when converting an interval to the interchange format? (see 14.4).
- 15) What interchange formats if any are provided for non IEEE 754 interval formats and on non IEEE 754 systems? (see 14.4). How are these provided to the user?

Does the implementation support the **compressed arithmetic** sub-profile of the set-based profile? If so answer the following set of questions.

d) Implementation-defined behavior

- 1) Which compressed interval types are provided? (see 11.10).

Does the implementation provide **non-standard flavors** not defined in this standard? (see 7.1). If so answer the following questions for each additional flavor.

e) Flavor definition

- 1) What is the set \mathfrak{F} of intervals of the flavor? And what is the embedding map describing the relation to the common intervals? (see 7.2).
- 2) What decorations does the flavor provide, and what is their mathematical definition?
- 3) What operations besides the required operations does the implementation provide? (see 7.2). How are the \mathfrak{F} -versions of all provided operations—required and flavor-specific—defined? (see 7.4).
- 4) What interval types are provided for the intervals of \mathfrak{F} ? (see 7.5.1).
- 5) What is the result of applying the flavor’s **convertType** operation to a non-interval datum? (see 7.5.2).
- 6) For what interval types \mathbb{T} of the flavor shall the implementation provide a \mathbb{T} -version of the flavor-specific operations? (see 7.5.3). For what interval types should the implementation provide a \mathbb{T} -version?
- 7) If possible, give a recommendation or requirements on the relation between the Level 2 number format \mathbb{F} and the interval type \mathbb{T} when passing from a Level 1 operation to a \mathbb{T} -version, see 7.5.3.
- 8) How are exceptional cases in the evaluation of a \mathbb{T} -version of an operation handled? (see 7.5.3).
- 9) Does the flavor define accuracy modes in addition to valid and tightest? Where do these accuracy modes apply? (see 7.5.4).

5. Structure of the standard in levels

For each flavor, the standard is structured into four levels, matching those defined in IEEE Std 754-2008 Table 3.1. They are summarized in Table 5.1.

Level 1, *mathematics*, defines the flavor’s underlying theory. The entities at this level are mathematical intervals and operations on them. An implementation of the flavor shall implement this theory. In addition to an ordinary (bare) interval, this level defines a *decorated* interval, comprising a bare interval and a *decoration*. In all flavors, decorations implement the standard’s way of handling exceptional conditions in interval operations.

Level 2, *discretization*, is the central part of the standard, approximating the mathematical theory by an implementation-defined finite set of entities and operations. A Level 2 entity is called a *datum*⁶.

Interval datums are organized into finite sets called *interval types*. An interval datum is a mathematical interval tagged by a symbol that indicates its type: an interval that “knows its type.” The type abstracts a particular way of representing intervals (e.g., by storing their lower and upper bounds as IEEE **binary64** numbers). Most Level 2 arithmetic operations act on intervals of a given type to produce an interval of the same type.

⁶Plural “datums” in this standard, since “data” is often misleading.

Table 5.1. Specification levels for interval arithmetic

Relationships between specification levels for interval arithmetic for a given flavor \mathfrak{F} and a given finite-precision bare interval type \mathbb{T} of \mathfrak{F} .		
Level 1	Number system used by flavor \mathfrak{F} . Set of mathematical \mathfrak{F} -intervals. Principles of how $+$, $-$, \times , \div and other arithmetic operations are extended to \mathfrak{F} -intervals.	Mathematical interval model
	$\downarrow \mathbb{T}$ -interval hull total, many-to-one b)	<i>identity map</i> \uparrow total, one-to-one a)
Level 2	A finite subset \mathbb{T} of the \mathfrak{F} -intervals—the \mathbb{T} -interval datums—and operations on them.	Discretization
		<i>represents</i> \uparrow partial, many-to-one, onto c)
Level 3	Representations of \mathbb{T} -interval datums, e.g., by two floating-point datums.	Representation
		<i>encodes</i> \uparrow partial, many-to-one, onto d)
Level 4	Bit strings 0111000...	Encoding

Level 3 is about *representation* of interval datums—usually but not necessarily in terms of floating-point values. A Level 3 entity is an *interval object*. Representations of decorations, hence of decorated intervals, are also defined at this level.

Level 4 is about *encoding* of interval objects into bit strings.

The Level 3 and 4 requirements in this standard are few, and mainly concern mappings from internal representations to external ones, such as interchange types.

The arrows in Table 5.1 denote mappings between levels. The phrases in italics name these mappings. Each phrase “total, many-to-one,” etc., labeled with a letter (a) to (d), is descriptive of the mapping and is equivalent to the corresponding labeled fact below.

- a) Ignoring the type-tag, an interval datum *is* a mathematical interval.
- b) For each type \mathbb{T} , each mathematical interval has a unique interval datum as its \mathbb{T} -*hull*—a minimal enclosing interval of that type. This is with respect to a meaning of “contain.” For the set-based flavor this is normal set inclusion, but in other flavors, e.g., Kaucher, might not always mean the same as set inclusion.
- c) Not every interval object necessarily represents an interval datum, but when it does, that datum is unique. Each interval datum has at least one representation and might have more than one.
- d) Not every interval encoding necessarily encodes an interval object, but when it does, that object is unique. Each interval object has at least one encoding and might have more than one.

NOTE—Items c) and d) are standard and necessary properties of representations. By contrast, the properties a) and b) of the maps from Level 1 to Level 2, and back, are fundamental design decisions of the standard.

6. Functions and expressions

6.1 Function definitions

In this document the terms **function**, **map** and **mapping** are synonymous and have the usual mathematical meaning (see 1.7 for general considerations) while **operation** has a specialized meaning. The following summarizes usage.

- a) A (partial) **function** from a set X to a set Y associates to each point x in some subset of X , called the **domain** $D = \text{Dom } f$ of f , a unique point $f(x)$ in Y called the value of f at x . At points outside D the

function has no value. The notation $f : X \rightarrow Y$ means that f is a function from X to Y . A **total** function is one for which $D = X$.

NOTE—Various formal set-theoretic foundations exist. Where relevant, this document uses a model where the notion of *ordered pair* (x, y) is a primitive; a cartesian product of $X \times Y$ of two sets is the set of all (x, y) with $x \in X, y \in Y$; a function is formally identical with its *graph*, the set $\{(x, f(x)) \mid x \in \text{Dom } f\} \subseteq X \times Y$.

Then $(X \times Y) \times Z$ is shortened to $X \times Y \times Z$ and its elements $((x, y), z)$ to (x, y, z) , and so on inductively to define $X_1 \times X_2 \times \cdots \times X_n$ with elements (x_1, \dots, x_n) (called “tuples” or “vectors”). This product may be identified with (but formally is different from) the set of all maps x defined on $\{1, 2, \dots, n\}$ and such that $x(i) \in X_i$ for each i ; and x_i is an alternative notation for $x(i)$.

- b) The **range** of f over a subset S of X is the set $\text{Rge}(f \mid S) = \{f(x) \mid x \in S \cap \text{Dom } f\}$. Points outside the domain are ignored—e.g., the range of the real square root function over $[-4, 4]$ is $[0, 2]$.

Vector notation $f(x) = f(x_1, \dots, x_n)$ may be used when X is a cartesian product $X = X_1 \times \cdots \times X_n$ as above. The range may then be written $\text{Rge}(f \mid S_1, \dots, S_n)$ when $S = S_1 \times \cdots \times S_n$ with $S_j \subseteq X_j$ for each j .

- c) Meaning of \mathbb{R}^0 . A *cartesian power* X^n means $\underbrace{X \times \cdots \times X}_{n \text{ times}}$. Take $X = \mathbb{R}$. As noted in item 1, \mathbb{R}^n may

be identified with the set of all maps x from the *index set* $\{1, 2, \dots, n\}$ to \mathbb{R} , equivalently of all tuples (x_1, \dots, x_n) of reals. Extended to the case $n = 0$, \mathbb{R}^0 is the set of maps x from the empty index set to \mathbb{R} ; the only such map has as its graph the empty set of pairs (i, x_i) . In this interpretation, \mathbb{R}^0 is the singleton set $\{\emptyset\}$ whose only member is \emptyset . In the linear algebra interpretation, \mathbb{R}^0 is the (unique up to isomorphism) 0-dimensional real linear space, whose one member is usually called 0. In the tuple interpretation, it is the set whose only member is the empty tuple $()$.

- d) A **point function** is a mathematical (Level 1) real function of real variables, $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ for some integers $n \geq 0, m > 0$.

It is a **scalar** function when $m = 1$, otherwise a **vector-valued** function. When not otherwise specified, scalar is assumed. A function with $n = 0$ and $m = 1$ —i.e., $\mathbb{R}^0 \rightarrow \mathbb{R}$ —is a **real constant**. The unique such function with empty domain is by definition the **Not a Number** constant, NaN.

- e) A **generic function** (more precisely, a “generic name” for functions) is a name that denotes any of several related functions called **versions** of the generic name. In each flavor a generic function has a Level 1 **primary version**—which may be the same in all flavors, or flavor-defined—from which all other versions are derived.

E.g., **add** or “+” denotes mathematical addition of reals (the primary version), or Level 1 interval addition in some flavor, or finite precision (Level 2) addition in some interval type.

- f) For a given flavor, a generic function is **supported**, and is an **operation**, if the flavor specifies its primary version; see 7.4. The operations in Clause 9 are so called because they are supported in all flavors.
- g) For a given implementation of a flavor, an operation is **provided** if the implementation provides one or more Level 2 versions of it; see 7.5.3. Since such a version must be derived from a primary version, a function can only be provided if it is supported, i.e., if it is an operation.
- h) An **arithmetic operation** is an operation whose primary version is a point function. This point function shall have a nonempty domain (so Level 1 NaN cannot be an arithmetic operation). Any other operation is a **non-arithmetic operation**.
- i) A **constructor** is a function that creates an interval from non-interval data (9.8).

6.2 Expression definitions

An expression is some symbolic form that can be used to define a function—in general not a static object within program code, but derived dynamically from a particular program execution. Expressions are central to interval computation, because the Fundamental Theorem of Interval Arithmetic (FTIA) is about interpreting an expression in different ways:

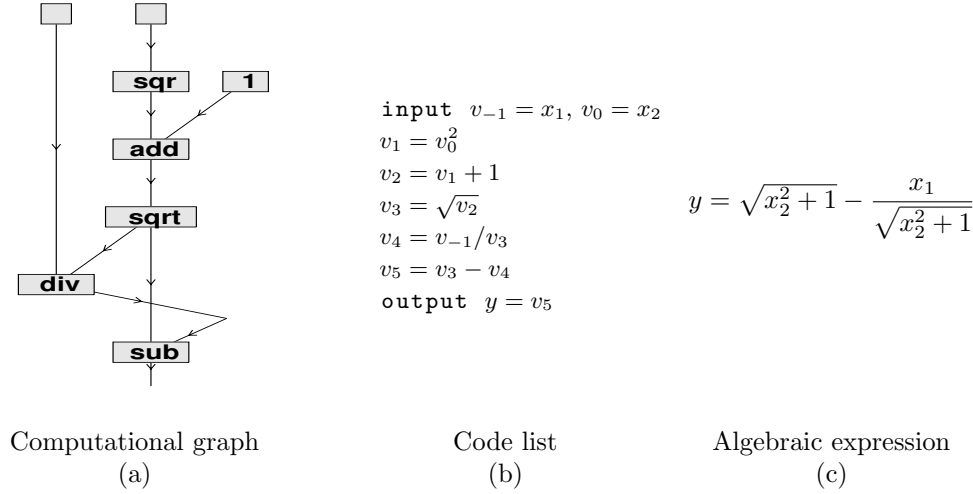


Figure 6.1. Essentially equivalent notations for an expression.

Notes: In (a), structure is shown by labeling nodes with operations only; the order of arguments is shown by reading incoming edges left to right, e.g., the inputs to **sub** are the results of the preceding **sqrt** and **div**, in that order. Similarly the input nodes are x_1 and x_2 left to right.

Form (c) has redundancy in the sense of repeated subexpressions.

Constant 1 denotes the zero-argument constant function $1()$ wherever it occurs.

- a) as defining a Level 1 point function f ;
- b) as defining various (depending on the finite precision interval types used) interval functions that give proven enclosures for the range of f over an input box \mathbf{x} ;
- c) as defining corresponding decorated interval functions that can give a stronger conclusion: e.g., in the set-based flavor, that f is everywhere defined, or everywhere continuous, on \mathbf{x} —enabling, say, an automatic check of the hypotheses of the Brouwer Fixed Point Theorem [B17].

The meaning of a) is flavor-independent; that of b) and c) is flavor-defined. The standard specifies behavior, at the individual operation level, that enables such conclusions, whether or not the notion “expression” exists in a programming language.

A *formal expression* defines a relation between certain mathematical variables—the *inputs*—and others—the *outputs*—via the application of named *operations*. It is by definition an acyclic (having an acyclic computational graph, see below) finite set of dependences between mathematical variables, defined by equations

$$v = \varphi(u_1, \dots, u_k), \quad (1)$$

where v and the u_i come from a set \mathcal{X} of *variable-symbols*; φ comes from a set \mathcal{F} of generic operations called the **library**; and distinct equations have distinct v ’s—the *single assignment* property.

An **arithmetic expression** is a formal expression, all of whose operations are arithmetic operations. Evaluating it using in turn the point version, an interval version and a decorated interval version of each of its operations, with appropriate inputs, gives the interpretations a) to c) above, to which an FTIA appropriate to the flavor can be applied.

Non-arithmetic expressions (containing non-arithmetic operations) do not allow these three interpretations, e.g., the expression $\text{mid}(\text{intersection}(\mathbf{x}, \mathbf{y}))$ in this standard can only be interpreted as a function with two interval inputs and a numeric output.

Three descriptions of an expression are outlined here. To apply the FTIA, it suffices to consider expressions that are *scalar*, with a single output.

- a) Drawing an edge from each u_i to v for each dependence-equation (1) defines the *computational graph* \mathcal{G} —Figure 6.1(a)—a directed graph over the node set \mathcal{X} . The dependences define an expression if and only if \mathcal{G} is *acyclic*. There is then a nonempty set of *output* nodes having no outgoing edge, and a possibly empty set of *input* nodes having no incoming edge.
- b) Since \mathcal{G} is acyclic, the equations can be ordered so that each one only depends on already known (input, or previously computed) values, thus representing the expression as a *code list*—Figure 6.1(b). In the notation of A. Griewank [B4], the inputs are written v_{1-n}, \dots, v_0 where $n \geq 0$, conventionally given the aliases x_1, \dots, x_n , so x_i is the same as v_{i-n} . The operations are

$$v_r = \varphi_r(u_{r,1}, \dots, u_{r,k_r}), \quad (r = 1, \dots, m),$$

where $\varphi_r \in \mathcal{F}$ with arity (number of arguments) k_r , and each $u_{r,i}$ is a known v_j , that is $j = j(r, i) < r$. (Constants, which are operations of arity 0, may be referred to directly instead of assigned to a v_j .) Without loss, the outputs can be placed last so if there are p of them they are v_{m-p+1}, \dots, v_m , aliased to y_1, \dots, y_p , so that the code list defines a (vector) formal function $(y_1, \dots, y_p) = f(x_1, \dots, x_n)$; in case $p = 1$, it is written as a scalar function $y = f(x_1, \dots, x_n)$.

Either m or n , but not both, can be zero. The case $n = 0$ and $m \geq 1$ gives a *constant expression*. For the scalar case $p = 1$, if $m = 0$ and $n = 1$ there are no operations, and y is the same as x_1 , defining the *identity function* $y = f(x_1) = x_1$; while for $p = 1$, $m = 0$ and $n \geq 1$ there are n possibilities, the *coordinate projections* $y = \pi_j(x_1, \dots, x_n) = x_j$ ($j = 1, \dots, n$).

- c) By allowing redundancy, an expression can always be converted to a normal (scalar) *algebraic expression*—Figure 6.1(c)—over the variable-set \mathcal{X} and library \mathcal{F} , defined recursively as follows:
- if $x \in \mathcal{X}$ is a variable symbol, then x is an algebraic expression;
 - if $\varphi \in \mathcal{F}$ is a function symbol of arity k and if e_i is an expression for $i = 1, \dots, k$, then the *function symbol application* $\varphi(e_1, \dots, e_k)$ is an algebraic expression.

Multiple outputs may be represented by tuples of separate algebraic expressions, e.g., the vector function $f(\theta) = (\cos(\theta), \sin(\theta))$. Redundancy may occur because this form cannot refer to a subexpression by name and must repeat it in full at each use, as with $\sqrt{x_2^2 + 1}$ in Figure 6.1(c).

The three forms are semantically equivalent, both at Level 1 and at Level 2. Because of its simple recursive definition, (c) is the form used in the FTDIA proof in Annex B.

When an expression is evaluated in interval mode, multiple instances of the same variable can lead to excessive widening of the final result: e.g., evaluating $x - x$ with an interval input \mathbf{x} gives, not $[0, 0]$, but an interval twice the width of \mathbf{x} . The question of when it is valid to manipulate expressions (e.g., to replace $x - x$ by 0) is important for interval computation because of its potential to tighten enclosures, but is outside the scope of this standard.

6.3 Function libraries

An implementation's library for a flavor includes the following.

- The *Level 1 library* comprises the Level 1 operations, including those specified in Clause 9 and possible flavor-defined operations. These are the *primary versions* defined in 7.4 and, for arithmetic operations, their flavor-defined Level 1 bare and decorated interval versions.
- The *Level 2 library* comprises all Level 2 versions of the Level 1 operations using, e.g., different bare or decorated interval types or different number formats as the implementation may provide.

For arithmetic operations there is a further conceptual grouping of operations into the *point library* of primary versions at Level 1, and at Level 2 the *bare interval library* and the *decorated interval library*, to reflect different ways in which an arithmetic expression can be evaluated.

In this standard each generic operation $\varphi(u_1, \dots, u_k)$, see eqn (1), is presented as having a fixed number k of arguments, termed its **arity**. Also for each version of φ at both Level 1 and Level 2, the **datatype** of each argument u_i , namely the set from which values of this argument are taken, is fixed.

At Level 1 the following datatypes are used by operations required in all flavors:

- *interval*, a generic term that maps to a Level 1 interval of a given flavor;
- *decoration*, a generic term that maps to a decoration of a given flavor;
- *number*, denoting a member of the reals \mathbb{R} or extended reals $\overline{\mathbb{R}}$ depending on the flavor;
- *string*, a character sequence on a language- or implementation-defined alphabet;
- *integer*, used to describe parameterized sets of operations⁷;
- *boolean*, one of **false**, **true**.

Other flavor-defined operations may use other datatypes, e.g., the 16-valued state of the Allen extended interval comparisons in 10.6.4 of the set-based flavor.

At Level 2, the decoration, string and boolean datatypes shall have the same meaning as at Level 1; each instance of the interval and number datatypes maps to one of various implementation-defined interval types and number formats, respectively.

Just as the standard is not concerned with the actual names or invocation methods, this conceptual view of libraries is independent of how the operations are presented in an actual computing environment: it could be via programming libraries, language primitives, infix operators, or other means. Also an implementation might permit the arity of an operation, or the datatype of any of its arguments, to be variable—possibly determined at run time. An implementation shall document how the formal operations are mapped to language entities.

6.4 The FTIA

The required arithmetic operations of Clause 9 are flavor-independent in the sense that the point version is the same in all flavors. Hence the point function $f(x) = f(x_1, \dots, x_n)$, defined by an arithmetic expression f made up of required operations, is the same in all flavors; in particular f has a flavor-independent *natural domain* $\text{Dom}(f)$ determined by its constituent operations: the set of points $x \in \mathbb{R}^n$ where f has a value in the sense that the whole expression can be evaluated to give a real result.

[Example. Knowing the definitions of division x/y , addition $x + y$ and square root \sqrt{x} , including that their domains are $\{(x, y) \in \mathbb{R}^2 \mid y \neq 0\}$, all of \mathbb{R}^2 , and the interval $0 \leq x < +\infty$, respectively, one may deduce the natural domain of $\sqrt{1 + 1/x}$ is the union of intervals $-\infty < x \leq -1$ and $0 < x < +\infty$.]

Each flavor has a version of Moore’s Fundamental Theorem of Interval Arithmetic (FTIA): under suitable conditions, the range of such an f over a box is contained in the result of interval-evaluating the expression over the box. Typically the theorem has sub-cases where extra conditions give stronger conclusions, as illustrated by the version for the set-based flavor:

THEOREM 6.1 (FTIA in set-based flavor). *Let $\mathbf{y} = f(\mathbf{x})$ be the result of interval-evaluation of f over a box $\mathbf{x} = (x_1, \dots, x_n)$ using any interval versions of its component library functions. Then*

- a) (“Basic” form of FTIA.) *In all cases, \mathbf{y} contains the range of f over \mathbf{x} , that is, the set of $f(x)$ at points of \mathbf{x} where it is defined:*

$$\mathbf{y} \supseteq \text{Rge}(f \mid \mathbf{x}) = \{f(x) \mid x \in \mathbf{x} \cap \text{Dom}(f)\}. \quad (2)$$

- b) (“Defined” form of FTIA.) *If also each library operation in f is everywhere defined on its inputs, while evaluating \mathbf{y} , then f is everywhere defined on \mathbf{x} , that is $\text{Dom}(f) \supseteq \mathbf{x}$.*

- c) (“Continuous” form of FTIA.) *If in addition to b), each library operation in f is everywhere continuous on its inputs, while evaluating \mathbf{y} , then f is everywhere continuous on \mathbf{x} .*

- d) (“Undefined” form of FTIA.) *If some library operation in f is nowhere defined on its inputs, while evaluating \mathbf{y} , then f is nowhere defined on \mathbf{x} , that is $\text{Dom}(f) \cap \mathbf{x} = \emptyset$.*

⁷ Examples are the standard function $\text{pown}(x, p) = (x \text{ raised to integer power } p)$, which is treated as a family of functions $\varphi_p(x)$ of real x , parameterized by p . Similarly for any function having some non-real arguments. The array length in the reduction operations of 12.12.12 counts as such an argument.

The details depend on the flavor’s underlying theory, including its meaning of “contains.” The standard is constructed so that for any conforming implementation, the flavor’s FTIA holds in finite precision, not just at Level 1.

In each flavor the decoration system, Clause 8, shall make it possible, while evaluating an arithmetic expression, to determine that some sub-case of the FTIA holds. It shall be supported by a *Fundamental Theorem of Decorated Interval Arithmetic* (FTDIA) stating that if evaluating an expression using decorated intervals returns a certain decoration d on the result, then the conditions for a corresponding FTIA sub-case are verified, hence the corresponding FTIA conclusion follows.

NOTE—The converse is false, e.g., f may be “continuous” but the evaluated d may only say “defined” or “basic.” Finding optimal decorations d and tight enclosures y requires careful algorithm design.

6.5 Related issues

When program code contains conditionals (including loops), the run time data flow and hence the computed expression generally depends on the input data—for instance, Example (b) in 11.8 where a function is defined piecewise. The user is responsible for checking that a property such as global continuity holds as intended in such cases. The standard provides no way to check this automatically.

Though the set operations **intersection** and **convexHull** are not point-operations and cannot appear directly in an arithmetic expression, they are useful for efficiently *implementing* interval extensions of functions defined piecewise, see the 11.8 example mentioned in the last paragraph.

The standard requires that at Level 2, for all interval types, all operations and all inputs other than NaI (if NaI is provided by the flavor), the interval part of a decorated interval operation equal the corresponding bare interval operation. This ensures that converting bare interval program code to use decorated intervals leaves the data flow entirely unchanged (provided no conditionals depend on decoration values, and NaI does not occur)—hence the computed expression and the interval part of its result are unchanged. If this were not so, there might in principle be an arbitrarily large discrepancy between the bare and the decorated versions of a computation that contains conditionals.

If a reproducible mode is supported, it should be supported in the same spirit as that of IEEE Std 754-2008: provided that certain programming restrictions (see 1.8 item (7)) are adhered to, the order of operations and the resulting values should be predictable based on the given input values.

7. Flavors

7.1 Flavors overview

The standard permits different interval behaviors via the flavor concept described in this clause.⁸ An **interval model** means a particular foundational approach to interval arithmetic, in the sense of a Level 1 abstract datatype of entities called intervals and of operations on them. A **flavor** is an interval model that conforms to the core specification described below.

A **provided flavor** is a flavor that the implementation provides in finite precision, see 7.5.

A **standard flavor** is one that has a specification in this standard, which extends the core specification. The set-based flavor is currently the only standard flavor.

An implementation shall provide at least one standard flavor. The implementation as a whole is conforming if each standard flavor conforms to the specification of that flavor, and each non-standard flavor conforms to the core specification.

Flavor is a property of program execution context, not of an individual interval. Therefore, just one flavor shall be in force at any point of execution. It is recommended that at the language level, the flavor should be constant over a procedure/function or a compilation unit.

⁸Any person proposing a new flavor should submit it as a revision of this standard, see the item “Comments on standards” in the front matter of this document.

An implementation with more than one flavor should provide means for conversion between suitably chosen pairs of interval types of different flavors. How this is provided—e.g., whether it is an import operation to, or an export operation from, the current flavor—is implementation-defined.

For brevity, phrases such as “A flavor shall provide, or document, a feature” mean that an implementation of that flavor shall provide the feature, or its documentation describe it.

The core specification of a flavor is summarized in the list below, and detailed in the following subclauses. The entities in the list are part of the Level 1 mathematical theory unless said otherwise.

- a) There is a flavor-independent set of *common intervals*, defined in 7.2.
- b) There is a flavor-independent set of named *common operations*, defined in Clause 9.
- c) There is a flavor-independent set of *common evaluations* of common operations, defined in 7.3. They have common intervals as input and, in the sense of 7.2, give the same result in any flavor.
- d) There is a flavor-defined set of *intervals* of the flavor that, in the sense defined in 7.2, contains the common intervals as a subset. There is a flavor-defined finite set of *decorations* as described in Clause 8; in particular it includes the `com` decoration. A decorated interval is a pair (interval, decoration).
- e) There is a flavor-defined partial order called *contains* for the intervals, see 7.2, which for common intervals coincides with normal set containment.
- f) A flavor shall define the Level 1 bare and decorated interval version of each of its arithmetic operations, and the Level 1 version (primary version) of each other operation, see 7.4.
- g) At Level 2, intervals are organized into flavor-defined *types*. An (interval) type is essentially a finite set of Level 1 intervals; see 7.5.1.
- h) The relation between a Level 1 operation and a version of it at Level 2, see 7.5.3, is summarized as follows. The latter evaluates the Level 1 operation on the Level 1 values denoted by its inputs. If (at those inputs) the operation has no value, an exception is signaled, or some default value returned, or both, in a flavor-defined way. Otherwise the returned value is converted to a Level 2 result of an appropriate Level 2 datatype. If the result is of interval type, overflow may occur in some flavors, causing an exception to be signaled.

7.2 Flavor basic properties

A flavor is described at Level 1 by a set \mathfrak{F} of entities, the **intervals** of the flavor, a set of decorations, and a set of generic operations that includes the required operations of Clause 9. The symbol \mathfrak{F} is also used to refer to the flavor as a whole.

The (flavor-independent) **common intervals** are defined to be the set \mathbb{IR} of nonempty closed bounded real intervals used in classical Moore arithmetic [B8].

The relation of \mathfrak{F} to the common intervals is described by a one-to-one **embedding map** $f: \mathbb{IR} \rightarrow \mathfrak{F}$. Usually, $f(x)$ is abbreviated to $\mathfrak{f}x$. The set of all $\mathfrak{f}x$ for $x \in \mathbb{IR}$ forms the **common intervals of \mathfrak{F}** . Usually x is identified with $\mathfrak{f}x$ and \mathbb{IR} is treated as a subset of \mathfrak{F} , for every flavor. To emphasize that this has been done, a statement may be said to hold *modulo the embedding map*.

The set of decorations of \mathfrak{F} is finite, and includes the `com` decoration; see Clause 8.

[Examples.

- A *Kaucher interval* is defined to be a pair (a, b) of real numbers—equivalently, a point in the plane \mathbb{R}^2 —which for $a \leq b$ is “proper” and identified with the normal real interval $[a, b]$, and for $a > b$ is “improper.” Thus the embedding map is $x \mapsto (\inf x, \sup x)$ for $x \in \mathbb{IR}$.
- For the set-based flavor, every common interval is actually an interval of that flavor (\mathbb{IR} is a subset of $\overline{\mathbb{IR}}$), so the embedding is the identity map $x \mapsto x$ for $x \in \mathbb{IR}$.

]

For each flavor \mathfrak{F} , a relation \supseteq , called **contains** or **encloses**, shall be defined between intervals. It shall be a partial order: $\mathbf{x} \supseteq \mathbf{y}$ and $\mathbf{y} \supseteq \mathbf{z}$ imply $\mathbf{x} \supseteq \mathbf{z}$, and $\mathbf{x} = \mathbf{y}$ if and only if both $\mathbf{x} \supseteq \mathbf{y}$ and $\mathbf{y} \supseteq \mathbf{x}$ hold. For common intervals it shall have the normal meaning modulo the embedding map: if $\mathbf{x}, \mathbf{y} \in \mathbb{IR} \subseteq \mathfrak{F}$ then $\mathbf{x} \supseteq \mathbf{y}$ in \mathfrak{F} if and only if $\mathbf{x} \supseteq \mathbf{y}$ in the sense of set containment.

[Example. In the Kaucher flavor, (a, b) is defined to contain (a', b') if and only if $a \leq a'$ and $b \geq b'$; e.g., $[1, 2] \supseteq [2, 1]$ is true. For common (proper) intervals, this coincides with normal containment.]

7.3 Common evaluations

For each Level 1 version of an operation φ for which at least one input or output is of bare interval datatype, a set of k -tuples $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_k)$ is specified for which $\varphi(\mathbf{x}_1, \dots, \mathbf{x}_k)$ shall have the same Level 1 value \mathbf{y} in all flavors. Then \mathbf{x} is a **common input**, and \mathbf{y} is the **common value** of φ at \mathbf{x} . The $(k+1)$ -tuple $(\mathbf{x}_1, \dots, \mathbf{x}_k; \mathbf{y})$ is a **common evaluation** of φ , alternatively written $\mathbf{y} = \varphi(\mathbf{x}_1, \dots, \mathbf{x}_k)$. It may be called a common evaluation **instance** to emphasize that a specific input tuple is involved.

A tuple $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_k)$ whose components \mathbf{x}_i are common intervals—nonempty closed bounded intervals in \mathbb{R} —can be regarded as a nonempty closed bounded box in \mathbb{R}^k . A tuple whose \mathbf{x}_i are common intervals of some flavor \mathfrak{F} can be identified with such a box, modulo the embedding map.

There are two cases:

- a) If φ is one of the arithmetic operations in 9.1, its point version is a real function $y = \varphi(x_1, \dots, x_k)$ of $k \geq 0$ real variables⁹. The common inputs are those boxes $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_k)$ such that each \mathbf{x}_i is common and the point version of φ is defined and continuous at each point of \mathbf{x} . The common value \mathbf{y} is the range

$$\mathbf{y} = \text{Rge}(\varphi | \mathbf{x}) = \{ \varphi(x_1, \dots, x_n) \mid x_i \in \mathbf{x}_i \text{ for each } i \}, \quad (\text{here } \varphi \text{ is the point version}). \quad (3)$$

By theorems of real analysis, \mathbf{y} is nonempty, closed, bounded and connected, so it is a common interval.

- b) If φ is one of the non-arithmetic operations in 9.2 to 9.5, it has a direct definition unrelated to a point function. Its set of common inputs comprises those $(\mathbf{x}_1, \dots, \mathbf{x}_k)$ such that each \mathbf{x}_i that is of interval datatype is a common interval, and $\mathbf{y} = \varphi(\mathbf{x}_1, \dots, \mathbf{x}_k)$ is defined and, if of interval datatype, is a common interval. The common value is \mathbf{y} .

7.4 Primary versions and Level 1 interval versions

Let \mathfrak{F} be a flavor. With terms as defined in 6.1, a primary version shall be **fully specified**, which means that the following are defined, whether in Clause 9 or by the flavor: the arity k of the function φ , namely the number of arguments; the Level 1 datatype of each argument; the **domain**, namely the set of input tuples (x_1, \dots, x_k) , with x_i of the specified datatypes, at which the operation has a value; and its value $y = \varphi(x_1, \dots, x_k)$ at each such point. Equivalently, the flavor defines the primary version's function-graph: the set of all (x_1, \dots, x_k, y) such that $y = \varphi(x_1, \dots, x_k)$.

Full specification makes it possible to test objectively whether, in a flavor \mathfrak{F} , an implemented Level 2 version of an operation encloses the Level 1 result (if an interval), approximates it (if numeric) or equals it (if boolean), whenever the Level 1 result is defined in \mathfrak{F} .

The possible cases for a generic operation φ are listed in the following subclauses.

7.4.1 Arithmetic operations

Here the primary version is a point function, whether flavor-defined, or required in all flavors according to Clause 9.

The flavor shall define a mapping that takes an arbitrary point function $y = \varphi(x_1, \dots, x_k)$ to a function $\mathbf{y} = \varphi(\mathbf{x}_1, \dots, \mathbf{x}_k)$, called the **Level 1 (bare) interval version** of φ in \mathfrak{F} , for which the \mathbf{x}_i and \mathbf{y} are Level 1 bare intervals of \mathfrak{F} . Its evaluations shall include the common evaluations of 7.3, modulo the embedding map.

⁹For the treatment of “integer power” $\text{pown}(x, p)$ and similar functions, see the footnote on p.18.

The flavor shall also define a mapping that takes such a point function to a function $Y = \varphi(X_1, \dots, X_k)$, called the **Level 1 decorated interval version** of φ in \mathfrak{F} , for which the X_i and Y are Level 1 decorated intervals of \mathfrak{F} . If $y = \varphi(x_1, \dots, x_k)$ is a common evaluation of the bare interval version, then $y_{\text{com}} = \varphi((x_1)_{\text{com}}, \dots, (x_k)_{\text{com}})$ shall be an evaluation of the decorated interval \mathfrak{F} -version. It is termed a **Level 1 decorated common evaluation** of φ , and $((x_1)_{\text{com}}, \dots, (x_k)_{\text{com}})$ is a **Level 1 decorated common input**. These mappings from point function to Level 1 bare and decorated interval version shall ensure the validity of suitable flavor-defined Fundamental Theorems of interval arithmetic, and of decorated interval arithmetic, respectively.

The Level 1 interval versions shall be fully specified but may have no value for some combinations of input intervals.

[Example. In the set-based flavor, the Level 1 bare interval version of φ is the natural interval extension of the point function, and is defined for all combinations of input intervals. However in classical Moore arithmetic and in Kaucher arithmetic, the Level 1 interval version of division x/y , for instance, is not defined when $0 \in y$.]

7.4.2 Nonarithmetic operations required in all flavors

For a nonarithmetic operation required in all flavors according to Clause 9, the primary version in \mathfrak{F} shall be a Level 1 function for which each interval input or output becomes an interval of \mathfrak{F} , and whose evaluations include the common evaluations specified in Clause 9, modulo the embedding map.

[Example. The interval constructors `numsToInterval` and `textToInterval` are required in all flavors. A flavor will in general extend these beyond the common evaluations. E.g., in the set-based flavor `numsToInterval(l, u)` constructs unbounded intervals by allowing non-common inputs with $l = -\infty$ and/or $u = +\infty$; in a Kaucher flavor it might construct improper intervals by allowing $l > u$.]

7.4.3 Flavor-defined nonarithmetic operations

Besides the operations required in all flavors, a flavor may define *flavor-specific* operations that are required in each implementation of that flavor. Each such operation shall be fully specified at Level 1.

[Examples.

- *The set-based flavor defines reduction operations, one of which is the dot-product. The primary version of this is $s = \text{dot}(x, y) = \sum_{i=1}^n x_i y_i$ where x, y are two real vectors of length n . Though this may be regarded as a point function (or a family of them, parameterized by n), it is not an arithmetic operation because it is not intended to have an interval version.*
- *A flavor might define a decorated interval version of a non-arithmetic operation.*
- *A flavor might define an operation on decorations, adapted to its decoration model.*
- *A flavor might provide other constructors besides the required ones, or give extra inputs to the required constructors. E.g., in a flavor that provides open and half-open as well as closed intervals, optional extra input(s) to `numsToInterval` might specify which bounds are open.*

]

NOTE—Apart from the decorated common evaluations the relation between the bare and decorated interval versions of an operation φ is flavor-defined, but barring special cases such as “Not an Interval” input, the interval part of the latter should equal the former. That is, if φ is defined at decorated interval inputs (X_1, \dots, X_k) , with decorated interval value Y , and if the interval parts of X_i and Y are x_i and y , then $\varphi(x_1, \dots, x_k)$ should in most cases be defined and have value y .

7.5 The relation of Level 1 to Level 2

Let \mathfrak{F} be a flavor.

7.5.1 Types

At Level 2, bare intervals of \mathfrak{F} shall be organized into finite sets called (bare interval) **types**. A type is an abstraction of a particular way to represent intervals. What types are provided by an implementation is both flavor-defined and language- or implementation-defined.

There shall be a one-to-one correspondence between bare interval types and decorated interval types (*bare types* and *decorated types* for short). An element of a decorated type is a pair¹⁰ $\mathbf{X} = (\mathbf{x}, dx)$ where \mathbf{x} and dx —the *interval-part* and the *decoration-part* of \mathbf{X} —belong to the corresponding bare type, and to the finite set of decorations of the flavor, respectively.

An implementation shall ensure that the decoration `com` is applied only to common intervals; a common interval may however have a decoration other than `com`.

Level 2 entities are called **datums**. A type may contain some exceptional datums, e.g., a “Not an Interval” datum, besides those that denote intervals. A **T-datum** means a general member of a type \mathbb{T} ; a **T-interval** means a \mathbb{T} -datum that denotes an interval. If the type has no exceptional datums, these terms are synonymous.

Each Level 2 bare interval shall denote a unique Level 1 interval \mathbf{x} of \mathfrak{F} and is normally regarded as being that interval. However \mathfrak{F} may formally define it as a pair (\mathbf{x}, t) where t (the type name) is a symbol that uniquely identifies the type, thus making datums of different types different even if they denote the same Level 1 interval.

Hence, each Level 2 decorated interval denotes a unique Level 1 decorated interval of \mathfrak{F} but may formally be a triple (\mathbf{x}, dx, t) of interval, decoration and type name.

7.5.2 Hull

Each bare type \mathbb{T} has a **T-hull** operation. At Level 1 it shall be defined by an algorithm that maps an arbitrary interval \mathbf{x} of \mathfrak{F} to a minimal \mathbb{T} -interval \mathbf{y} such that $\mathbf{y} \supseteq \mathbf{x}$ in the flavor’s “contains” order, if such a \mathbf{y} exists, and otherwise is undefined. Minimal means that if \mathbf{z} is another \mathbb{T} -interval and $\mathbf{y} \supseteq \mathbf{z} \supseteq \mathbf{x}$ then $\mathbf{y} = \mathbf{z}$.

At Level 2, an implementation should provide the \mathbb{T} -hull for each supported bare type \mathbb{T} , as an operation `convertType` that maps an arbitrary interval of any other supported bare type of \mathfrak{F} to its \mathbb{T} -hull if this exists, and signals the flavor-independent `Intv1Overflow` exception otherwise. Its action on non-interval datums is flavor-defined.

7.5.3 Level 2 operations

For each bare or decorated interval version of an arithmetic operation in 9.1 and for each operation in the rest of Clause 9, a \mathbb{T} -version of φ shall be provided for each bare interval type \mathbb{T} of the implementation of \mathfrak{F} . For flavor-defined operations the flavor shall specify whether an implementation shall or should provide a \mathbb{T} -version.

When passing from a Level 1 operation to a \mathbb{T} -version—whether required in all flavors or flavor-specific:

- Inputs or output of bare interval datatype change to type \mathbb{T} .
- Those of decorated interval datatype change to the decorated type of \mathbb{T} .
- Those of real datatype change to a type-dependent Level 2 number format \mathbb{F} . The flavor should make recommendations or requirements on the relation between \mathbb{F} and \mathbb{T} .

In the description below, sometimes a *flavor-defined Level 2 value* may be returned in cases where no Level 1 value exists, instead of or alongside signaling an exception.

¹⁰ For readability, the decorations attached to intervals $\mathbf{u}, \mathbf{v}, \dots$ are often named du, dv, \dots , and each resulting decorated interval named by the corresponding uppercase letter, e.g., $\mathbf{U} = (\mathbf{u}, du) = \mathbf{u}_{du}$, etc. The d ’s here have no connection with differentials in calculus.

[Example. For the `mid()` function at Level 2 in the set-based flavor, returning the midpoint of `Empty` as `NaN`, and of `Entire` as 0, illustrate flavor-defined values. Both values are undefined at Level 1. It was considered that no numeric value of `mid()` makes sense, but that some algorithms are simplified by returning a default value 0 for `mid(Entire)`.]

For some operations and inputs the implementation might be unable to effectively compute some value or determine whether some statement is true or false—e.g., owing to constraints on time, space or algorithmic complexity. Below, the term *is found* means that the implementation is able to compute such a value or determine such a fact; *is not found* means the opposite.

[Examples.

– A constructor in the set-based flavor needs to ensure $l \leq u$ when forming an interval $[l, u]$. This may be hard to check, for some implementations and types.

– In a flavor where it is a fatal error to evaluate an arithmetic operation outside its domain, it may be hard to decide if the required operation $\varphi(x) = \tan(x)$ is defined on $x = [\underline{x}, \bar{x}]$ when \underline{x}, \bar{x} are large and differ by less than π . Let format \mathbb{F} be IEEE 754 `binary64` and \mathbb{T} the inf-sup type based on \mathbb{F} . The integer $a = 214112296674652$ differs from $136308121570117\pi/2$ by less than $2.6\text{e-}16$; and $a - 1, a, a + 1$ are exact \mathbb{F} -numbers. One of the \mathbb{T} -intervals $[a - 1, a]$ and $[a, a + 1]$ (but not both) contains a singularity of φ ; but which?]

Evaluation of a \mathbb{T} -version $\varphi_{\mathbb{T}}$ of an operation φ shall be as if the following is done. Let $X = (X_1, \dots, X_k)$ be the tuple of Level 1 values denoted by the inputs to $\varphi_{\mathbb{T}}$.

- a) If the Level 1 operation is found to be undefined at X then $\varphi_{\mathbb{T}}$ signals the flavor-independent `UndefinedOperation` exception, or returns a flavor-defined value, or both. This includes the case where some input has no Level 1 value, e.g., a numeric input is `NaN`.
- b) If it is not found whether the Level 1 operation is defined at X then $\varphi_{\mathbb{T}}$ signals the flavor-independent `PossiblyUndefinedOperation` exception, or returns a flavor-defined value, or both.
- c) Otherwise, it is found that the Level 1 value $Y = \varphi(X_1, \dots, X_k)$ is defined in \mathfrak{F} .
 - 1) If Y is a bare interval \mathbf{y} there are two cases.
 - If a \mathbb{T} -interval \mathbf{z} containing \mathbf{y} is found (in particular if `Entire` exists in \mathfrak{F} and is a \mathbb{T} -interval), then $\varphi_{\mathbb{T}}$ returns such a \mathbf{z} .
 - Otherwise, $\varphi_{\mathbb{T}}$ signals the flavor-independent `IntvlOverflow` exception and the returned result, if any, is flavor-defined.
 - 2) If Y is a decorated interval \mathbf{y}_{dy} there are three cases.
 - If φ is an arithmetic operation, and (X_1, \dots, X_k) is a decorated common input (this implies $dy = \text{com}$, see 7.4.1), and a common \mathbb{T} -interval \mathbf{z} containing \mathbf{y} is found, then $\varphi_{\mathbb{T}}$ returns such a \mathbf{z} with the decoration `com`.
 - Otherwise, if a \mathbb{T} -interval \mathbf{z} containing \mathbf{y} is found (in particular if `Entire` exists in \mathfrak{F} and is a \mathbb{T} -interval), then $\varphi_{\mathbb{T}}$ returns such a \mathbf{z} with a flavor-defined decoration.
 - Otherwise, no such \mathbf{z} is found. Then $\varphi_{\mathbb{T}}$ signals the `IntvlOverflow` exception and the returned result, if any, is flavor-defined.
 - 3) If Y is numeric, $\varphi_{\mathbb{T}}$ returns a value of an appropriate number format, approximating Y in a type- and operation-dependent way.
 - 4) If Y is boolean, $\varphi_{\mathbb{T}}$ returns Y .

7.5.4 Measures of accuracy

Two **accuracy modes** for an interval-valued operation are defined for all flavors. An accuracy mode is in the first instance a property of an individual evaluation of an operation φ , over an input box \mathbf{x} , returning a result of type \mathbb{T} . If the evaluation does not have an interval value at Level 1, its accuracy mode is undefined. The basic property of enclosure in the sense of the flavor, defined in 7.5.3 and required for conformance, is

called *valid* accuracy mode. The property that the result equals the T-hull of the Level 1 value is called *tightest* accuracy mode.

A flavor may define other accuracy modes, and may make requirements on the accuracy achieved by an implementation. To simplify the user interface, modes should be linearly ordered by strength, with *tightest* the strongest and *valid* the weakest, where mode M is stronger than mode M' if M implies M' .

NOTE—Flavors and the Fundamental Theorem. For a common evaluation of an arithmetic expression, each library operation is, modulo the embedding map, defined and continuous on its inputs so that it satisfies the conditions of the strongest, “continuous” form of the FTIA in Theorem 6.1. At Level 1, the range enclosure obtained by a common evaluation is the same, independent of flavor.

It is possible in principle for an implementation to make this true also at Level 2 by providing shared number formats and interval types that represent the same sets of reals or intervals in each flavor; and library operations on these types and formats that have identical numerical behavior in each flavor. For example, both set-based and Kaucher flavors might use intervals stored as two IEEE 754 **binary64** numbers representing the lower and upper bounds, and might ensure that operations, when applied to the intervals recognized by both flavors, behave identically. Such shared behavior might be useful for testing correctness of an implementation.

Beyond common evaluations, versions of the FTIA in different flavors can be strictly incomparable. For example, the set-based FTIA handles unbounded intervals, which the classical Kaucher flavor does not; conversely, Kaucher intervals have an extended FTIA applicable to reversed-bound intervals, which has no simple interpretation in the set-based flavor.

8. Decoration system

8.1 Decorations overview

A decoration is information attached to an interval; the combination is called a decorated interval. Interval calculation has two main objectives:

- obtaining correct range enclosures for real-valued functions of real variables;
- verifying the assumptions of existence, uniqueness, or nonexistence theorems.

Traditional interval analysis targets the first objective; the decoration system, as defined in this standard, targets the second.

A decoration primarily describes a property, not of the interval it is attached to, but of the function defined by some code that produced the interval by evaluating over some input box.

The function f is assumed to be represented by an expression in the sense of Clause 6. For instance, if a section of code defines the expression $\sqrt{y^2 - 1} + xy$, then decorated-interval evaluation of this code with suitably initialized input intervals \mathbf{x}, \mathbf{y} gives information about the definedness, continuity, etc. of the point function $f(x, y) = \sqrt{y^2 - 1} + xy$ over the box (\mathbf{x}, \mathbf{y}) in the plane.

The decoration system is designed in a way that naive users of interval arithmetic do not notice anything about decorations, unless they inquire explicitly about their values. For example, in the set-based flavor, they only need

- call the **newDec** operation on the inputs of any function evaluation used to invoke an existence theorem,
- explicitly convert relevant floating-point constants (but not integer parameters such as the p in **pown** $(x, p) = x^p$) to intervals,

and have the full rigor of interval calculations available. A smart implementation might even relieve users from these tasks. Expert users can inspect, set and modify decorations to improve code efficiency, but are responsible for checking that computations done in this way remain rigorously valid.

Especially in the set-based flavor, decorations are based on the desire that, from an interval evaluation of a real function f on a box \mathbf{x} , one should get not only a range enclosure $f(\mathbf{x})$ but also a guarantee that the pair (f, \mathbf{x}) has certain important properties, such as $f(x)$ being defined for all $x \in \mathbf{x}$, f restricted to \mathbf{x} being

continuous, etc. This goal is achieved, in parts of a program that require it, by performing *decorated interval evaluation*, whose semantics is summarized as follows:

Each intermediate step of the original computation depends on some or all of the inputs, so it can be viewed as an intermediate function of these inputs. The interval obtained on each intermediate step is an enclosure for the range of the corresponding intermediate function. The decoration attached to this intermediate interval reflects the available knowledge about whether this intermediate function is everywhere defined, continuous, bounded, etc., on the given inputs.

In some flavors, certain interval operations ignore decorations, i.e., give undecorated interval output. Users are responsible for the appropriate propagation of decorations by these operations.

This standard's decoration system—in contrast with IEEE 754's way of reporting on the outcome of an operation—has no status flags. It provides a fully local handling of exceptional conditions in interval calculations—important in a concurrent computing environment. A general aim, as in IEEE 754's use of NaN and flags, is not to interrupt the flow of computation: rather, to collate information while evaluating f , that can be inspected afterwards.

The following language- or implementation-defined features may be provided: (i) status flags that are raised in the event of certain decoration values being produced by an operation; (ii) means for the user to specify that such an event signals an exception, and to invoke a system- or user-defined handler as a result. *[Example. The user might be able to specify that execution be terminated if an arithmetic operation is evaluated on a box that is not wholly inside its domain—an interval version of IEEE 754's "invalid operation" exception.]*

8.2 Decoration definition and propagation

Each flavor shall document its set of provided decorations, their mathematical definitions, and their associated Fundamental Theorems, the FTIA and the FTDIA, as described in 6.4. These are flavor-defined, except for the decoration `com`, see 8.3.

The implementation makes the decoration system of each flavor available to the user via *decorated interval extensions* of relevant library operations. Such an operation φ , with interval inputs $\mathbf{x}_1, \dots, \mathbf{x}_k$ carrying decorations¹¹ dx_1, \dots, dx_k , shall compute the same interval output \mathbf{y} as the corresponding bare interval extension of φ —hence, dependent on the \mathbf{x}_i but not on the dx_i . It shall compute a *local decoration* d , dependent on the \mathbf{x}_i and possibly on \mathbf{y} , but not on the dx_i . It shall combine d with the dx_i by a flavor-defined *propagation rule* to give an output decoration dy , and return \mathbf{y} decorated by dy .

The local decoration d might convey purely Level 1 information—e.g., that φ is everywhere continuous on the box $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_k)$. It might convey Level 2 information related to the particular finite-precision interval types being used—e.g., that \mathbf{y} , though mathematically a bounded interval, became unbounded by overflow. For diagnostic use it might convey Level 3 or 4 information, e.g., how an interval is represented, or how memory is used.

If f is an expression, decorated interval evaluation of an expression means evaluation of f with decorated interval inputs and using decorated interval extensions of the expression's library operations. Those inputs generally need to be given initial decorations that lead to the most informative output-decoration. A flavor should provide a function that gives this initial decoration to a bare interval.

It is the responsibility of each flavor to document the meaning of its decorations and the correct use of these decorations within programs.

8.3 Recognizing common evaluation

A flavor shall provide the decoration `com`, with the following propagation rule for library arithmetic operations.

Here φ denotes a Level 2 version of a point library arithmetic operation, see 7.5.3.

If each input to φ , and the result, is common, and if each input is decorated `com`, then the result shall be decorated `com`.

¹¹See footnote 10 on the \mathbf{u}, du notation.

NOTE—`com` is the only decoration that shall have the same meaning in all flavors.

Informally, it records that the individual operation φ took bounded nonempty input intervals and produced a bounded (necessarily nonempty) output interval. This can be interpreted as “overflow did not occur.” The propagation rule ensures that if the initial inputs to an arithmetic expression f are common, and are initially decorated `com`, then the final result $y = f(x_1, \dots, x_k)$ is decorated `com` if and only if the evaluation of the whole expression was common as defined in 7.5.3.

[Examples. Reasons why an individual evaluation of φ with common inputs $x = (x_1, \dots, x_k)$ might not return `com` include the following.

- Outside domain: The implementation finds φ is not defined and continuous everywhere on x . Examples: $\sqrt{[-4, 4]}$, $\text{sign}([0, 2])$.
- Overflow: The Level 1 result is too large to be represented. Example: Consider an interval type \mathbb{T} whose intervals are represented by their lower and upper bounds in some floating-point format, let `REALMAX` be the largest finite number in that format, and x be the common \mathbb{T} -datum $[0, \text{REALMAX}]$. Then $x + x$ cannot be enclosed in a common \mathbb{T} -datum.
- Cost: It is too expensive to determine whether the result can be represented. A possible example is $\tan([a, b])$ where $[a, b]$ is of a high-precision interval type, and one of its bounds happens to be very close to a singularity of $\tan(x)$.

]

9. Operations and related items required in all flavors

This clause defines the required library arithmetic and non-arithmetic operations, see 6.2, of the standard. An implementation shall provide each required operation in each provided flavor. It also defines interval literals, which are operands of the `textToInterval` operation.

For each arithmetic operation, this clause gives the mathematical formula for the point function, its domain of definition, and the set where it is continuous if different from the domain. The common evaluations of these operations are defined by this information according to case a) in 7.3, and are not described for each individual operation.

For each non-arithmetic operation, the common evaluations are described explicitly.

The behavior of each operation outside the set of common evaluations is flavor-defined, subject to the general flavor requirements in 7.4, 7.5.

9.1 Arithmetic operations

Table 9.1 lists required arithmetic operations, including those normally written in function notation $f(x, y, \dots)$ and those normally written in unary or binary operator notation, $\bullet x$ or $x \bullet y$.

Each listed function is continuous at each point of its domain, except where stated in the Notes. Square and round brackets are used to include or exclude an interval bound, e.g., $(-\pi, \pi]$ denotes $\{x \in \mathbb{R} \mid -\pi < x \leq \pi\}$.

NOTE—The list includes: all general-computational operations in 5.4 of IEEE Std 754-2008 except `convertFromInt`; and some recommended functions in 9.2 of IEEE Std 754-2008.

Table 9.1. Required forward elementary functions.

Name	Definition	Point function domain	Point function range	Table Footnotes
<i>Basic operations</i>				
neg(x)	$-x$	\mathbb{R}	\mathbb{R}	
add(x, y)	$x + y$	\mathbb{R}^2	\mathbb{R}	
sub(x, y)	$x - y$	\mathbb{R}^2	\mathbb{R}	
mul(x, y)	xy	\mathbb{R}^2	\mathbb{R}	
div(x, y)	x/y	$\mathbb{R}^2 \setminus \{y = 0\}$	\mathbb{R}	a
recip(x)	$1/x$	$\mathbb{R} \setminus \{0\}$	$\mathbb{R} \setminus \{0\}$	
sqr(x)	x^2	\mathbb{R}	$[0, \infty)$	
sqrt(x)	\sqrt{x}	$[0, \infty)$	$[0, \infty)$	
fma(x, y, z)	$(x \times y) + z$	\mathbb{R}^3	\mathbb{R}	
<i>Power functions</i>				
pown(x, p)	$x^p, p \in \mathbb{Z}$	$\begin{cases} \mathbb{R} & \text{if } p \geq 0 \\ \mathbb{R} \setminus \{0\} & \text{if } p < 0 \end{cases}$	$\begin{cases} \mathbb{R} & \text{if } p > 0 \text{ odd} \\ [0, \infty) & \text{if } p > 0 \text{ even} \\ \{1\} & \text{if } p = 0 \\ \mathbb{R} \setminus \{0\} & \text{if } p < 0 \text{ odd} \\ (0, \infty) & \text{if } p < 0 \text{ even} \end{cases}$	b
pow(x, y)	x^y	$\{x > 0\} \cup \{x = 0, y > 0\}$	$[0, \infty)$	a, c
exp, exp2, exp10(x)	b^x	\mathbb{R}	$(0, \infty)$	d
log, log2, log10(x)	$\log_b x$	$(0, \infty)$	\mathbb{R}	d
<i>Trigonometric/hyperbolic</i>				
sin(x)		\mathbb{R}	$[-1, 1]$	
cos(x)		\mathbb{R}	$[-1, 1]$	
tan(x)		$\mathbb{R} \setminus \{(k + \frac{1}{2})\pi k \in \mathbb{Z}\}$	\mathbb{R}	
asin(x)		$[-1, 1]$	$[-\pi/2, \pi/2]$	e
acos(x)		$[-1, 1]$	$[0, \pi]$	e
atan(x)		\mathbb{R}	$(-\pi/2, \pi/2)$	e
atan2(y, x)		$\mathbb{R}^2 \setminus \{(0, 0)\}$	$(-\pi, \pi]$	e, f, g
sinh(x)		\mathbb{R}	\mathbb{R}	
cosh(x)		\mathbb{R}	$[1, \infty)$	
tanh(x)		\mathbb{R}	$(-1, 1)$	
asinh(x)		\mathbb{R}	\mathbb{R}	
acosh(x)		$[1, \infty)$	$[0, \infty)$	
atanh(x)		$(-1, 1)$	\mathbb{R}	
<i>Integer functions</i>				
sign(x)		\mathbb{R}	$\{-1, 0, 1\}$	h
ceil(x)		\mathbb{R}	\mathbb{Z}	i
floor(x)		\mathbb{R}	\mathbb{Z}	i
trunc(x)		\mathbb{R}	\mathbb{Z}	i
roundTiesToEven(x)		\mathbb{R}	\mathbb{Z}	j
roundTiesToAway(x)		\mathbb{R}	\mathbb{Z}	j
<i>Absmax functions</i>				
abs(x)	$ x $	\mathbb{R}	$[0, \infty)$	
min(x, y)		\mathbb{R}^2	\mathbb{R}	k
max(x, y)		\mathbb{R}^2	\mathbb{R}	k

Footnotes to Table 9.1

- a. In describing the domain, notation such as $\{y = 0\}$ is short for $\{(x, y) \in \mathbb{R}^2 \mid y = 0\}$, etc.
- b. Regarded as a family of functions of one real variable x , parameterized by the integer argument p .
- c. Defined as $e^{y \ln x}$ for real $x > 0$ and all real y , and 0 for $x = 0$ and $y > 0$, else has no value. It is continuous at each point of its domain, including the positive y axis which is on the boundary of the domain.
- d. $b = e, 2$ or 10 , respectively.
- e. The ranges shown are the mathematical range of the point function. To ensure containment, an interval result may include values outside the mathematical range.
- f. $\text{atan2}(y, x)$ is the principal value of the argument (polar angle) of (x, y) in the plane. It is discontinuous on the half-line $y = 0, x < 0$ contained within its domain.
- g. To avoid confusion with notation for open intervals, in this table coordinates in \mathbb{R}^2 are delimited by angle brackets $\langle \rangle$.
- h. $\text{sign}(x)$ is -1 if $x < 0$; 0 if $x = 0$; and 1 if $x > 0$. It is discontinuous at 0 in its domain.
- i. $\text{ceil}(x)$ is the smallest integer $\geq x$. $\text{floor}(x)$ is the largest integer $\leq x$. $\text{trunc}(x)$ is the nearest integer to x in the direction of zero. ceil and floor are discontinuous at each integer. trunc is discontinuous at each nonzero integer. (As defined in 7.12.9 of the C standard [B3].)
- j. $\text{roundTiesToEven}(x)$, $\text{roundTiesToAway}(x)$ are the nearest integer to x , with ties rounded to the even integer or away from zero, respectively. They are discontinuous at each $x = n + \frac{1}{2}$ where n is an integer. (As defined in 7.12.9 of the C standard [B3].)
- k. Smallest, or largest, of its real arguments.

9.2 Cancellative addition and subtraction

For common intervals $\mathbf{x} = [\underline{x}, \bar{x}]$, $\mathbf{y} = [\underline{y}, \bar{y}]$, the operation $\text{cancelMinus}(\mathbf{x}, \mathbf{y})$ is defined if and only if the width of \mathbf{x} is not less than that of \mathbf{y} , i.e., $\bar{x} - \underline{x} \geq \bar{y} - \underline{y}$, and is then the unique interval \mathbf{z} such that $\mathbf{y} + \mathbf{z} = \mathbf{x}$, with formula $\mathbf{z} = [\underline{x} - \underline{y}, \bar{x} - \bar{y}]$. The operation $\text{cancelPlus}(\mathbf{x}, \mathbf{y})$ is equivalent to $\text{cancelMinus}(\mathbf{x}, -\mathbf{y})$.

9.3 Set operations

For common intervals $\mathbf{x} = [\underline{x}, \bar{x}]$, $\mathbf{y} = [\underline{y}, \bar{y}]$:

- $\text{intersection}(\mathbf{x}, \mathbf{y})$ is the intersection $\mathbf{x} \cap \mathbf{y}$ if this is nonempty, with formula $[\max(\underline{x}, \underline{y}), \min(\bar{x}, \bar{y})]$. If the intersection is empty, no common value is defined.
- $\text{convexHull}(\mathbf{x}, \mathbf{y})$ is the tightest interval containing \mathbf{x} and \mathbf{y} , with formula $[\min(\underline{x}, \underline{y}), \max(\bar{x}, \bar{y})]$.

9.4 Numeric functions of intervals

The operations in Table 9.2 are defined for all common intervals, with the formula shown.

Table 9.2. Required numeric functions of intervals.

Name	Definition
$\text{inf}(\mathbf{x})$	\underline{x}
$\text{sup}(\mathbf{x})$	\bar{x}
$\text{mid}(\mathbf{x})$	$(\underline{x} + \bar{x})/2$
$\text{wid}(\mathbf{x})$	$\bar{x} - \underline{x}$
$\text{rad}(\mathbf{x})$	$(\bar{x} - \underline{x})/2$
$\text{mag}(\mathbf{x})$	$\sup\{ x \mid x \in \mathbf{x}\} = \max(\underline{x} , \bar{x})$
$\text{mig}(\mathbf{x})$	$\inf\{ x \mid x \in \mathbf{x}\} = \begin{cases} \min(\underline{x} , \bar{x}) & \text{if } \underline{x}, \bar{x} \text{ have the same sign} \\ 0 & \text{otherwise} \end{cases}$

9.5 Boolean functions of intervals

The comparison relations in Table 9.3 have a boolean ($1 = \text{true}$, $0 = \text{false}$) result.

9.6 Operations on/with decorations

The function newDec adds a decoration to a bare interval \mathbf{x} :

$$\text{newDec}(\mathbf{x}) = \mathbf{x}_d$$

Table 9.3. Comparisons for intervals a and b . Notation \forall_a means “for all a in a ,” and so on. Column 4 gives formulae when $a=[\underline{a}, \bar{a}]$ and $b=[\underline{b}, \bar{b}]$ are common.

Name	Symbol	Defining predicate	Common a, b	Description
<code>equal(a, b)</code>	$a = b$	$\forall_a \exists_b a = b \wedge \forall_b \exists_a b = a$	$\underline{a} = \underline{b} \wedge \bar{a} = \bar{b}$	a equals b
<code>subset(a, b)</code>	$a \subseteq b$	$\forall_a \exists_b a = b$	$\underline{b} \leq \underline{a} \wedge \bar{a} \leq \bar{b}$	a is a subset of b
<code>interior(a, b)</code>	$a \Subset b$	$\forall_a \exists_b a < b \wedge \forall_a \exists_b b < a$	$\underline{b} < \underline{a} \wedge \bar{a} < \bar{b}$	a is interior to b
<code>disjoint(a, b)</code>	$a \not\cap b$	$\forall_a \forall_b a \neq b$	$\bar{a} < \underline{b} \vee \bar{b} < \underline{a}$	a and b are disjoint

where d depends on x and the flavor, such that the result is suitable input for decorated-interval evaluation of an expression in that flavor.

For a decorated interval x_d , the operations `intervalPart(x_d)` and `decorationPart(x_d)` have value x and d , respectively.

9.7 All-flavor interval and number literals

This subclause defines a flavor-independent syntax of literals for common intervals, which may be extended by a flavor to include its non-common intervals.

9.7.1 Overview

An **interval literal** of a flavor is a (text) string that denotes a Level 1 interval of the flavor. It is a **bare interval literal** or a **decorated interval literal** according as it denotes a bare or a decorated interval. A **number literal** is a string that denotes an extended-real number; a (decimal) **integer literal** is a particular case. A **decoration literal** is an alphanumeric string that denotes a decoration; these shall be in one to one correspondence with the decorations of the flavor. The string `com` shall denote the decoration `com` in all flavors.

Bare and decorated interval literals are used as input to bare and decorated versions of `textToInterval` in 9.8. In this standard, number literals are only used within interval literals. The definitions of literals are not intended to constrain the syntax and semantics that a language might use to denote numbers and intervals in other contexts.

The value of an interval literal is a bare or decorated Level 1 interval x . Level 2 operations with interval literal inputs are evaluated following 7.5.3; typically they return the T-hull of x for some interval type T.

[Example. The interval denoted by the bare literal `[1.2345]` is the Level 1 single-point bare interval $x = [1.2345, 1.2345]$. However, the result of `T-textToInterval("1.2345")`, where T is the IEEE 754 `infsup binary64` type of the set-based flavor, is the interval, approximately `[1.2344999999999999, 1.2345000000000002]`, whose bounds are the nearest `binary64` numbers on either side of 1.2345.]

This subclause defines **all-flavor** number and interval literals. An all-flavor literal shall have the value defined here (modulo the embedding map if it is an interval literal) in all flavors. An all-flavor interval literal denotes a common interval; if decorated it has the decoration `com`.

A flavor may define a set of literals, including the all-flavor literals, that shall have the same value in each implementation of the flavor; these are called literals **of the flavor**. (E.g., `inf` and `[1,inf]` are literals of the set-based flavor that shall have the values $+\infty$ and $[1, +\infty]$ in each implementation.)

An implementation may support an extended form of literals, e.g., using number literals in the syntax of the host language of the implementation. It may restrict the support of literals at Level 2, by relaxing conversion accuracy of hard cases: rational number literals, long strings, etc. What extensions and restrictions of this kind are permitted is flavor-defined.

The case of alphabetic characters in interval and number literals is ignored (e.g., `[1,1e3]_com` is equivalent to `[1,1E3]_COM`.) By default number syntax shall be that of the default locale (C locale); locale-specific variants may be provided.

9.7.2 All-flavor number literals

A sign is a plus sign + or a minus sign -. An integer literal comprises an optional sign and (i.e., followed by) a nonempty sequence of decimal digits, with the usual integer value. It is a natural-number literal if the sign is absent. A positive-natural literal is a natural-number literal whose value is not zero.

An all-flavor number literal denotes a real number. It has one of the following forms.

- a) A decimal number. This comprises an optional sign, a nonempty sequence of decimal digits optionally containing a point, and an optional exponent field comprising **e** and an integer literal exponent. The value of a decimal number is the value of the sequence of decimal digits with optional point multiplied by ten raised to the power of the value of the exponent, negated if there is a leading minus sign.
- b) A number in the hexadecimal-floating-constant form of the C99 standard (ISO/IEC9899, N1256 (6.4.4.2)), equivalently hexadecimal-significand form of IEEE Std 754-2008 (5.12.3). This comprises an optional sign, the string 0x, a nonempty sequence of hexadecimal digits optionally containing a point, and an exponent field comprising **p** and an integer literal exponent. The value of a hexadecimal number is the value of the sequence of hexadecimal digits with optional point multiplied by two raised to the power of the value of the exponent, negated if there is a leading minus sign.
- c) A rational literal p/q . This comprises an integer literal p , the / character, and a positive-natural literal q . Its value is the value of p divided by the value of q .

9.7.3 Unit in last place

The “uncertain form” of interval literal, below, uses the notion of the *unit in the last place* of a number literal s of some radix b , possibly containing a point but without an exponent field. Ignoring the sign and any radix-specifying code (such as 0x for hexadecimal), s is a nonempty sequence of radix- b digits optionally containing a point. Its *last place* is the integer $p = -d$ where $d = 0$ if s contains no point, otherwise d is the number of digits after the point. Then $\text{ulp}(s)$ is defined to equal b^p . When context makes clear, “ x ulps of s ” or just “ x ulps” means $x \times \text{ulp}(s)$. [Example. For the decimal strings 123 and 123., as well as 0 and 0., the last place is 0 and one ulp is 1. For .123 and 0.123, as well as .000 and 0.000, the last place is -3 and one ulp is 0.001.]

9.7.4 All-flavor bare interval literals

An all-flavor bare interval literal has one of the following forms. To simplify stating the needed constraints, e.g., $l \leq u$, the number literals l, u, m, r are identified with their values.

- a) Inf-sup form: A string $[l, u]$ where l and u are all-flavor number literals with $l \leq u$. Its common value is the common interval $[l, u]$. A string $[m]$ with all-flavor number literal m is equivalent to $[m, m]$.
- b) Uncertain form: a string $m ? r v E$ where: m is a decimal number literal of form a) in 9.7.2, without exponent; r is empty or is a natural-number literal *ulp-count*; v is empty or is a *direction character*, either **u** (up) or **d** (down); and E is empty or is an *exponent field* comprising the character **e** followed by an integer literal *exponent* e . No whitespace is permitted within the string.

With ulp meaning $\text{ulp}(m)$, the literal $m?$ by itself denotes m with a symmetrical uncertainty of half an ulp, that is the interval $[m - \frac{1}{2}\text{ulp}, m + \frac{1}{2}\text{ulp}]$. The literal $m?r$ denotes m with a symmetrical uncertainty of r ulps, that is $[m - r \times \text{ulp}, m + r \times \text{ulp}]$. Adding **d** (down) or **u** (up) converts this to uncertainty in one direction only, e.g., $m?\text{d}$ denotes $[m - \frac{1}{2}\text{ulp}, m]$ and $m?\text{ru}$ denotes $[m, m + r \times \text{ulp}]$. The exponent field if present multiplies the whole interval by 10^e , e.g., $m?\text{ru}ee$ denotes $10^e \times [m, m + r \times \text{ulp}]$.

[Examples. Table 9.4 illustrates all-flavor bare interval literals. These strings are not all-flavor bare interval literals: [1.000.000], [1.0 e3], [1,2!comment], [2,1], [5?1], @5 ?1@, 5??u, [], [empty], [ganz], [1,], [1,inf].]

9.7.5 All-flavor decorated interval literals

An all-flavor decorated interval literal is a string comprising an all-flavor bare interval literal with value x , an underscore “_” and the decoration literal **com**. Its value is x_{com} .

Table 9.4. All-flavor bare interval literal examples.

Form	Literal	Exact value
Inf-sup	[1.e-3, 1.1e-3]	[0.001, 0.0011]
	[-0x1.3p-1, 2/3]	[-19/32, 2/3]
	[3.56]	[3.56, 3.56]
Uncertain	3.56?1	[3.55, 3.57]
	3.56?1e2	[355, 357]
	3.560?2	[3.558, 3.562]
	3.56?	[3.555, 3.565]
	3.560?2u	[3.560, 3.562]
	-10?	[-10.5, -9.5]
	-10?u	[-10.0, -9.5]
	-10?12	[-22.0, 2.0]

Table 9.5. Grammar for literals, using the notation of 5.12.3 of IEEE Std 754-2008.

Integer literal is `integerLiteral`, number literal is `numberLiteral`, bare interval literal is `bareIntvlLiteral` and decorated interval literal is `intervalLiteral`. `\t` denotes the TAB character.

<code>decDigit</code>	[0123456789]
<code>nonzeroDecDigit</code>	[123456789]
<code>hexDigit</code>	[0123456789abcdef]
<code>spaceChar</code>	[\t]
<code>natural</code>	{ <code>decDigit</code> } +
<code>sign</code>	[+-]
<code>integerLiteral</code>	{ <code>sign</code> } ? { <code>natural</code> }
<code>decSignificand</code>	{ <code>decDigit</code> } * "." { <code>decDigit</code> } + { <code>decDigit</code> } + "." { <code>decDigit</code> } +
<code>hexSignificand</code>	{ <code>hexDigit</code> } * "." { <code>hexDigit</code> } + { <code>hexDigit</code> } + "." { <code>hexDigit</code> } +
<code>decNumLit</code>	{ <code>sign</code> } ? { <code>decSignificand</code> } ("e" { <code>integerLiteral</code> }) ?
<code>hexNumLit</code>	{ <code>sign</code> } ? "0x" { <code>hexSignificand</code> } "p" { <code>integerLiteral</code> }
<code>positiveNatural</code>	("0") * { <code>nonzeroDecDigit</code> } { <code>decDigit</code> } *
<code>ratNumLit</code>	{ <code>integerLiteral</code> } "/" { <code>positiveNatural</code> }
<code>numberLiteral</code>	{ <code>decNumLit</code> } { <code>hexNumLit</code> } { <code>ratNumLit</code> }
<code>sp</code>	{ <code>spaceChar</code> } *
<code>dir</code>	"d" "u"
<code>pointIntvl</code>	"[" { <code>sp</code> } { <code>numberLiteral</code> } { <code>sp</code> } "]"
<code>infSupIntvl</code>	"[" { <code>sp</code> } { <code>numberLiteral</code> } { <code>sp</code> } ", " { <code>sp</code> } { <code>numberLiteral</code> } { <code>sp</code> } "]"
<code>radius</code>	{ <code>natural</code> }
<code>uncertIntvl</code>	{ <code>sign</code> } ? { <code>decSignificand</code> } "?" { <code>radius</code> } ? { <code>dir</code> } ? ("e" { <code>integerLiteral</code> }) ?
<code>bareIntvlLiteral</code>	{ <code>pointIntvl</code> } { <code>infSupIntvl</code> } { <code>uncertIntvl</code> }
<code>decorationLit</code>	"com"
<code>decoratedIntvlLiteral</code>	{ <code>bareIntvlLiteral</code> } "_" { <code>decorationLit</code> }

9.7.6 Grammar for all-flavor literals

The syntax of all-flavor integer and number literals and of all-flavor bare and decorated interval literals is defined by `integerLiteral`, `numberLiteral`, `bareIntvlLiteral` and `decoratedIntvlLiteral`, respectively, in the grammar in Table 9.5. An all-flavor literal of any of these four kinds is a string that after conversion to lowercase is accepted by this grammar.

9.8 Constructors

An interval constructor by definition is an operation that creates a bare or decorated interval from non-interval data. Constructors of bare intervals are defined in each flavor as follows.

- The bare operation `numsToInterval(l,u)` takes real values *l* and *u*. If the condition $l \leq u$ holds, the operation's value is the common bare interval $[l, u] = \{x \in \mathbb{R} \mid l \leq x \leq u\}$.
- The bare operation `textToInterval(s)` takes a text string *s*. If *s* is a bare interval literal of the flavor with value *x*, see 9.7.1, the operation's value is *x*.

These constructors shall have decorated versions as follows.

- The decorated operation `numsToInterval(l,u)` has value `newDec(x)` if the bare operation has value *x*.
- The decorated operation `textToInterval(s)` has:
 - value *x_d* if *s* is a decorated interval literal of the flavor with value *x_d*;
 - value `newDec(x)` if *s* is a bare interval literal of the flavor with value *x*.

A flavor may define extensions of these constructors and/or provide other constructors.

*[Example. A flavor might support `numsToInterval(l,u)` when *l* and/or *u* is infinite, or when $l > u$. It might support `textToInterval(s)` for similarly extended literals *s*, such as `[1,inf]` or `[2,1]`. It might provide extra arguments, e.g., letting `textToInterval("[1,2]",dac)` denote the interval $[1, 2]$ with decoration *dac*.]*

An implementation may relax at Level 2 the accuracy mode of some input strings (too long strings or strings with a rational number literal). Nevertheless, the constructor shall either return a Level 2 result containing the Level 1 value, or signal an exception, see 7.5.3.

PART 2

Set-Based Intervals

10. Level 1 description

In this clause, subclauses 10.1 to 10.4 describe the theory of mathematical intervals and interval functions that underlies this flavor. The relation between expressions and the point or interval functions that they define is specified, since it is central to the Fundamental Theorem of Interval Arithmetic. Subclauses 10.5, 10.6 list the required and recommended *arithmetic operations* with their mathematical specifications.

10.1 Non-interval Level 1 entities

In addition to intervals, the required operations of this flavor handle entities of the following kinds, as inputs or outputs; see the general considerations in 6.3.

- The set $\overline{\mathbb{R}} = \mathbb{R} \cup \{-\infty, +\infty\}$ of **extended reals**. Following the terminology of IEEE Std 754-2008 (e.g., 2.1.25), any member of $\overline{\mathbb{R}}$ is called a number: it is a **finite number** if it belongs to \mathbb{R} , else an **infinite number**.

An interval's members are finite numbers, but its bounds can be infinite. Finite or infinite numbers can be inputs to interval constructors, as well as outputs from operations, e.g., the interval width operation.

- The set of **(text) strings**, namely finite sequences of **characters** chosen from some alphabet. Since Level 1 is primarily for human communication, there are no Level 1 restrictions on the alphabet used. Strings may be inputs to interval constructors, as well as inputs or outputs of read/write operations.
- The set of **integers**, used mainly to describe parameterized sets of operations.
- The **boolean** values **false**, **true**.
- The set of **decorations** defined in 11.

10.2 Intervals

The set of mathematical intervals provided by this flavor is denoted $\overline{\mathbb{IR}}$. It comprises those subsets x of the real line \mathbb{R} that are closed and connected in the topological sense: that is, the empty set (denoted \emptyset or Empty) together with all the nonempty intervals, denoted $[\underline{x}, \overline{x}]$, defined by

$$[\underline{x}, \overline{x}] = \{x \in \mathbb{R} \mid \underline{x} \leq x \leq \overline{x}\}, \quad (4)$$

where $\underline{x} = \inf x$ and $\overline{x} = \sup x$ are extended-real numbers satisfying $\underline{x} \leq \overline{x}$, $\underline{x} < +\infty$ and $\overline{x} > -\infty$.

One calls \underline{x} the **lower bound** and \overline{x} the **upper bound** of the interval. Together they are its **bounds**. Conventionally \emptyset has bounds $\inf \emptyset = +\infty$, $\sup \emptyset = -\infty$.

NOTE 1—The definition implies $-\infty$ and $+\infty$ can be bounds of an interval, but are never members of it. In particular, (4) defines $[-\infty, +\infty]$ to be the set of all **real** numbers satisfying $-\infty \leq x \leq +\infty$, which is the whole real line \mathbb{R} —not the whole extended real line $\overline{\mathbb{R}}$.

NOTE 2—Mathematical literature generally uses a round bracket, or reversed square bracket, to show that a bound is excluded from an interval, e.g., $(a, b]$ or $]a, b]$ to denote $\{x \mid a < x \leq b\}$. Where it is convenient to use this notation, it is pointed out, e.g., in the tables of function domains and ranges in 10.5, 10.6.

NOTE 3—The set of intervals $\overline{\mathbb{IR}}$ could be described more concisely as comprising all sets $\{x \in \mathbb{R} \mid \underline{x} \leq x \leq \overline{x}\}$ for *arbitrary* extended-real $\underline{x}, \overline{x}$. However, this obtains Empty in many ways, as $[\underline{x}, \overline{x}]$ for any bounds satisfying $\underline{x} > \overline{x}$,

and also as $[-\infty, -\infty]$ or $[+\infty, +\infty]$. The description (4) was preferred as it makes a one-to-one mapping between valid pairs \underline{x}, \bar{x} of bounds and the nonempty intervals they specify.

A **box** or **interval vector** is an n -tuple $(\mathbf{x}_1, \dots, \mathbf{x}_n)$ whose components \mathbf{x}_i are intervals, that is a member of $\overline{\mathbb{R}}^n$. Usually \mathbf{x} is identified with the cartesian product $\mathbf{x}_1 \times \dots \times \mathbf{x}_n$ of its components, a subset of \mathbb{R}^n ; however, the correspondence is one-to-one only when all the \mathbf{x}_j are nonempty.

In particular $x \in \mathbf{x}$, for $x \in \mathbb{R}^n$, means by definition $x_i \in \mathbf{x}_i$ for all $i = 1, \dots, n$; and \mathbf{x} is empty if and only if any of its components \mathbf{x}_i is empty.

10.3 Hull

The (interval) **hull** of an arbitrary subset \mathbf{s} of \mathbb{R}^n , written $\text{hull}(\mathbf{s})$, is the tightest member of $\overline{\mathbb{R}}^n$ that contains \mathbf{s} . (The **tightest** set with a given property is the intersection of all sets having that property, provided the intersection itself has this property.)

10.4 Functions and expressions

The terms *function*, *domain*, *range*, *point function*, *supported*, *operation*, *provided*, *arithmetic operation* have the meanings defined in 6.1, and *expression* has the meaning in 6.2.

Subclause 7.4 requires each flavor to define the Level 1 bare interval version and decorated interval version of any point function. In this flavor, the former is the natural interval extension defined here.

Given an n -variable scalar point function f , an **interval extension** of f is a (total) mapping \mathbf{f} from n -dimensional boxes to intervals, that is $\mathbf{f} : \overline{\mathbb{R}}^n \rightarrow \overline{\mathbb{R}}$, such that $f(x) \in \mathbf{f}(\mathbf{x})$ whenever $x \in \mathbf{x}$ and $f(x)$ is defined, equivalently

$$\mathbf{f}(\mathbf{x}) \supseteq \text{Rge}(f | \mathbf{x})$$

for any box $\mathbf{x} \in \overline{\mathbb{R}}^n$, regarded as a subset of \mathbb{R}^n . The **natural interval extension** of f is the mapping \mathbf{f} defined by

$$\mathbf{f}(\mathbf{x}) = \text{hull}(\text{Rge}(f | \mathbf{x})).$$

Equivalently, using multiple-argument notation for f , an interval extension satisfies

$$\mathbf{f}(\mathbf{x}_1, \dots, \mathbf{x}_n) \supseteq \text{Rge}(f | \mathbf{x}_1, \dots, \mathbf{x}_n),$$

and the natural interval extension is defined by

$$\mathbf{f}(\mathbf{x}_1, \dots, \mathbf{x}_n) = \text{hull}(\text{Rge}(f | \mathbf{x}_1, \dots, \mathbf{x}_n))$$

for any intervals $\mathbf{x}_1, \dots, \mathbf{x}_n$.

In some contexts, it is useful for \mathbf{x} to be a general subset of \mathbb{R}^n , or the \mathbf{x}_i to be general subsets of \mathbb{R} ; the definition is unchanged. In this case, we refer to each of the above extensions as a *full Level 1 extension* of f .

The natural extension is automatically defined for all interval or set arguments. The decoration system, Clause 11, gives a way of diagnosing when the underlying point function has been evaluated outside its domain.

When f is a binary operator \bullet written in infix notation, this gives the usual definition of its natural interval extension as

$$\mathbf{x} \bullet \mathbf{y} = \text{hull}(\{x \bullet y \mid x \in \mathbf{x}, y \in \mathbf{y}, \text{ and } x \bullet y \text{ is defined}\}).$$

[Example. With these definitions, the relevant natural interval extensions satisfy $\sqrt{[-1, 4]} = [0, 2]$ and $\sqrt{[-2, -1]} = \emptyset$; also $\mathbf{x} \times [0, 0] = [0, 0]$ for any nonempty \mathbf{x} , and $\mathbf{x}/[0, 0] = \emptyset$, for any \mathbf{x} .]

When f is a vector point function, a vector interval function with the same number of inputs and outputs as f is called an interval extension of f if each of its components is an interval extension of the corresponding component of f .

From the above definition, an interval extension of a **real constant**—a zero-argument point function returning a real value c —is any zero-argument interval function that returns an interval containing c . Its *natural extension* returns the single-point interval $[c, c]$.

10.5 Required operations

The operations listed in this subclause include those required in all flavors, see Clause 9. An implementation shall provide interval versions of them appropriate to its supported interval types. For constants and the forward and reverse arithmetic operations in 10.5.2, 10.5.3, 10.5.4, each such version shall be an interval extension (10.4) of the corresponding point function—for a constant, that means any constant interval enclosing the point value. The required rounding behavior of these operations, and of the numeric functions of intervals in 10.5.9, is detailed in 12.9, 12.12.

The names of operations, as well as symbols used for operations (e.g., for the comparisons in 10.5.10), might not correspond to those that a particular programming language would use.

10.5.1 Interval literals

An **interval literal** is a text string that denotes an interval. Level 1, which is mainly for human communication, merely assumes there exist some agreed rules on the form and meaning of interval literals. A specified form and meaning is used in Level 2 onward: the definition is in 12.11, extending 9.7. This definition is also used in Level 1 of this document for examples, where relevant. *[Example. This includes the inf-sup form $[1.234e5, \text{Inf}]$; the uncertain form $3.1416?1$; and the named interval constant $[\text{Empty}]$.]*

10.5.2 Interval constants

The constant functions `empty()` and `entire()` have values `Empty` and `Entire`, respectively.

10.5.3 Forward-mode elementary functions

Table 9.1 lists required arithmetic operations. The term *operation* includes functions normally written in function notation $f(x, y, \dots)$, as well as those normally written in unary or binary operator notation, $\bullet x$ or $x \bullet y$.

10.5.4 Reverse-mode elementary functions

Constraint-satisfaction algorithms use the operations in this subclause for iteratively tightening an enclosure of a solution to a system of equations. They are listed in Table 10.1.

Given a unary arithmetic operation φ , a **reverse interval extension** of φ is a binary interval function φRev such that

$$\varphi\text{Rev}(\mathbf{c}, \mathbf{x}) \supseteq \{x \in \mathbf{x} \mid \varphi(x) \text{ is defined and in } \mathbf{c}\}, \quad (5)$$

for any intervals \mathbf{c}, \mathbf{x} .

Similarly, a binary arithmetic operation \bullet has two forms of reverse interval extension, which are ternary interval functions $\bullet\text{Rev}_1$ and $\bullet\text{Rev}_2$ such that

$$\bullet\text{Rev}_1(\mathbf{b}, \mathbf{c}, \mathbf{x}) \supseteq \{x \in \mathbf{x} \mid b \in \mathbf{b} \text{ exists such that } x \bullet b \text{ is defined and in } \mathbf{c}\}, \quad (6)$$

$$\bullet\text{Rev}_2(\mathbf{a}, \mathbf{c}, \mathbf{x}) \supseteq \{x \in \mathbf{x} \mid a \in \mathbf{a} \text{ exists such that } a \bullet x \text{ is defined and in } \mathbf{c}\}. \quad (7)$$

If \bullet is commutative, then $\bullet\text{Rev}_1$ and $\bullet\text{Rev}_2$ agree and may be implemented simply as $\bullet\text{Rev}$.

In each of (5, 6, 7), the unique **natural reverse interval extension** is the one whose value is the interval hull of the right-hand side. Clearly, any reverse interval extension encloses this hull.

The last argument \mathbf{x} in each of (5, 6, 7) is optional, with default $\mathbf{x} = \mathbb{R}$ if absent.

NOTE—The argument \mathbf{x} can be thought of as giving prior knowledge about the range of values taken by a point-variable x , which is then sharpened by applying the reverse function: see the example below.

Table 10.1. Required reverse elementary functions.

From unary functions	From binary functions
$\text{sqrRev}(\mathbf{c}, \mathbf{x})$	$\text{mulRev}(\mathbf{b}, \mathbf{c}, \mathbf{x})$
$\text{absRev}(\mathbf{c}, \mathbf{x})$	$\text{powRev1}(\mathbf{b}, \mathbf{c}, \mathbf{x})$
$\text{pownRev}(\mathbf{c}, \mathbf{x}, p)$	$\text{powRev2}(\mathbf{a}, \mathbf{c}, \mathbf{x})$
$\text{sinRev}(\mathbf{c}, \mathbf{x})$	$\text{atan2Rev1}(\mathbf{b}, \mathbf{c}, \mathbf{x})$
$\text{cosRev}(\mathbf{c}, \mathbf{x})$	$\text{atan2Rev2}(\mathbf{a}, \mathbf{c}, \mathbf{x})$
$\text{tanRev}(\mathbf{c}, \mathbf{x})$	
$\text{coshRev}(\mathbf{c}, \mathbf{x})$	

Reverse versions of division and reciprocal are not provided. In relevant applications, their effect can be obtained by changing $z = x/y$ to $x = y \times z$ and using reverse multiplication. $\text{pownRev}(x, p)$ is regarded as a family of unary functions parameterized by p .

[Example.

- Consider the function $\text{sqr}(x) = x^2$. Evaluating $\text{sqrRev}([1, 4])$ answers the question: given that $1 \leq x^2 \leq 4$, what interval can we restrict x to? Using the natural reverse extension, we have

$$\text{sqrRev}([1, 4]) = \text{hull}\{x \in \mathbb{R} \mid x^2 \in [1, 4]\} = \text{hull}([-2, -1] \cup [1, 2]) = [-2, 2].$$

- If we can add the prior knowledge that $x \in \mathbf{x} = [0, 1.2]$, then using the optional second argument gives the tighter enclosure

$$\text{sqrRev}([1, 4], [0, 1.2]) = \text{hull}\{x \in [0, 1.2] \mid x^2 \in [1, 4]\} = \text{hull}([0, 1.2] \cap ([-2, -1] \cup [1, 2])) = [1, 1.2].$$

- One might think it suffices to apply the operation without the optional argument and intersect the result with \mathbf{x} . This is less effective because “hull” and “intersect” do not commute. E.g., in the above, this method evaluates

$$\text{sqrRev}([1, 4]) \cap \mathbf{x} = [-2, 2] \cap [0, 1.2] = [0, 1.2],$$

so no tightening of the enclosure \mathbf{x} is obtained.

]

10.5.5 Two-output division

Some algorithms exploit the fact that the result of interval division, considered as a set, can have zero, one or two disjoint connected components. An important example is the one-dimensional Interval Newton Method, which also treats such division as reverse multiplication (so that $0/0$ is “any real number” instead of “undefined”):

$$\mathbf{c} /_{\text{set}} \mathbf{b} = \{x \in \mathbb{R} \mid b \in \mathbf{b} \text{ exists such that } xb \in \mathbf{c}\}.$$

[Examples. Using the classical notation where a square or round bracket means a closed or open bound, respectively, then

$$\begin{aligned} [1, 2] /_{\text{set}} [0, 0] &= \emptyset, \\ [0, 2] /_{\text{set}} [0, 0] &= \text{Entire}, \\ [1, 2] /_{\text{set}} [1, 1] &= [1, 2], \\ [1, 1] /_{\text{set}} [1, +\infty) &= (0, 1], \\ [1, 2] /_{\text{set}} [-1, 1] &= (-\infty, -1] \cup [1, +\infty), \\ [1, 1] /_{\text{set}} \text{Entire} &= (-\infty, 0) \cup (0, +\infty), \end{aligned}$$

are cases with 0, 1, 1, 1, 2 and 2 output components, respectively. The fourth and sixth cases have components that are not closed.]

The `mulRevToPair` operation outputs the components separately. Its value is an ordered pair of closed intervals:

$$\text{mulRevToPair}(\mathbf{b}, \mathbf{c}) = \begin{cases} (\emptyset, \emptyset) & \text{if } \mathbf{c} /_{\text{set}} \mathbf{b} \text{ is empty,} \\ (\overline{\mathbf{u}}, \emptyset) & \text{if } \mathbf{c} /_{\text{set}} \mathbf{b} \text{ has one component } \mathbf{u}, \\ (\overline{\mathbf{u}}, \overline{\mathbf{v}}) & \text{if } \mathbf{c} /_{\text{set}} \mathbf{b} \text{ has two components } \mathbf{u}, \mathbf{v}, \text{ ordered so that } \mathbf{u} < \mathbf{v}, \end{cases}$$

where $\overline{\mathbf{a}}$ denotes the topological closure of \mathbf{a} , which in this case is equivalent to $\text{hull}(\mathbf{a})$.

NOTE—`mulRevToPair` is not regarded as an arithmetic operation, since if it appears in an arithmetic expression, containment might be lost unless its two outputs are handled with care.

10.5.6 Cancellative addition and subtraction

Cancellative subtraction solves the problem: Recover interval \mathbf{z} from intervals \mathbf{x} and \mathbf{y} , given that one knows \mathbf{x} was obtained as the sum $\mathbf{y} + \mathbf{z}$.

[Example. In some applications, one has a list of intervals $\mathbf{a}_1, \dots, \mathbf{a}_n$, and needs to form each interval \mathbf{s}_k which is the sum of all the \mathbf{a}_i except \mathbf{a}_k , that is $\mathbf{s}_k = \sum_{i=1, i \neq k}^n \mathbf{a}_i$, for $k = 1, \dots, n$. Evaluating all these sums independently costs $O(n^2)$ work. However, if one forms the sum \mathbf{s} of all the \mathbf{a}_i , one can obtain each \mathbf{s}_k from \mathbf{s} and \mathbf{a}_k by cancellative subtraction. This method only costs $O(n)$ work.

This example illustrates that in finite precision, computing \mathbf{x} (as a sum of terms) typically incurs at least one roundoff error, and might incur many. Thus the assumption underlying these cancellative operations is that \mathbf{x} is an enclosure of an unknown true sum \mathbf{x}_0 , whereas \mathbf{y} is “exact.” The computed \mathbf{z} is an enclosure of an unknown true \mathbf{z}_0 such that $\mathbf{y} + \mathbf{z}_0 = \mathbf{x}_0$.]

The operation `cancelPlus`(\mathbf{x}, \mathbf{y}) is equivalent to `cancelMinus`($\mathbf{x}, -\mathbf{y}$) and therefore not specified separately.

For any two bounded intervals \mathbf{x} and \mathbf{y} , the value of the operation `cancelMinus`(\mathbf{x}, \mathbf{y}) is the tightest interval \mathbf{z} such that

$$\mathbf{y} + \mathbf{z} \supseteq \mathbf{x}, \quad (8)$$

if such a \mathbf{z} exists. Otherwise `cancelMinus`(\mathbf{x}, \mathbf{y}) has no value at Level 1.

This specification leads to the following Level 1 algorithm. If $\mathbf{x} = \emptyset$ and \mathbf{y} is bounded, then $\mathbf{z} = \emptyset$. If $\mathbf{x} \neq \emptyset$ and $\mathbf{y} = \emptyset$, then \mathbf{z} has no value. If $\mathbf{x} = [\underline{x}, \overline{x}]$ and $\mathbf{y} = [\underline{y}, \overline{y}]$ are both nonempty and bounded, define $\underline{z} = \underline{x} - \underline{y}$ and $\overline{z} = \overline{x} - \overline{y}$. Then \mathbf{z} is defined to be $[\underline{z}, \overline{z}]$ if $\underline{z} \leq \overline{z}$ (equivalently if $\text{width}(\mathbf{x}) \geq \text{width}(\mathbf{y})$), and has no value otherwise. If either \mathbf{x} or \mathbf{y} is unbounded, \mathbf{z} has no value.

NOTE—Because of the cancellative nature of these operations, care is needed in finite precision to determine whether the result is defined or not.

10.5.7 Set operations

The value of the operation `intersection`(\mathbf{x}, \mathbf{y}) is the intersection $\mathbf{x} \cap \mathbf{y}$ of the intervals \mathbf{x} and \mathbf{y} .

The value of the operation `convexHull`(\mathbf{x}, \mathbf{y}) is the interval hull of the union $\mathbf{x} \cup \mathbf{y}$ of the intervals \mathbf{x} and \mathbf{y} .

10.5.8 Constructors

An interval constructor by definition is an operation that creates a bare or decorated interval from non-interval data. The following bare interval constructors are defined.

The operation `numsToInterval`(l, u), takes extended-real values l and u . If (see 10.2) the conditions $l \leq u$, $l < +\infty$ and $u > -\infty$ hold, its value is the nonempty interval

$$[l, u] = \{x \in \mathbb{R} \mid l \leq x \leq u\}.$$

Otherwise, it has no value.

The operation `textToInterval`(\mathbf{s}) takes a text string \mathbf{s} . If \mathbf{s} is an interval literal, see 9.7, 12.11, its value is the interval denoted by \mathbf{s} . Otherwise, it has no value.

10.5.9 Numeric functions of intervals

The operations in Table 10.2 are defined, the argument being an interval and the result a number, which for some of the operations might be infinite.

Implementations should provide an operation that returns $\text{mid}(\mathbf{x})$ and $\text{rad}(\mathbf{x})$ simultaneously.

Table 10.2. Required numeric functions of an interval $\mathbf{x} = [\underline{x}, \bar{x}]$.

Note sup can have value $-\infty$; each of inf , wid , rad and mag can have value $+\infty$.

Name	Definition
$\text{inf}(\mathbf{x})$	$\begin{cases} \text{lower bound of } \mathbf{x}, \text{ if } \mathbf{x} \text{ is nonempty} \\ \infty, \text{ if } \mathbf{x} \text{ is empty} \end{cases}$
$\text{sup}(\mathbf{x})$	$\begin{cases} \text{upper bound of } \mathbf{x}, \text{ if } \mathbf{x} \text{ is nonempty} \\ -\infty, \text{ if } \mathbf{x} \text{ is empty} \end{cases}$
$\text{mid}(\mathbf{x})$	$\begin{cases} \text{midpoint } (\underline{x} + \bar{x})/2, \text{ if } \mathbf{x} \text{ is nonempty bounded} \\ \text{no value, if } \mathbf{x} \text{ is empty or unbounded} \end{cases}$
$\text{wid}(\mathbf{x})$	$\begin{cases} \text{width } \bar{x} - \underline{x}, \text{ if } \mathbf{x} \text{ is nonempty} \\ \text{no value, if } \mathbf{x} \text{ is empty} \end{cases}$
$\text{rad}(\mathbf{x})$	$\begin{cases} \text{radius } (\bar{x} - \underline{x})/2, \text{ if } \mathbf{x} \text{ is nonempty} \\ \text{no value, if } \mathbf{x} \text{ is empty} \end{cases}$
$\text{mag}(\mathbf{x})$	$\begin{cases} \text{magnitude } \sup\{ x \mid x \in \mathbf{x}\}, \text{ if } \mathbf{x} \text{ is nonempty} \\ \text{no value, if } \mathbf{x} \text{ is empty} \end{cases}$
$\text{mig}(\mathbf{x})$	$\begin{cases} \text{mignitude } \inf\{ x \mid x \in \mathbf{x}\}, \text{ if } \mathbf{x} \text{ is nonempty} \\ \text{no value, if } \mathbf{x} \text{ is empty} \end{cases}$

10.5.10 Boolean functions of intervals

The following operations are defined, whose value is a boolean (1 = true, 0 = false) result.

There is a function $\text{isEmpty}(\mathbf{x})$, with value 1 if \mathbf{x} is the empty set, 0 otherwise. There is a function $\text{isEntire}(\mathbf{x})$, with value 1 if \mathbf{x} is the whole line, 0 otherwise.

There are eight boolean-valued comparison relations, which take two interval inputs. These are defined in Table 10.3, in which column three gives the set-theoretic definition, and column four gives an equivalent specification when both intervals are nonempty. Table 10.4 shows what the definitions imply when at least one interval is empty.

Table 10.3. Comparisons for intervals \mathbf{a} and \mathbf{b} . Notation \forall_a means “for all a in \mathbf{a} ,” and so on. In column 4, $\underline{a}=[\underline{a}, \bar{a}]$ and $\underline{b}=[\underline{b}, \bar{b}]$, where $\underline{a}, \underline{b}$ may be $-\infty$, and \bar{a}, \bar{b} may be $+\infty$; and $<'$ is the same as $<$ except that $-\infty <' -\infty$ and $+\infty <' +\infty$ are true.

Name	Symbol	Defining predicate	For $\mathbf{a}, \mathbf{b} \neq \emptyset$	Description
$\text{equal}(\mathbf{a}, \mathbf{b})$	$\mathbf{a} = \mathbf{b}$	$\forall_a \exists_b a = b \wedge \forall_b \exists_a b = a$	$\underline{a} = \underline{b} \wedge \bar{a} = \bar{b}$	\mathbf{a} equals \mathbf{b}
$\text{subset}(\mathbf{a}, \mathbf{b})$	$\mathbf{a} \subseteq \mathbf{b}$	$\forall_a \exists_b a = b$	$\underline{b} \leq \underline{a} \wedge \bar{a} \leq \bar{b}$	\mathbf{a} is a subset of \mathbf{b}
$\text{less}(\mathbf{a}, \mathbf{b})$	$\mathbf{a} \leq \mathbf{b}$	$\forall_a \exists_b a \leq b \wedge \forall_b \exists_a a \leq b$	$\underline{a} \leq \underline{b} \wedge \bar{a} \leq \bar{b}$	\mathbf{a} is weakly less than \mathbf{b}
$\text{precedes}(\mathbf{a}, \mathbf{b})$	$\mathbf{a} \prec \mathbf{b}$	$\forall_a \forall_b a \leq b$	$\bar{a} \leq \underline{b}$	\mathbf{a} is to left of but may touch \mathbf{b}
$\text{interior}(\mathbf{a}, \mathbf{b})$	$\mathbf{a} \oslash \mathbf{b}$	$\forall_a \exists_b a < b \wedge \forall_a \exists_b b < a$	$\underline{b} <' \underline{a} \wedge \bar{a} <' \bar{b}$	\mathbf{a} is interior to \mathbf{b}
$\text{strictLess}(\mathbf{a}, \mathbf{b})$	$\mathbf{a} < \mathbf{b}$	$\forall_a \exists_b a < b \wedge \forall_b \exists_a a < b$	$\underline{a} <' \underline{b} \wedge \bar{a} <' \bar{b}$	\mathbf{a} is strictly less than \mathbf{b}
$\text{strictPrecedes}(\mathbf{a}, \mathbf{b})$	$\mathbf{a} < \mathbf{b}$	$\forall_a \forall_b a < b$	$\bar{a} < \underline{b}$	\mathbf{a} is strictly to left of \mathbf{b}
$\text{disjoint}(\mathbf{a}, \mathbf{b})$	$\mathbf{a} \not\cap \mathbf{b}$	$\forall_a \forall_b a \neq b$	$\bar{a} < \underline{b} \vee \bar{b} < \underline{a}$	\mathbf{a} and \mathbf{b} are disjoint

NOTE 1—Column two of Table 10.3 gives suggested symbols for use in typeset algorithms.

NOTE 2—All these relations, except $\mathbf{a} \not\cap \mathbf{b}$, are transitive for *nonempty* intervals.

NOTE 3—The first three are reflexive.

NOTE 4—The rules of set theory imply some counter-intuitive facts, e.g., the empty interval is both a subset of any interval and disjoint from it. Also, **interior** uses the topological definition: **b** is a neighbourhood of each point of **a**. This implies, for instance, that **interior**(Entire,Entire) is true.

NOTE 5—All occurrences of $<$ in column 4 of Table 10.3 can be replaced by $<'$.

Table 10.4. Comparisons with empty intervals

	$a = \emptyset$ $b \neq \emptyset$	$a \neq \emptyset$ $b = \emptyset$	$a = \emptyset$ $b = \emptyset$
$a = b$	0	0	1
$a \subseteq b$	1	0	1
$a \leq b$	0	0	1
$a \prec b$	1	1	1
$a \odot b$	1	0	1
$a < b$	0	0	1
$a \prec b$	1	1	1
$a \not\prec b$	1	1	1

10.6 Recommended operations

An implementation should provide interval versions of the functions listed in this subclause. If such an interval version is provided, it shall behave as specified here.

10.6.1 Forward-mode elementary functions

The list of recommended functions is in Table 10.5. Each interval version shall be an interval extension of the point function.

Table 10.5. Recommended elementary functions. Normal mathematical notation is used to include or exclude an interval bound, e.g., $[-1, 1]$ denotes $\{x \in \mathbb{R} \mid -1 \leq x \leq 1\}$.

Name	Definition	Point function domain	Point function range	Table Footnotes
$\text{rootn}(x, q)$	real $\sqrt[q]{x}$, $q \in \mathbb{Z} \setminus \{0\}$	$\begin{cases} \mathbb{R} & \text{if } q > 0 \text{ odd} \\ [0, \infty) & \text{if } q > 0 \text{ even} \\ \mathbb{R} \setminus \{0\} & \text{if } q < 0 \text{ odd} \\ (0, \infty) & \text{if } q < 0 \text{ even} \end{cases}$	same as domain	a
$\begin{cases} \text{expm1}(x) \\ \text{exp2m1}(x) \\ \text{exp10m1}(x) \end{cases}$	$b^x - 1$	\mathbb{R}	$(-1, \infty)$	b, c
$\begin{cases} \text{logp1}(x) \\ \text{log2p1}(x) \\ \text{log10p1}(x) \end{cases}$	$\log_b(x+1)$	$(-1, \infty)$	\mathbb{R}	b, c
$\text{compoundm1}(x, y)$	$(1+x)^y - 1$	$\{x > -1\} \cup \{x = -1, y > 0\}$	$[-1, \infty)$	c, d
$\text{hypot}(x, y)$	$\sqrt{x^2 + y^2}$	\mathbb{R}^2	$[0, \infty)$	
$\text{rSqrt}(x)$	$1/\sqrt{x}$	$(0, \infty)$	$(0, \infty)$	
$\text{sinPi}(x)$	$\sin(\pi x)$	\mathbb{R}	$[-1, 1]$	e
$\text{cosPi}(x)$	$\cos(\pi x)$	\mathbb{R}	$[-1, 1]$	e
$\text{tanPi}(x)$	$\tan(\pi x)$	$\mathbb{R} \setminus \{k + \frac{1}{2} \mid k \in \mathbb{Z}\}$	\mathbb{R}	e
$\text{asinPi}(x)$	$\arcsin(x)/\pi$	$[-1, 1]$	$[-1/2, 1/2]$	e
$\text{acosPi}(x)$	$\arccos(x)/\pi$	$[-1, 1]$	$[0, 1]$	e
$\text{atanPi}(x)$	$\arctan(x)/\pi$	\mathbb{R}	$(-1/2, 1/2)$	e
$\text{atan2Pi}(y, x)$	$\text{atan2}(y, x)/\pi$	$\mathbb{R}^2 \setminus \{(0, 0)\}$	$(-1, 1]$	e, f

- a. Regarded as a family of functions of one real variable x , parameterized by the integer argument q .
- b. $b = e, 2$ or 10 , respectively.
- c. Mathematically unnecessary, but included to let implementations give better numerical behavior for small values of the arguments.
- d. In describing domains, notation such as $\{y = 0\}$ is short for $\{(x, y) \in \mathbb{R}^2 \mid y = 0\}$, and so on.
- e. These functions avoid a loss of accuracy due to π being irrational, cf. Table 9.1, note e.
- f. To avoid confusion with notation for open intervals, here coordinates in \mathbb{R}^2 are delimited by angle brackets $\langle \rangle$.

10.6.2 Slope functions

The functions in Table 10.6 are the commonest ones needed for efficient implementation of improved range enclosures via *first- and second-order slope* algorithms. They are analytic at $x = 0$ after filling in the removable singularity there, where each has the value 1.

Table 10.6. Recommended slope functions

Name	Definition	Point function domain	Point function range
$\text{expSlope1}(x)$	$\frac{1}{x}(e^x - 1)$	\mathbb{R}	$(0, \infty)$
$\text{expSlope2}(x)$	$\frac{2}{x^2}(e^x - 1 - x)$	\mathbb{R}	$(0, \infty)$
$\text{logSlope1}(x)$	$-\frac{2}{x^2}(\log(1+x) - x)$	\mathbb{R}	$(0, \infty)$
$\text{logSlope2}(x)$	$\frac{3}{x^3}(\log(1+x) - x + \frac{x^2}{2})$	\mathbb{R}	$(0, \infty)$
$\text{cosSlope2}(x)$	$-\frac{2}{x^2}(\cos x - 1)$	\mathbb{R}	$[0, 1]$
$\text{sinSlope3}(x)$	$-\frac{6}{x^3}(\sin x - x)$	\mathbb{R}	$(0, 1]$
$\text{asinSlope3}(x)$	$\frac{6}{x^3}(\arcsin x - x)$	$[-1, 1]$	$[1, 3\pi - 6]$
$\text{atanSlope3}(x)$	$-\frac{3}{x^3}(\arctan x - x)$	\mathbb{R}	$(0, 1]$
$\text{coshSlope2}(x)$	$\frac{2}{x^2}(\cosh x - 1)$	\mathbb{R}	$[1, \infty)$
$\text{sinhSlope3}(x)$	$\frac{3}{x^3}(\sinh x - x)$	\mathbb{R}	$[\frac{1}{2}, \infty)$

10.6.3 Boolean functions of intervals

The following operations should be provided, whose value is a boolean ($1 = \text{true}$, $0 = \text{false}$) result.

There is a function $\text{isCommonInterval}(\mathbf{x})$, with value 1 if \mathbf{x} is a common interval (7.2) and 0 otherwise. There is a function $\text{isSingleton}(\mathbf{x})$, with value 1 if \mathbf{x} is the set of a single real and 0 otherwise.

There is a function $\text{isMember}(m, \mathbf{x})$, where m is an extended real and \mathbf{x} is an interval. Its value is 1 if m is a finite real and m is a member of the set \mathbf{x} , 0 otherwise.

10.6.4 Extended interval comparisons

The **interval overlapping function** $\text{overlap}(\mathbf{a}, \mathbf{b})$, also written $\mathbf{a} \oslash \mathbf{b}$, arises from the work of J.F. Allen [B2] on temporal logic. It may be used as an infrastructure for other interval comparisons. If implemented, it should also be available at user level; how this is done is implementation-defined or language-defined.

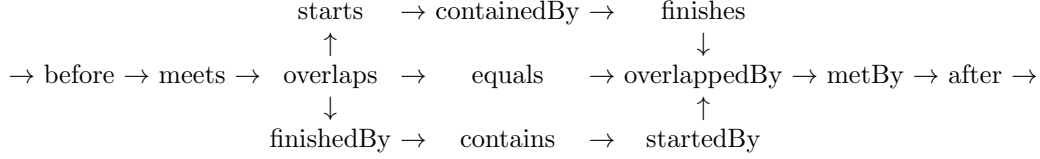
Table 10.7. The 16 states of interval overlapping situations for intervals a, b .

Notation \forall_a means “for all a in \mathbf{a} ,” and so on. Phrases within a cell are joined by “and,” e.g., **starts** is specified by $(\underline{a} = \underline{b} \wedge \bar{a} < \bar{b})$.

State $\mathbf{a} \oslash \mathbf{b}$ is	Set specification	Bound specification	Diagram
States with either interval empty			
bothEmpty	$\mathbf{a} = \emptyset \wedge \mathbf{b} = \emptyset$		
firstEmpty	$\mathbf{a} = \emptyset \wedge \mathbf{b} \neq \emptyset$		
secondEmpty	$\mathbf{a} \neq \emptyset \wedge \mathbf{b} = \emptyset$		
States with both intervals nonempty			
before	$\forall_a \forall_b a < b$	$\bar{a} < \underline{b}$	
meets	$\forall_a \forall_b a \leq b$ $\exists_a \forall_b a < b$ $\exists_a \exists_b a = b$	$\underline{a} < \bar{a}$ $\bar{a} = \underline{b}$ $\underline{b} < \bar{b}$	
overlaps	$\exists_a \forall_b a < b$ $\exists_b \forall_a a < b$ $\exists_a \exists_b b < a$	$\underline{a} < \underline{b}$ $\underline{b} < \bar{a}$ $\bar{a} < \bar{b}$	
starts	$\forall_b \exists_a a \leq b$ $\forall_a \exists_b b \leq a$ $\exists_b \forall_a a < b$	$\underline{a} = \underline{b}$ $\bar{a} < \bar{b}$	
containedBy	$\exists_b \forall_a b < a$ $\exists_b \forall_a a < b$	$\underline{b} < \underline{a}$ $\bar{a} < \bar{b}$	
finishes	$\exists_b \forall_a b < a$ $\forall_b \exists_a b \leq a$ $\forall_a \exists_b a \leq b$	$\underline{b} < \underline{a}$ $\bar{a} = \bar{b}$	
equals	$\forall_a \exists_b a = b$ $\forall_b \exists_a b = a$	$\underline{a} = \underline{b}$ $\bar{a} = \bar{b}$	
finishedBy	$\exists_a \forall_b a < b$ $\forall_a \exists_b a \leq b$ $\forall_b \exists_a b \leq a$	$\underline{a} < \underline{b}$ $\bar{b} = \bar{a}$	
contains	$\exists_a \forall_b a < b$ $\exists_a \forall_b b < a$	$\underline{a} < \underline{b}$ $\bar{b} < \bar{a}$	
startedBy	$\forall_a \exists_b b \leq a$ $\forall_b \exists_a a \leq b$ $\exists_a \forall_b b < a$	$\underline{b} = \underline{a}$ $\bar{b} < \bar{a}$	
overlappedBy	$\exists_b \forall_a b < a$ $\exists_a \forall_b b < a$ $\exists_b \exists_a a < b$	$\underline{b} < \underline{a}$ $\underline{a} < \bar{b}$ $\bar{b} < \bar{a}$	
metBy	$\forall_b \forall_a b \leq a$ $\exists_b \exists_a b = a$ $\exists_b \forall_a b < a$	$\underline{b} < \bar{b}$ $\bar{b} = \underline{a}$ $\underline{a} < \bar{a}$	
after	$\forall_b \forall_a b < a$	$\bar{b} < \underline{a}$	

Allen identified 13 states of a pair (\mathbf{a}, \mathbf{b}) of nonempty intervals, which are ways in which they can be related with respect to the usual order $a < b$ of the reals. Together with three states for when either interval is empty, these define the 16 possible values of $\text{overlap}(\mathbf{a}, \mathbf{b})$.

To describe the states for nonempty intervals of positive width, it is useful to think of $\mathbf{b} = [\underline{b}, \bar{b}]$ (with $\underline{b} < \bar{b}$) as fixed, while $\mathbf{a} = [\underline{a}, \bar{a}]$ (with $\underline{a} < \bar{a}$) starts far to its left and moves to the right. Its bounds move continuously with strictly positive velocity. Then, depending on the relative sizes of \mathbf{a} and \mathbf{b} , the value of $\mathbf{a} \oslash \mathbf{b}$ follows a path from left to right through the graph below, whose nodes represent Allen's 13 states.



For instance, “ \mathbf{a} overlaps \mathbf{b} ”—equivalently $\mathbf{a} \oslash \mathbf{b}$ has the value **overlaps**—is the case $\underline{a} < \underline{b} < \bar{a} < \bar{b}$.

The three extra values are: **bothEmpty** when $\mathbf{a} = \mathbf{b} = \emptyset$, else **firstEmpty** when $\mathbf{a} = \emptyset$, **secondEmpty** when $\mathbf{b} = \emptyset$.

Table 10.7 shows the 16 states, with the 13 “nonempty” states specified (a) in terms of set membership using quantifiers and (b) in terms of the bounds $\underline{a}, \bar{a}, \underline{b}, \bar{b}$, and also (c) shown diagrammatically.

The “set” and “bound” columns in Table 10.7 remove some ambiguities of the diagram view when one interval shrinks to a single point that coincides with a bound of the other. Such a case is allocated to **equal** when all four bounds coincide; else to **starts**, **finishes**, **finishedBy** or **startedBy** as appropriate; never to **meets** or **metBy**.

NOTE—The 16 state values can be encoded in four bits. However, if they are then translated into patterns P in a 16-bit word, having one position equal to 1 and the rest zero, one can easily implement interval comparisons by using bit-masks.

For instance, suppose we make the states s in Table 10.7’s order correspond to the 16 bits in the word, left-to-right, so $s = \text{bothEmpty}$ maps to $P(s) = 1000000000000000$, $s = \text{firstEmpty}$ maps to $P(s) = 0100000000000000$, and so on. Consider the relation $\text{disjoint}(\mathbf{a}, \mathbf{b})$. This is true if and only if one or both of \mathbf{a} or \mathbf{b} is empty, or \mathbf{a} is “before” \mathbf{b} , or \mathbf{a} is “after” \mathbf{b} . That is, iff the logical “and” of $P(s)$ with the mask $\text{disjointMask} = 1111000000000001$ is not identically zero.

This scheme can be efficiently implemented in hardware, see for instance M. Nehmeier, S. Siegel and J. Wolff von Gudenberg [B9]. All the required comparisons in this standard can be implemented in this way, as can be, e.g., the “possibly” and “certainly” comparisons of Sun’s interval Fortran [B16]. Thus the overlap operation is a primitive from which it is simple to derive all interval comparisons commonly found in the literature.

11. The decoration system at Level 1

11.1 Decorations and decorated intervals overview

The decoration system of the set-based flavor conforms to the principles of 8.1. An implementation makes the decoration system available by providing:

- a decorated version of each interval extension of an arithmetic operation, of each interval constructor, and of some other operations;
- various auxiliary functions, e.g., to extract a decorated interval’s interval and decoration parts, and to apply a standard initial decoration to an interval.

The system is specified here at a mathematical level, with the finite-precision aspects throughout Clause 12. Subclauses 11.2, 11.3, 11.4 give the basic concepts. 11.5, 11.6 define how intervals are given an initial decoration, and the binding of decorations to library interval arithmetic operations to give correct propagation through expressions. 11.7 is about non-arithmetic operations. 11.5 describes housekeeping operations on decorations, including comparisons, and conversion between a decorated interval and its interval and decoration parts. 11.8 discusses the decoration of user-defined arithmetic operations. The decoration **com** makes it possible to verify, under fairly restrictive conditions, whether a given computation gives the same result in different flavors; 11.9 gives explanatory notes on the **com** decoration. 11.10 defines a restricted decorated arithmetic that suffices for some important applications and is easier to implement efficiently.

The Annex B contains a rigorous theoretical foundation, including a proof of the Fundamental Theorem of Interval Arithmetic for this flavor.

11.2 Definitions and basic properties

Formally, a decoration d is a property (that is, a boolean-valued function) $p_d(f, \mathbf{x})$ of pairs (f, \mathbf{x}) , where f is a real-valued function with domain $\text{Dom}(f) \subseteq \mathbb{R}^n$ for some $n \geq 0$ and $\mathbf{x} \in \overline{\mathbb{IR}}^n$ is an n -dimensional box, regarded as a subset of \mathbb{R}^n . The notation (f, \mathbf{x}) unless said otherwise denotes such a pair, for arbitrary n , f and \mathbf{x} . Equivalently, d is identified with the set of pairs for which the property holds:

$$d = \{ (f, \mathbf{x}) \mid p_d(f, \mathbf{x}) \text{ is true} \}. \quad (9)$$

The set \mathbb{D} of decorations has five members:

Value	Short description	Property	Definition
com	common	$p_{\text{com}}(f, \mathbf{x})$	\mathbf{x} is a bounded, nonempty subset of $\text{Dom}(f)$; f is continuous at each point of \mathbf{x} ; and the computed interval $f(\mathbf{x})$ is bounded.
dac	defined & continuous	$p_{\text{dac}}(f, \mathbf{x})$	\mathbf{x} is a nonempty subset of $\text{Dom}(f)$, and the restriction of f to \mathbf{x} is continuous;
def	defined	$p_{\text{def}}(f, \mathbf{x})$	\mathbf{x} is a nonempty subset of $\text{Dom}(f)$;
trv	trivial	$p_{\text{trv}}(f, \mathbf{x})$	always true (so gives no information);
ill	ill-formed	$p_{\text{ill}}(f, \mathbf{x})$	Not an Interval; formally $\text{Dom}(f) = \emptyset$, see 11.3.

These are listed according to the propagation order (19), which may also be thought of as a quality-order of (f, \mathbf{x}) pairs—decorations above **trv** are “good” and those below are “bad.”

A **decorated interval** is a pair, written interchangeably as (\mathbf{u}, d) or \mathbf{u}_d , where $\mathbf{u} \in \overline{\mathbb{IR}}$ is a real interval and $d \in \mathbb{D}$ is a decoration. (\mathbf{u}, d) may also denote a decorated box $((\mathbf{u}_1, d_1), \dots, (\mathbf{u}_n, d_n))$, where \mathbf{u} and d are the vectors of interval parts \mathbf{u}_i and decoration parts d_i , respectively. The set of decorated intervals is denoted by $\overline{\mathbb{DIR}}$, and the set of decorated boxes with n components is denoted by $\overline{\mathbb{DIR}}^n$.

For readability, the decorations attached to intervals $\mathbf{u}, \mathbf{v}, \dots$ are often named du, dv, \dots , and each resulting decorated interval named by the corresponding uppercase letter, e.g., $\mathbf{U} = (\mathbf{u}, du) = \mathbf{u}_{du}$, etc. The d ’s here have no connection with differentials in calculus.

An interval or decoration may be called a **bare** interval or decoration to emphasize that it is not a decorated interval.

Treating the decorations as sets as in (9), \mathbf{trv} is the set of all (f, \mathbf{x}) pairs, and the others are nonempty subsets of \mathbf{trv} . By design they satisfy the **exclusivity rule**

$$\text{For any two decorations, either one contains the other or they are disjoint.} \quad (11)$$

Namely the definitions (10) give:

$$\mathbf{com} \subset \mathbf{dac} \subset \mathbf{def} \subset \mathbf{trv} \supset \mathbf{ill}, \quad \text{note the change from } \subset \text{ to } \supset; \quad (12)$$

$$\mathbf{com}, \mathbf{dac} \text{ and } \mathbf{def} \text{ are disjoint from } \mathbf{ill}. \quad (13)$$

Property (11) implies that for any (f, \mathbf{x}) there is a unique tightest (in the containment order (12)), decoration such that $p_d(f, \mathbf{x})$ is true, called the **strongest decoration of (f, \mathbf{x})** , or of f over \mathbf{x} , and written $\text{Dec}(f | \mathbf{x})$. That is:

$$\text{Dec}(f | \mathbf{x}) = d \iff p_d(f, \mathbf{x}) \text{ holds, but } p_e(f, \mathbf{x}) \text{ fails for all } e \subset d. \quad (14)$$

NOTE—Like the exact range $\text{Rge}(f | \mathbf{x})$, the strongest decoration is theoretically well-defined, but its value for a particular f and \mathbf{x} might be impractically expensive to compute, or even undecidable.

11.3 The ill-formed interval

The \mathbf{ill} decoration results from invalid constructions, and propagates unconditionally through arithmetic expressions. Namely, if a constructor call does not return a valid decorated interval, it returns an ill-formed one (i.e., decorated with \mathbf{ill}); and the decorated interval result of a library arithmetic operation is ill-formed, if and only if one of its inputs is ill-formed. Formally, \mathbf{ill} may be identified with the property $\text{Dom}(f) = \emptyset$ of (f, \mathbf{x}) pairs.

An ill-formed decorated interval is also called NaI, **Not an Interval**. Except as described in the next paragraph, an implementation shall behave as if there is only one NaI, whose interval part has no value at Level 1.

Information may be stored in an NaI in an implementation-defined way (like the payload of an IEEE 754 NaN), and functions may be provided for a user to set and read this for diagnostic purposes. An implementation may provide means for an exception to be signaled when an NaI is produced.

[Example. The constructor call `numsToInterval(2, 1)` is invalid in this flavor, so its decorated version returns NaI.]

11.4 Permitted combinations

A decorated interval \mathbf{y}_{dy} shall always be such that $\mathbf{y} \supseteq \text{Rge}(f | \mathbf{x})$ and $p_{dy}(f, \mathbf{x})$ holds, for some (f, \mathbf{x}) as in 11.2—informally, it must tell the truth about some conceivable evaluation of a function over a box. If $dy = \mathbf{def}$, \mathbf{dac} or \mathbf{com} then by definition \mathbf{x} is nonempty, and f is everywhere defined on it, so that $\text{Rge}(f | \mathbf{x})$ is nonempty, implying \mathbf{y} is nonempty. Hence, these decorated intervals are contradictory: $\emptyset_{\mathbf{dac}}$ and $\emptyset_{\mathbf{def}}$; and $\mathbf{x}_{\mathbf{com}}$ if \mathbf{x} is empty or unbounded. Implementations shall not produce them.

No other combinations are essentially forbidden.

11.5 Operations on/with decorations

This subclause contains operations to initialize the decoration on a bare interval and to disassemble and reassemble a decorated interval; and comparisons for decorations.

11.5.1 Initializing

Correct use of decorations when evaluating an expression has two parts: correctly initialize the input intervals; and evaluate using decorated interval extensions of library operations.

The simplest expression with one argument x is the trivial expression “ x ” with no operations. It defines the identity function Id that maps a real x to itself, $\text{Id}(x) = x$. For interval-evaluation of this expression

over some bare interval \mathbf{x} , the appropriate initial decoration for \mathbf{x} is the strongest decoration d that makes $p_d(\text{Id}, \mathbf{x})$ true, that is

$$d = \text{Dec}(\text{Id} \mid \mathbf{x}) = \begin{cases} \text{com} & \text{if } \mathbf{x} \text{ is nonempty and bounded,} \\ \text{dac} & \text{if } \mathbf{x} \text{ is unbounded,} \\ \text{trv} & \text{if } \mathbf{x} \text{ is empty.} \end{cases}$$

The function **newDec()** constructs a decorated interval from a bare one by adding such an initial decoration:

$$\text{newDec}(\mathbf{x}) = \mathbf{x}_d \quad \text{where} \quad d = \text{Dec}(\text{Id} \mid \mathbf{x}). \quad (15)$$

Initializing each input thus, before evaluating an expression, provides the most informative decoration on the output.

11.5.2 Disassembling and assembling

For a decorated interval \mathbf{x}_{dx} , the operations **intervalPart**(\mathbf{x}_{dx}) and **decorationPart**(\mathbf{x}_{dx}) shall be provided, with value \mathbf{x} and dx , respectively. For the case of **NaI**, **decorationPart**(**NaI**) has the value **ill**, but **intervalPart**(**NaI**) has no value at Level 1.

Given an interval \mathbf{x} and a decoration dx , the operation **setDec**(\mathbf{x}, dx) returns the decorated interval \mathbf{x}_{dx} if this is an allowed combination. The cases of forbidden combinations are as follows:

- **setDec**(\emptyset, dx) where dx is one of **def**, **dac** or **com** returns \emptyset_{trv} .
- **setDec**(\mathbf{x}, com), for any unbounded \mathbf{x} , returns \mathbf{x}_{dac} .
- **setDec**(\mathbf{x}, ill) for any \mathbf{x} , whether empty or not, returns **NaI**.

NOTE—Careless use of the **setDec** function can negate the aims of the decoration system and lead to false conclusions that violate the FTDIA. It is provided for expert users, who might need it, e.g., to decorate the output of functions whose definition involves the **intersection** and **convexHull** operations.

11.5.3 Comparisons

For decorations, comparison operations for equality $=$ and its negation \neq shall be provided, as well as comparisons $>, <, \geq, \leq$ with respect to the propagation order (19).

11.6 Decorations and arithmetic operations

Given a scalar point function φ of k variables, a **decorated interval extension** of φ —denoted here by the same name φ —adds a decoration component to a bare interval extension of φ . It has the form $\mathbf{w}_{dw} = \varphi(\mathbf{v}_{dv})$, where $\mathbf{v}_{dv} = (\mathbf{v}, dv)$ is a k -component decorated box $((\mathbf{v}_1, dv_1), \dots, (\mathbf{v}_k, dv_k))$. By the definition of a bare interval extension, the interval part \mathbf{w} depends only on the input intervals \mathbf{v} ; the decoration part dw generally depends on both \mathbf{v} and dv . In this context, **NaI** is regarded as being \emptyset_{ill} .

The definition of a bare interval extension implies

$$\mathbf{w} \supseteq \text{Rge}(\varphi \mid \mathbf{v}), \quad (\text{enclosure}). \quad (16)$$

The decorated interval extension of φ determines a dv_0 such that

$$p_{dv_0}(\varphi, \mathbf{v}) \text{ holds,} \quad (\text{a “local decoration”}). \quad (17)$$

It then evaluates the output decoration dw by

$$dw = \min\{dv_0, dv_1, \dots, dv_k\}, \quad (\text{the “min-rule”}), \quad (18)$$

where the minimum is taken with respect to the **propagation order**:

$$\text{com} > \text{dac} > \text{def} > \text{trv} > \text{ill}. \quad (19)$$

NOTE 1—Because **NaI** is treated as \emptyset_{ill} , this definition implies (without treating it as a special case) that $\varphi(\mathbf{v}_{dv})$ is **NaI** if, and only if, some component of \mathbf{v}_{dv} is **NaI**.

NOTE 2—Let $f(z_1, \dots, z_n)$ be an expression defining a real point function $f(x_1, \dots, x_n)$. Then decorated interval evaluation of f on a correctly initialized input decorated box \mathbf{x}_{dx} gives a decorated interval \mathbf{y}_{dy} such that not only does one have

$$\mathbf{y} \supseteq \text{Rge}(f \mid \mathbf{x}) \quad (20)$$

but also

$$p_{dy}(f, \mathbf{x}) \text{ holds.} \quad (21)$$

For instance, if the computed dy equals **def**, then f is proven to be everywhere defined on the box \mathbf{x} . This is an instance of the FTDIA for this flavor—see Annex B for a rigorous statement and proof.

NOTE 3—In the same way as the enclosure requirement (16) is compatible with many bare interval extensions, typically coming from different interval types at Level 2, so there might be several dv_0 satisfying the local decoration requirement (17). The ideal choice is the strongest decoration d such that $p_d(\varphi, \mathbf{v})$ holds, that is to take

$$dv_0 = \text{Dec}(\varphi \mid \mathbf{v}). \quad (22)$$

This is easily computable in finite precision for the arithmetic operations in 10.5, 10.6. However, functions may be added to the library in the future for which (22) is impractical to compute for some arguments \mathbf{v} . Hence, the weaker requirement (17) is made.

11.7 Decoration of non-arithmetic operations

11.7.1 Interval-valued operations

These give interval results but are not interval extensions of point functions:

- the reverse-mode operations of 10.5.4;
- the cancellative operations **cancelPlus**(\mathbf{x}, \mathbf{y}) and **cancelMinus**(\mathbf{x}, \mathbf{y}) of 10.5.6;
- The set-oriented operations **intersection**(\mathbf{x}, \mathbf{y}) and **convexHull**(\mathbf{x}, \mathbf{y}) of 10.5.7.

No one way of decorating these operations gives useful information in all contexts. Therefore, a *trivial* decorated interval version is provided as follows. If any input is **NaI**, the result is **NaI**; otherwise the corresponding operation is applied to the interval parts of the inputs, and its result decorated with **trv**. The user may replace this by a nontrivial decoration via **setDec**(), see 11.5, where this can be deduced in a given application.

11.7.2 Non-interval-valued operations

These give non-interval results:

- the numeric functions of 10.5.9;
- the boolean-valued functions of 10.5.10;
- the overlap function of 10.6.4.

For each such operation, if any input is **NaI**, the result has no value at Level 1. Otherwise, the operation acts on decorated intervals by discarding the decoration and applying the corresponding bare interval operation.

11.8 User-supplied functions

A user may define a decorated interval extension of some point function, as defined in 11.6, to be used within expressions as if it were a library operation.

[Examples.

a) In an application, an interval extension of the function

$$f(x) = x + 1/x$$

was required. Evaluated as written, it gives unnecessarily pessimistic enclosures: e.g., with $\mathbf{x} = [\frac{1}{2}, 2]$, one obtains

$$f(\mathbf{x}) = [\frac{1}{2}, 2] + 1/[\frac{1}{2}, 2] = [\frac{1}{2}, 2] + [\frac{1}{2}, 2] = [1, 4],$$

much wider than $\text{Rge}(f \mid \mathbf{x}) = [2, 2\frac{1}{2}]$.

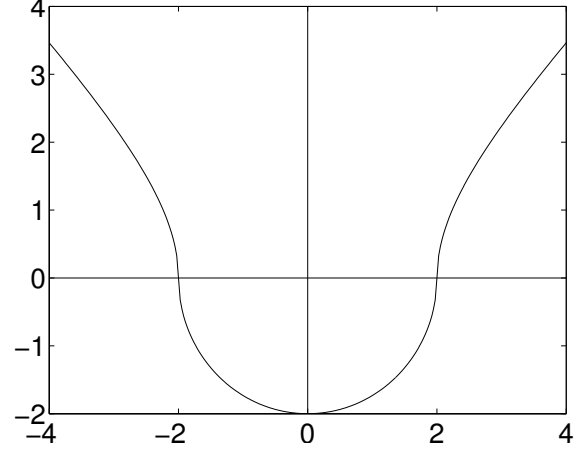
Thus it is useful to code a tight interval extension by special methods, e.g., monotonicity arguments, and to provide this as a new library function. Suppose this has been done. To convert it to a decorated interval extension just entails adding code to provide a local decoration and combine this with the input decoration by the min-rule (18). In this case, it is straightforward to compute the strongest local decoration $d = \text{Dec}(f \mid x)$, as follows.

$$d = \begin{cases} \text{com} & \text{if } 0 \notin x \text{ and } x \text{ is nonempty and bounded,} \\ \text{dac} & \text{if } 0 \notin x \text{ and } x \text{ is unbounded,} \\ \text{trv} & \text{if } 0 \in x \text{ or } x \text{ is empty.} \end{cases}$$

b)

The next example shows how an expert might manipulate decorations explicitly to give a function, defined piecewise by different formulas in different regions of its domain, the best possible decoration. Suppose that

$$f(x) = \begin{cases} f_1(x) := \sqrt{x^2 - 4} & \text{if } |x| > 2, \\ f_2(x) := -\sqrt{4 - x^2} & \text{otherwise,} \end{cases}$$



where $:=$ means “defined as,” see the diagram.

The function consists of three pieces on regions $x \leq -2$, $-2 \leq x \leq 2$ and $x \geq 2$ that join continuously at region boundaries, but the standard gives no way to determine this continuity, at run time or otherwise. However, a user-defined decorated interval function as defined below provides the best possible decorations.

```
function  $y_{dy} = f(x_{dx})$ 
 $u = f_1(x \cap [-\infty, -2])$ 
 $v = f_2(x \cap [-2, 2])$ 
 $w = f_1(x \cap [2, +\infty])$ 
 $y = u \cup v \cup w$ 
 $dy = dx$ 
```

Here \cup denotes the convexHull operation. The user’s knowledge that f is everywhere defined and continuous is expressed by the statement $dy = dx$, propagating the input decoration unchanged. f , thus defined, can safely be used within a larger decorated interval evaluation.

]

11.9 Notes on the com decoration

NOTE—The force of **com** is the Level 2 property that the *computed* interval $f(x)$ is bounded. Equivalently, overflow did not occur, where overflow has the generalized meaning that a finite-precision operation could not enclose a mathematically bounded result in a bounded interval of the required output type. Briefly, for a single operation, “**com** is **dac** plus bounded inputs and no overflow.”

Thus the result of interval-evaluating an arithmetic expression in finite precision is decorated **com** if and only if the evaluation is *common at Level 2*, meaning: each input that affects the result is nonempty and bounded, and each individual operation that affects the result is everywhere defined and continuous on its inputs and does not overflow.

A tempting alternative is to make **com** record whether the evaluation is *common at Level 1*, meaning that all the relevant intervals are mathematically bounded, even if overflow occurred in finite precision. E.g., one might drop the “bounded inputs” requirement and require “mathematically bounded” instead of “actually bounded” on the output of an operation.

However, the **dac** decoration already provides such information, and the suggested change gives nothing extra. Namely, if the inputs \mathbf{x} to $f(\mathbf{x})$ are bounded, and the output decoration is **dac**, it follows, from the fact that a continuous function on a compact set is bounded, that the point function f is mathematically bounded on \mathbf{x} , and all its individual operations are mathematically bounded on their inputs even if overflow might have occurred in finite precision.

For example, consider $f(x) = 1/(2x)$ evaluated at $\mathbf{x} = [1, M]$ using an inf-sup type where M is the largest representable real. This gives

$$\mathbf{y}_{dy} = f(\mathbf{x}_{com}) = 1/(2 * [1, M]_{com}) = 1/[2, +\infty]_{dac} = [0, \frac{1}{2}]_{dac}.$$

Despite the overflow, one can deduce from the final **dac** that the result of the multiplication was mathematically bounded.

This might be of limited use: consider $g(x) = 1/f(x) = 1/(1/(2x))$, evaluated at the same $\mathbf{x} = [1, M]$ giving \mathbf{z}_{dz} . The standard has no way to record that the lower bound of \mathbf{y} is mathematically positive, i.e., $1/(2M)$. Thus the Level 2 result is $\mathbf{z}_{dz} = [2, +\infty]_{trv}$, compare $[2, 2M]_{com}$ at Level 1.

11.10 Compressed arithmetic with a threshold (optional)

11.10.1 Motivation

The **compressed decorated interval arithmetic** (compressed arithmetic for short) described here lets experienced users obtain more efficient execution in applications where the use of decorations is limited to the context described below. An implementation need not provide it; if it does so, the behavior described in this subclause is required. Which compressed interval types are provided is implementation-defined.

Each Level 2 instance of compressed arithmetic is based on a supported Level 2 bare interval type \mathbb{T} , but is a distinct “compressed type,” with its own datums and library of operations.

The context is that of evaluating an arithmetic expression, where the use made of a decorated interval evaluation $\mathbf{y}_{dy} = f(\mathbf{x}_{dx})$ depends on a check of the result decoration dy against an application-dependent **threshold** τ , where $\tau \geq \mathbf{trv}$ in the propagation order (19):

$dy \geq \tau$: represents normal computation. The decoration is not used, but one exploits the range enclosure given by the interval part and the knowledge that dy remained $\geq \tau$.

$dy < \tau$: declares an exceptional condition to have occurred. The interval part is not used, but one exploits the information given by the decoration.

11.10.2 Compressed interval types

For such uses, one needs to record an interval’s value, or its decoration, but never both at once. The **compressed type** of threshold τ , **associated with** \mathbb{T} , is the type each of whose datums is either a (bare) \mathbb{T} -interval or a decoration less than τ . It is denoted \mathbb{T}_τ . Two such types are the same if and only if they have the same \mathbb{T} and the same τ . A \mathbb{T}_τ datum can be any \mathbb{T} datum or any decoration except that:

- Only decorations $< \tau$ occur; in particular **com** is never used.
- The empty interval \emptyset is replaced by—equivalently, is regarded by the implementation as being—a new decoration **emp** added to the table in (10), whose defining property is

Value	Short description	Property	Definition
emp	empty	$p_{emp}(f, \mathbf{x})$	$\mathbf{x} \cap \text{Dom}(f)$ is empty;

(23)

emp lies between **trv** and **ill** in the containment order (12) and the propagation order (19):

$$\begin{aligned} \mathbf{com} &\subset \mathbf{dac} \subset \mathbf{def} \subset \mathbf{trv} \supset \mathbf{emp} \supset \mathbf{ill}, \\ \mathbf{com} &> \mathbf{dac} > \mathbf{def} > \mathbf{trv} > \mathbf{emp} > \mathbf{ill}. \end{aligned} \tag{24}$$

Since $\tau \geq \mathbf{trv}$, it is always true that **emp** $< \tau$, which means that as soon as an empty result is produced while evaluating an expression, the $dy < \tau$ case has occurred.

NOTE—The reason for treating \emptyset as a decoration $< \tau$ is that obtaining an empty result (e.g., by doing something like $\sqrt{[-2, -1]}$ while evaluating a function) is one of the exceptional conditions that compressed interval computation should detect.

The only way to use compressed arithmetic with a threshold τ is to construct \mathbb{T}_τ datums. Conversion between compressed types, say from a \mathbb{T}_τ -interval to a \mathbb{T}'_τ -interval, shall be equivalent to converting first to a normal decorated interval by `normalInterval()`, then between decorated interval types if $\mathbb{T} \neq \mathbb{T}'$, and finally to the output type by `τ' -compressedInterval()`.

NOTE—Since, for any practical interval type \mathbb{T} , a decoration fits into less space than an interval, one can implement arithmetic on compressed interval datums that take up the same space as a bare interval of that type. For instance, if \mathbb{T} is the IEEE 754 `binary64` inf-sup type, a compressed interval uses 16 bytes, the same as a bare \mathbb{T} -interval; a full decorated \mathbb{T} -interval needs at least 17 bytes.

Because compressed intervals must behave exactly like bare intervals as long as one does not fall below the threshold, and take up the same space, there is no room to encode τ as part of the interval's value.

11.10.3 Operations

The enquiry function `isInterval(x)` returns true if the compressed interval x is an interval, false if it is a decoration.

The constructor `τ -compressedInterval()` is provided for each threshold value τ . The result of `τ -compressedInterval(X)`, where $X = x_{dx}$ is a decorated \mathbb{T} -interval, is a \mathbb{T}_τ -interval as follows:

```
if  $dx \geq \tau$ , return the  $\mathbb{T}_\tau$ -interval with value  $x$ 
else return the  $\mathbb{T}_\tau$ -interval with value  $dx$ .
```

`τ -compressedInterval(x)` for a bare interval x is equivalent to `τ -compressedInterval(newDec(x))`.

The function `normalInterval(x)` converts a \mathbb{T}_τ -interval to a decorated interval of the parent type, as follows:

```
if  $x$  is an interval, return  $x_\tau$ 
if  $x$  is a decoration  $d$ 
    if  $d = \text{ill}$ , return Emptyill = NaN
    elseif  $d = \text{emp}$ , return Emptytrv
    else return Entired
```

Arithmetic operations on compressed intervals shall follow *worst case semantics* rules that treat a decoration in `{trv, def, dac}` as representing a set of decorated intervals, and are necessary if the fundamental theorem is to remain valid. Namely, inputs to each operation behave as follows:

- Operations purely on bare intervals are performed as if each x is the decorated interval x_τ , resulting in a decorated interval y_{dy} that is then converted back into a compressed interval. If $dy < \tau$, the result is the decoration dy , otherwise the bare interval y .
- For operations with at least one decoration input, the result is always a decoration. A bare interval input is treated as in the previous item. A decoration d in `{emp, ill}` is treated as \emptyset_d . A decoration d in `{trv, def, dac}` is treated (conceptually) as x_d with an arbitrary nonempty interval x . The decoration `com` cannot occur. Performing the resulting decorated interval operation on all such possible inputs leads to a set of all possible results y_{dy} . The tightest decoration (in the containment order (24)) enclosing all resulting dy is returned.

As a result, each operation returns an actual or implied decoration compatible with its input, so that in an extended evaluation, the final decoration using compressed arithmetic is never stronger than that produced by full decorated interval arithmetic.

[Example. Assuming $\tau > \text{def}$,

- The division `def/[1, 2]` becomes `$x_{\text{def}}/[1, 2]_\tau$` with arbitrary nonempty interval x . The result is always decorated `def`, so returns `def`.
- But `[1, 2]/def` becomes `$[1, 2]_\tau/x_{\text{def}}$` with arbitrary nonempty interval x . The result can be decorated `def`, `trv` or `emp`, so returns the tightest decoration containing these, namely `trv`.

/

Since there are only a few decorations, one can prepare complete operation tables according to this rule, and only these tables need to be implemented.

If compressed arithmetic is implemented, it shall provide versions of all the required operations of 10.5, and it should provide the recommended operations of 10.6.

12. Level 2 description

12.1 Level 2 introduction

Entities and operations at Level 2 are said to have **finite precision**. From them, implementable interval algorithms may be constructed. Level 2 entities are called **datums**¹². Since the standard deals with numeric functions of intervals (such as the midpoint) and interval functions of numbers (such as the construction of an interval from its lower and upper bounds), this clause involves both numeric and interval datums, as well as decoration, string and boolean datums.

12.1.1 Types and formats

Following IEEE 754 terminology, numeric (usually but not necessarily floating-point) datums are organized into **formats**. Interval datums are organized into **types**. Each format or type is a finite set of datums, with associated operations. If \mathbb{F} denotes a format, an **F-number** means a member of \mathbb{F} , possibly infinite but not the “not a number” value NaN; to allow it to be NaN, it is called an **F-datum**. If \mathbb{T} denotes a bare or decorated interval type, a **T-interval** means a member of \mathbb{T} , possibly empty but (in the decorated case) not the “not an interval” value NaI; to allow it to be NaI, it is called a **T-datum**. A **T-box** means a box with T-datum components—equivalently, with T-interval components, if \mathbb{T} is a bare interval type.

The standard defines three kinds of interval type:

- **Bare interval types**, see 12.5, are named finite sets of (mathematical, Level 1) intervals.
- **Decorated interval types**, also see 12.5, are named finite sets of decorated intervals.
- **Compressed interval types** (optional) are named finite sets of compressed intervals. They are described in 11.10, and Level 2 makes no further requirements on them.

For each bare interval type there shall be a corresponding **derived** decorated interval type, and each decorated interval type shall be derived from a bare interval type, see 12.5.1. An implementation shall provide at least one supported bare interval type. If 754-conforming, it shall provide the inf-sup type, see 12.5.2, of at least one of the five basic formats in 3.3 of IEEE Std 754-2008. Beyond this, which types are provided is language- or implementation-defined.

It is language- or implementation-defined whether the format or type of a datum can be determined at run time.

12.1.2 Operations

A Level 2 operation is a finite-precision approximation to the corresponding Level 1 operation. To describe the required functionality, the standard treats each Level 1 operation as having a number of Level 2 **versions** in which each interval input or output is given a specific interval type, and each numeric input or output is given a specific numeric format.

NOTE—An implementation might provide the functionality via differently named operations for each version, or by overloading a single operation name, or by a single operation that accepts all the required type/format combinations, or in other ways.

For example, let \mathbb{T}_1 and \mathbb{T}_2 be two supported interval types. The standard requires a version of addition $z = x + y$ where each of x, y, z has type \mathbb{T}_1 , and another where each of them has type \mathbb{T}_2 . An implementation might provide

¹²Not “data,” whose common meaning could cause confusion.

this functionality by having two separate operations; another might have a single operation that takes inputs of types \mathbb{T}_1 and \mathbb{T}_2 in any combination, with some rule to determine the type of the output z ; etc.

The term **\mathbb{T} -version** of a Level 1 operation denotes one in which any input or output that is an interval, is a \mathbb{T} -datum. For bare interval types this includes the following:

- a) A \mathbb{T} -interval extension (12.9) of one of the required or recommended arithmetic operations of 10.5, 10.6.
- b) A set operation, such as intersection and convex hull of \mathbb{T} -intervals, returning a \mathbb{T} -interval.
- c) A function such as the midpoint, whose input is a \mathbb{T} -interval and output is numeric.
- d) A constructor, whose input is numeric or text and output is a \mathbb{T} -datum.
- e) The operations of 13.3, 13.4, whose input is a \mathbb{T} -interval and output is a string.

Additionally, a version of `convertType` with output of type \mathbb{T} , see 12.12.10, is a \mathbb{T} -version irrespective of the type of the input interval. Generically these comprise the **operations of the type \mathbb{T}** , for the implementation.

12.1.3 Exception behavior

For some operations, and some particular inputs, there might not be a valid result. At Level 1 there are several cases when no value exists. However, a Level 2 operation always returns a value. When the Level 1 result does not exist, the operation returns either

- a special value indicating this event (e.g., NaN for most of the numeric functions in 12.12.8); or
- a value considered reasonable in practice. For example, `mid(Entire)` returns 0; a constructor given invalid input returns `Empty`; and one of the comparisons of 12.12.9, if any input is NaI, returns `false`.

The standard defines the following **exceptions** that may be signaled, causing a handler to be invoked:

- For the interval constructors of 12.12.7, and for `exactToInterval` in 13.4, it is possible that at Level 2, using finite precision, the implementation cannot decide whether a Level 1 value exists or not, see *Difficulties in implementation* in 12.12.7. If the constructor is certain that no Level 1 value exists, the exception `UndefinedOperation` is signaled; if it cannot decide, the exception `PossiblyUndefinedOperation` is signaled.
- If `intervalPart()` is called with NaI as input, the exception `IntvlPartOfNaI` is signaled, see 12.12.11.
- If some inputs of an operation are invalid, the exception `InvalidOperand` might be signaled, see 14.3.

These exceptions are separate from the exceptional conditions handled by the decoration system. The mechanism of how a handler deals with an exception is language- or implementation-defined. An implementation may provide exception handling additional to that above.

12.2 Naming conventions for operations

An operation is generally given a name that suits the context. For example, the addition of two interval datums x, y might be written in generic algebra notation $x + y$; or with a generic text name `add(x, y)`; or giving full type information such as *infsup-decimal64-add(x, y)*. It might also be written as \mathbb{T} -`add(x, y)` to show it is an operation of a particular but unspecified type \mathbb{T} .

In a specific language or programming environment, the names used for types might differ from those used in this document.

12.3 Tagging, and the meaning of equality at Level 2

A Level 2 format or type is an abstraction of a particular way to represent numbers or intervals—e.g., “IEEE 64 bit binary floating-point” for numbers—focusing on the Level 1 entities denoted, and hiding the Level 3 representation.

However, a datum is more than just the Level 1 value: for instance, the number 3.75 represented in 32 bit binary floating-point is a different datum from the same number represented in 64 bit binary floating-point (“single” and “double” precision respectively in typical implementations).

This is achieved by formally regarding each datum as a pair:

$$\begin{aligned}\text{number datum} &= (\text{Level 1 number, format name}), \\ \text{bare interval datum} &= (\text{Level 1 interval, type name}),\end{aligned}$$

where the name is a symbol that uniquely identifies the format or type. Since a decorated interval combines a bare interval and a decoration it thus becomes a triple at Level 2:

$$\text{decorated interval datum} = (\text{Level 1 interval, type name, decoration}).$$

The Level 1 value is said to be **tagged** by the format or type name. It follows that distinct formats or types are disjoint sets. By convention, such names are omitted from datums except when clarity requires.

[Example. Level 2 interval addition within a type named t is normally written $z = x + y$, though the full correct form is $(z, t) = (x, t) + (y, t)$. The full form might be used, for instance, to indicate that mixed-type addition is forbidden between types s and t but allowed between types s and u . Namely, one can say that $(x, s) + (y, t)$ is undefined, but $(x, s) + (y, u)$ is defined.]

The interval comparison operations of 10.5.10, including comparison for equality, are provided between datums x, y of the same type. Additionally they are provided between datums of different types provided the types are *comparable*, see 12.5.1.

Therefore it is necessary to distinguish kinds of equality. x and y are **identical** datums if they have the same Level 1 value and the same type. If their types are comparable then **equal**(x, y) is defined for them; if they have the same Level 1 value, **equal**(x, y) returns **true** and they are called equal. If their types are not comparable, **equal**(x, y) is undefined; they are not equal even if they have the same Level 1 value.

NOTE—This is like the situation for IEEE 754 floating-point numbers. For instance, the number 3.75 is representable exactly by datums x, y, z in **binary32**, **binary64** and **decimal64**, respectively. They are non-identical datums; but the comparison $x = y$ (equivalently **compareQuietEqual**(x, y)) is defined since x and y have the same radix, and returns **true** because they have the same Level 1 value. However, $x = z$ is not defined within IEEE Std 754-2008, because x has a different radix from z .

Similarly, let x, y and z be the datums in the inf-sup types of **binary32**, **binary64** and **decimal64**, respectively, for an interval that they all represent exactly, such as $[1, 3.75]$. They are non-identical datums; but the comparison $x = y$ (equivalently **equal**(x, y)) is defined and returns **true**; while $x = z$ is not defined in this standard, though x and z have the same Level 1 value, because their types are not comparable.

For a decorated interval type \mathbb{T} , the unique NaI datum is equal to itself as a datum, but compares unequal to any \mathbb{T} -datum, including itself, with the **equal** relation. This follows the behavior of NaN among IEEE 754 floating-point datums.

12.4 Number formats

In view of 12.3, a **number format**, or just format, is formally the set of all pairs (x, f) such that x belongs to a given finite set \mathbb{F} of numbers and symbols, and f is a name for the format. A format is **provided** if the implementation provides a representation of it as in 14.2, and is **supported** (by this standard) if in addition:

- \mathbb{F} comprises a subset of the extended reals $\overline{\mathbb{R}}$ that includes $-\infty$ and $+\infty$, together with a value NaN, and optionally signed zeros -0 and $+0$.
- \mathbb{F} contains 0, or contains both -0 and $+0$, or contains all three of these values.
- \mathbb{F} contains at least one nonzero finite number.
- \mathbb{F} is symmetric: if x is in \mathbb{F} , so is its additive inverse $-x$, where ± 0 are additive inverses of each other, as are $\pm\infty$.

Following the convention of omitting names, the format is normally identified with the set \mathbb{F} , and one may say a supported format is a set comprising NaN together with a finite subset of \mathbb{R} , and possibly ± 0 , subject to the above rules. A member of \mathbb{F} is called an **\mathbb{F} -datum**. Every member except NaN is called **numeric**, or an **\mathbb{F} -number**.

As an aid to notation, if x is an \mathbb{F} -number then $\text{Val}(x)$, the (extended-real) **value** of x , is the member of $\overline{\mathbb{R}}$ that x denotes: $\text{Val}(-0) = \text{Val}(+0) = 0$, $\text{Val}(x) = x$ for other \mathbb{F} -numbers, and $\text{Val}(\text{NaN})$ is undefined. Thus $\text{Val}(\mathbb{F})$ is the whole set $\{\text{Val}(x) \mid x \in \mathbb{F}\}$ of extended reals denoted by \mathbb{F} .

A floating-point format in the IEEE 754 sense, such as **binary64**, is identified with the number format for which \mathbb{F} is the set of datums representable in that format. At Level 2, different kinds of NaN map to the unique NaN of the number format.

In this document the five basic formats in 3.3 of IEEE Std 754-2008 are named **binary32**, **binary64**, **binary128**, **decimal64**, **decimal128**. Abbreviated names such as **b64** instead of **binary64** are sometimes used, and refer to the same format.

A number format \mathbb{F} is said to be **compatible** with an interval type \mathbb{T} if each non-empty \mathbb{T} -interval contains at least one finite \mathbb{F} -number.

For a Level 2 operation, an input or output that in the corresponding Level 1 operation is an extended real becomes a datum of some format \mathbb{F} . On input, if this datum x is not NaN, it is replaced by $\text{Val}(x)$ before applying the Level 1 operation; if NaN, it is handled by rules specific to the operation. On output, an extended-real Level 1 result is mapped (**rounded**) to an \mathbb{F} -datum according to rules specific to the operation (see 12.12.8).

12.5 Bare and decorated interval types

12.5.1 Definition In view of 12.3, a **bare interval type**, or just type, is formally the set of all pairs (x, t) such that x belongs to a given finite subset \mathbb{T} of the mathematical intervals \mathbb{IR} , and t is a name for the type. It is **provided** if the implementation provides a representation of it as in 14.2, and is **supported** (by this standard) if in addition:

- \mathbb{T} contains Empty and Entire.
- \mathbb{T} is symmetric: if x is in \mathbb{T} , so is $-x$.

The **decorated interval type derived** from \mathbb{T} is formally the set of triples (x, t, d) such that (x, t) is a \mathbb{T} -interval, and $d \in \mathbb{D}$ is a decoration that follows the rule for permitted combinations (x, d) in 11.4.

Following the convention of omitting names, the type is normally regarded as being the set \mathbb{T} , and one may say a bare interval type is an arbitrary set of intervals subject to the above rules. The derived decorated type is then regarded as a set of pairs (x, d) , equivalently x_d , where $x \in \mathbb{T}$ and $d \in \mathbb{D}$.

Following IEEE Std 754-2008’s terminology for formats (Definition 2.1.36), a type \mathbb{T}' is **wider**¹³ than a type \mathbb{T} (and \mathbb{T} is **narrower** than \mathbb{T}') if \mathbb{T} is a subset of \mathbb{T}' when they are regarded as sets of Level 1 intervals, ignoring the type tags and possible decorations. Two types are **comparable** if either is wider than the other. [Example. The basic 754-conforming types of a given radix are comparable, see 12.6.]

Each decorated interval type shall contain a “Not an Interval” datum NaN, identified with (\emptyset, inl) . It shall appear to be unique at Level 2, but non-Level-2 operations may be provided to set and get a payload in an NaN for diagnostic purposes, in an implementation-defined way (see 11.3).

[Example. To illustrate the flexibility allowed in defining types, let S_1 and S_2 be the sets of inf-sup intervals using IEEE 754 single (binary32) and double (binary64) precision, respectively. That is, a member of S_1 [respectively S_2] is either empty, or an interval whose bounds are exactly representable in binary32 [respectively binary64].

An implementation usually would define these as different bare interval types, by tagging members of S_1 by one type name t_1 and members of S_2 by another name t_2 —represented at Level 3 by a pair of binary32 or of binary64 numbers,

¹³Wider means having more precision. In the IEEE 754 context, for a given radix, a wider format is one with a wider bit string for the exponent and/or significand in its Level 4 encoding.

respectively. However, it might treat them as one type, with the representation by a pair of `binary32`'s being a space-saving alternative to the pair of `binary64`'s, to be used, say, for some large arrays. The resulting Level 1 intervals are exactly the same those of `inf-sup binary64`, so this is just another way to store the latter type; but an implementation would give it a different name, to reflect the different storage and hence different operations at the code level.]

12.5.2 Inf-sup and mid-rad types

The **inf-sup type** derived from a supported number format \mathbb{F} (the type **inf-sup** \mathbb{F} , e.g., “inf-sup `binary64`”) is the bare interval type \mathbb{T} comprising all intervals whose bounds are in $\text{Val}(\mathbb{F})$, together with `Empty`. When \mathbb{F} is an IEEE 754 format, the **radix** of \mathbb{T} means the radix of \mathbb{F} .

NOTE—This implies `Entire` is in \mathbb{T} because $\pm\infty \in \mathbb{F}$ by the definition of a supported format, see 12.4, so \mathbb{T} satisfies the requirements for a bare interval type given in 12.5.1.

A **mid-rad** bare interval type is one whose nonempty bounded intervals comprise all intervals of the form $[m - r, m + r]$, where $m \in \text{Val}(\mathbb{F})$ for some number format \mathbb{F} , and $r \in \text{Val}(\mathbb{F}')$ for a possibly different format \mathbb{F}' , with m, r finite and $r \geq 0$. From the definition in 12.5.1 such a type shall contain `Empty` and `Entire` (so at Level 3 it shall have representations of these). It may also contain semi-bounded intervals.

12.6 754-conformance

The standard defines the notion of 754-conformance, whose stronger requirements improve accuracy and programming convenience.

12.6.1 Definition A **754-conforming type** is an inf-sup type derived from an IEEE 754 floating-point format (one of the five basic formats or an extended precision or extendable precision format) in the sense of 12.5.2 that meets the general requirements for conformance and whose operations meet the accuracy requirements in 12.10.

A **754-conforming implementation** is an implementation, all of whose types are 754-conforming, and whose operations meet the requirements for mixed-type arithmetic described in the next paragraphs. It might be a conforming part of an implementation in the sense of 4.1.

12.6.2 754-conforming mixed-type operations

IEEE Std 754-2008 requires a conforming floating-point system to provide mixed-format *formatOf* operations, where the output format is specified and the inputs may be of any format of the same radix as the output. The result is computed as if using the exact inputs and rounded to the required accuracy on output.

A 754-conforming implementation shall provide corresponding mixed-type interval operations. Namely, if it provides types $\mathbb{T}, \mathbb{T}_1, \mathbb{T}_2, \dots$ derived from IEEE 754 formats $\mathbb{F}, \mathbb{F}_1, \mathbb{F}_2, \dots$ all of the same radix, then for each *formatOf* operation with output format \mathbb{F} and accepting input formats chosen from $\mathbb{F}_1, \mathbb{F}_2, \dots$ there shall be an interval version of that operation with output type \mathbb{T} and input types chosen from $\mathbb{T}_1, \mathbb{T}_2, \dots$. The result shall be computed as if using the exact inputs and shall meet the accuracy requirement for that operation, specified in 12.10.

NOTE—For inf-sup types this requirement may be met by implicit widening of inputs to a common widest format, as double rounding is not an issue for directed rounding.

12.7 Multi-precision interval types

Multi-precision systems—extendable precision in IEEE 754 terminology—generally provide an (at least conceptually) infinite sequence of levels of precision, where there is a finite set \mathbb{F}_n of numbers representable at the n th level ($n = 1, 2, 3, \dots$), and $\mathbb{F}_1 \subset \mathbb{F}_2 \subset \mathbb{F}_3 \dots$. These are typically used to define a corresponding infinite sequence of interval types \mathbb{T}_n with $\mathbb{T}_1 \subset \mathbb{T}_2 \subset \mathbb{T}_3 \dots$.

[Example. For multi-precision systems that define a nonempty \mathbb{T}_n -interval to be one whose bounds are \mathbb{F}_n -numbers, each \mathbb{T}_n is an inf-sup type with a unique interval hull operation—explicit, in the sense of 12.8.]

A conforming implementation defines such \mathbb{T}_n as a *parameterized sequence* of interval types. It cannot take the union over n of the sets \mathbb{T}_n as a *single* type, because this infinite set has no interval hull operation: there

is generally no tightest member of it enclosing a given set of real numbers. This constrains the design of conforming multi-precision interval systems.

12.8 Explicit and implicit types, and Level 2 hull operation

12.8.1 Hull in one dimension

For each bare interval type \mathbb{T} there shall be defined an **interval hull** operation

$$\mathbf{y} = \text{hull}_{\mathbb{T}}(\mathbf{s}),$$

also called the \mathbb{T} -hull, which is part of \mathbb{T} 's mathematical definition. For an implicit type, see below, the implementation's documentation shall specify the \mathbb{T} -hull, e.g., by an algorithm. For an explicit type, it is uniquely determined and need not be separately specified.

NOTE—An implementation provides $\text{hull}_{\mathbb{T}}$ as the operation **convertType** for conversion between any two supported types, see 12.12.10.

The \mathbb{T} -hull maps an arbitrary set of reals, \mathbf{s} , to a minimal \mathbb{T} -interval \mathbf{y} enclosing \mathbf{s} . Minimal is in the sense that

$$\mathbf{s} \subseteq \mathbf{y}, \text{ and for any other } \mathbb{T}\text{-interval } \mathbf{z}, \text{ if } \mathbf{s} \subseteq \mathbf{z} \subseteq \mathbf{y} \text{ then } \mathbf{z} = \mathbf{y}.$$

Since \mathbb{T} is a finite set and contains Entire, such a minimal \mathbf{y} exists for any \mathbf{s} . In general \mathbf{y} might not be unique. If it is unique for every subset \mathbf{s} of \mathbb{R} , then the type \mathbb{T} is called **explicit**, otherwise it is **implicit**.

Two types with different hull operations are different, even if they have the same set of intervals.

The \mathbb{T} -hull operation **overflows** if it can only find an unbounded \mathbb{T} -interval \mathbf{y} to enclose a bounded set \mathbf{s} . Similarly, any Level 2 interval-valued operation overflows when its Level 1 result is bounded but too large to be enclosed in a bounded interval of the output type.

[Examples. Every inf-sup type is explicit. A mid-rad type is typically implicit.]

As an example of the need for a specified hull algorithm, let \mathbb{T} be the mid-rad type (12.5.2), where m and r use the same floating-point format \mathbb{F} , say binary64, and let \mathbf{s} be the interval $[-1, 1 + \epsilon]$, where $1 + \epsilon$ is the next \mathbb{F} -number above 1. Clearly any minimal interval (m, r) enclosing \mathbf{s} has $r = 1 + \epsilon$, but m can be any of the many \mathbb{F} -numbers in the range 0 to ϵ ; each of these gives a minimal enclosure of \mathbf{s} .

A possible general algorithm, for a bounded set \mathbf{s} and a mid-rad type, is to choose $m \in \mathbb{F}$ as close as possible to the mathematical midpoint of the interval $[\underline{s}, \bar{s}] = [\inf \mathbf{s}, \sup \mathbf{s}]$ (with some way to resolve ties) and then the smallest $r \in \mathbb{F}'$ such that $r \geq \max(m - \underline{s}, \bar{s} - m)$. The cost of performing this depends on how the set \mathbf{s} is represented. If \mathbf{s} is a binary64 inf-sup interval, it is simple. If \mathbf{s} is defined as the range of a function, it might be expensive.]

12.8.2 Hull in several dimensions

In n dimensions the \mathbb{T} -hull, as defined mathematically in 12.8.1, is extended to act componentwise, namely for an arbitrary subset \mathbf{s} of \mathbb{R}^n it is $\text{hull}_{\mathbb{T}}(\mathbf{s}) = (\mathbf{y}_1, \dots, \mathbf{y}_n)$ where

$$\mathbf{y}_i = \text{hull}_{\mathbb{T}}(\mathbf{s}_i),$$

and $\mathbf{s}_i = \{s_i \mid s \in \mathbf{s}\}$ is the projection of \mathbf{s} on the i th coordinate dimension. It is easily seen that this is a minimal \mathbb{T} -box containing \mathbf{s} , and that if \mathbb{T} is explicit it equals the unique tightest \mathbb{T} -box containing \mathbf{s} .

12.9 Level 2 interval extensions

Let \mathbb{T} be a bare interval type and f an n -variable scalar point function. A **\mathbb{T} -interval extension** of f , also called a **\mathbb{T} -version** of f , is a mapping \mathbf{f} from n -dimensional \mathbb{T} -boxes to \mathbb{T} -intervals, that is $\mathbf{f} : \mathbb{T}^n \rightarrow \mathbb{T}$, such that $f(x) \in \mathbf{f}(\mathbf{x})$ whenever $x \in \mathbf{x}$ and $f(x)$ is defined. Equivalently

$$\mathbf{f}(\mathbf{x}) \supseteq \text{Rge}(f \mid \mathbf{x}) \tag{25}$$

for any \mathbb{T} -box $\mathbf{x} \in \mathbb{T}^n$, regarding \mathbf{x} as a subset of \mathbb{R}^n . Generically, such mappings are called Level 2 interval extensions.

Though only defined over a finite set of boxes, a Level 2 extension of f is equivalent to a full Level 1 extension of f (10.4) so that this document does not distinguish between Level 2 and Level 1 extensions. Namely define f^* by

$$f^*(s) = f(\text{hull}_{\mathbb{T}}(s))$$

for any subset s of \mathbb{R}^n . Then the interval $f^*(s)$ contains $\text{Rge}(f|s)$ for any s , making f^* a Level 1 extension; and $f^*(s)$ equals $f(s)$ whenever s is a \mathbb{T} -box.

12.10 Accuracy of operations

This subclause describes requirements and recommendations on the accuracy of operations. Here, *operation* denotes any Level 2 version, provided by the implementation, of a Level 1 operation with interval output and at least one interval input. Bare interval operations are described; the accuracy of a decorated operation is defined to be that of its interval part.

12.10.1 Measures of accuracy

Three **accuracy modes** are defined that indicate the quality of interval enclosure achieved by an operation: *tightest*, *accurate* and *valid* in order from strongest to weakest. Each mode is in the first instance a property of an individual evaluation of an operation f of type \mathbb{T} over an input box x . The term **tightness** means the strongest mode that holds uniformly for some set of evaluations.

[Example. For some one-argument function, an implementation might document the tightness of $f(x)$ as being *tightest* for all x contained in $[-10^{15}, 10^{15}]$ and at least *accurate* for all other x .]

The *tightest* and *valid* modes apply to all interval types and all operations. The *accurate* mode is defined only for inf-sup types (because it involves the **nextOut** function), and for interval forward and reverse arithmetic operations of 10.5.3, 10.5.4 (because it requires operations f that are monotone at Level 1: $u \subseteq v$ implies $f(u) \subseteq f(v)$).

Let f_{exact} denote the corresponding Level 1 operation. The weakest mode *valid* is just the property of enclosure:

$$f(x) \supseteq f_{\text{exact}}(x). \quad (26)$$

For a forward arithmetic operation, it is equivalent to (25).

The strongest mode *tightest* is the property that $f(x)$ equals $f_{\text{tightest}}(x)$, the \mathbb{T} -hull of the Level 1 result:

$$f_{\text{tightest}}(x) = \text{hull}_{\mathbb{T}}(f_{\text{exact}}(x)). \quad (27)$$

The intermediate mode *accurate* applies to an inf-sup type \mathbb{T} derived from a number format \mathbb{F} . It asserts that $f(x)$ is *valid*, (26), and is at most slightly wider than the result of applying the *tightest* version to a slightly wider input box:

$$f(x) \subseteq \text{nextOut}(f_{\text{tightest}}(\text{nextOut}(\text{hull}_{\mathbb{T}}(x)))). \quad (28)$$

Here the **nextOut** function is defined in terms of the functions

$$\text{nextUp}, \text{nextDown} : \text{Val}(\mathbb{F}) \rightarrow \text{Val}(\mathbb{F})$$

as follows. For $x \in \text{Val}(\mathbb{F})$, define **nextUp**(x) to be $+\infty$ if $x = +\infty$, and the least member of $\text{Val}(\mathbb{F})$ greater than x otherwise; since $+\infty \in \text{Val}(\mathbb{F})$, this is well-defined. Define **nextDown**(x) to be $-\text{nextUp}(-x)$. Then, if x is a \mathbb{T} -interval, define

$$\text{nextOut}(x) = \begin{cases} [\text{nextDown}(x), \text{nextUp}(x)] & \text{if } x = [x, x] \neq \emptyset, \\ \emptyset & \text{if } x = \emptyset. \end{cases} \quad (29)$$

When x is a \mathbb{T} -box, **nextOut** acts componentwise.

NOTE 1—For an IEEE 754 format, **nextUp** and **nextDown** are equivalent to the corresponding functions in 5.3.1 of IEEE Std 754-2008.

NOTE 2—In (28), the inner **nextOut**() aims to handle the problem of a function such as $\sin x$ evaluated at a very large argument, where a small relative change in the input can produce a large relative change in the result. The outer

`nextOut()` relaxes the requirement for correct (rather than, say, faithful) rounding, which might be hard to achieve for some special functions at some arguments.

NOTE 3—The input box x might have components of a different type from the result type \mathbb{T} , in the case of 754-conforming mixed-type operations 12.6.2. The `hullT` in (28) forces these to have type \mathbb{T} , so each component of x is widened by at least the local spacing of \mathbb{F} -numbers, at each finite bound.

[Example. Let \mathbb{T} be a 2-digit decimal inf-sup type. Then `nextOut` widens the \mathbb{T} -interval $x = [2.4, 3.7]$ to $[2.3, 3.8]$ —an *ulp* at each end, see 9.7.3. But an operation might accept 4-digit decimal inf-sup inputs, and x might be $[2.401, 3.699]$. Then `nextOut(x)` is `[nextDown(2.401), nextUp(3.699)] = [2.4, 3.7]`, giving an insignificant widening. But

$$\text{nextOut}(\text{hull}_{\mathbb{T}}([2.401, 3.699])) = \text{nextOut}([2.4, 3.7]) = [2.3, 3.8]$$

gives a widening comparable with the precision of \mathbb{T} .]

12.10.2 Accuracy requirements

Following the categories of functions in Table 9.1, the accuracy of the *basic operations*, the *integer functions* and the *absmax functions* shall be *tightest*. The accuracy of the *cancellative addition and subtraction* operations of 10.5.6 is specified in 12.12.5.

For all other operations in Table 9.1, for the reverse mode operations of Table 10.1, and for the recommended operations of Table 10.5 and Table 10.6, the accuracy shall be *valid*, and, for inf-sup types, should be *accurate*. For any operation in these four tables, if any input is Empty the result shall be Empty.

12.10.3 Documentation requirements

An implementation shall document the tightness of each of its interval operations for each supported bare interval type. This shall be done by dividing the set of possible inputs into disjoint subsets (“ranges”) and stating a tightness achieved in each range. This information may be supplemented by further detail, e.g., to give accuracy data in a more appropriate way for a non-inf-sup type.

[Example. Sample tightness information for the `sin` function might be

Operation	Type	Tightness	Range
<code>sin</code>	<code>infsup binary64</code>	<i>tightest</i> <i>accurate</i>	<i>for any $x \subseteq [-10^{15}, 10^{15}]$</i> <i>for all other x.</i>

]

Each operation should be identified by a language- or implementation-defined name of the Level 1 operation (which might differ from that used in this standard), its output type, its input type(s) if necessary, and any other information needed to resolve ambiguity.

12.11 Interval and number literals

12.11.1 Overview

This subclause extends the specifications of 9.7 to define interval literals in the set-based flavor. Interval literals are used as input to `textToInterval` in 12.12.7, and in Clause 13, Input/Output. The following definitions and usages from 9.7 are unchanged: integer literal; value of a literal; case insensitivity; unit in last place; the possibility of implementation-defined or locale-dependent variations and the requirement to document them.

A literal *of the flavor* (9.7.1), to be supported by each implementation of this flavor, is here called **portable**.

12.11.2 Number literals

In addition to those specified in 9.7, the following forms of number literal shall be provided.

- d) Either of the strings `inf` or `infinity` optionally preceded by a plus sign, with value $+\infty$; or preceded by a minus sign, with value $-\infty$.

Table 12.1. Portable interval literal examples.

Form	Literal	Exact decorated value
Special	[]	Empty _{trv}
	[entire]	$[-\infty, +\infty]_{\text{dac}}$
Inf-sup	[1.e-3, 1.1e-3]	$[0.001, 0.0011]_{\text{com}}$
	[-Inf, 2/3]	$[-\infty, 2/3]_{\text{dac}}$
	[0x1.3p-1,]	$[19/32, +\infty]_{\text{dac}}$
	[,]	Entire _{dac}
Uncertain	3.56?1	$[3.55, 3.57]_{\text{com}}$
	3.56?1e2	$[355, 357]_{\text{com}}$
	3.560?2	$[3.558, 3.562]_{\text{com}}$
	3.56?	$[3.555, 3.565]_{\text{com}}$
	3.560?2u	$[3.560, 3.562]_{\text{com}}$
	-10?	$[-10.5, -9.5]_{\text{com}}$
	-10?u	$[-10.0, -9.5]_{\text{com}}$
	-10?12	$[-22.0, 2.0]_{\text{com}}$
	-10??u	$[-10.0, +\infty]_{\text{dac}}$
	-10??	$[-\infty, +\infty]_{\text{dac}}$
NaI	[nai]	Empty _{ill}
Decorated	3.56?1_def	$[3.55, 3.57]_{\text{def}}$

12.11.3 Bare intervals

In addition to those specified in 9.7, the following forms of bare interval literal shall be supported. (A) marks that a new form of literal is Added; (C) marks an existing form, Changed by adding an extra feature.

- Inf-sup form (C): In the string $[l, u]$, the bound l may be $-\infty$ and u may be $+\infty$. Any of l and u may be omitted, with implied values $l = -\infty$ and $u = +\infty$, respectively, e.g., $[,]$ denotes Entire.
- Uncertain form (C): This form, with radius empty or *ulp-count*, is adequate for narrow (hence bounded) intervals, but is severely restricted otherwise. Uncertain form with radius ? is used for unbounded intervals, e.g., $m???d$ denotes $[-\infty, m]$, $m???u$ denotes $[m, +\infty]$ and $m??$ denotes Entire with m being like a comment.
- Special values (A): The strings $[]$ and $[\text{empty}]$, whose bare value is Empty; the string $[\text{entire}]$, whose bare value is Entire. Here and below, space shown between elements of a literal is optional: it denotes zero or more space characters. E.g., one may write $[\text{empty}]$ or $[\text{ empty}]$, etc.

12.11.4 Decorated intervals

The following forms of decorated interval literal shall be supported.

- A bare interval literal sx , an underscore “_”, and a 3-character decoration string sd ; where sd is one of **trv**, **def**, **dac** or **com**, denoting the corresponding decoration dx .

If sx has the bare value x , and if x_{dx} is a permitted combination according to 11.4, then sx_sd has the value x_{dx} . Otherwise sx_sd has no value as a decorated interval literal.

- The string $[\text{nai}]$, with the value Empty_{ill}.

[Examples. Table 12.1 illustrates portable interval literals. These strings are not portable interval literals: empty, [5?1], [1.000.000], [ganz], [entire!comment], [inf], 5???u, [nai]_ill, []_ill, []_def, [0,inf]_com.]

12.11.5 Grammar for portable literals

The syntax of portable integer and number literals and of portable bare and decorated interval literals in this flavor is defined by `integerLiteral`, `numberLiteral`, `bareIntervalLiteral` and `intervalLiteral`, respectively, in a variant of the grammar defined in Table 9.5. The differences are shown in Table 12.2.

Table 12.2. Differences from the general grammar for literals in Table 9.5.

In the left column, ‘A’ marks a new term that is Added, and ‘C’ marks a term that is Changed from its definition in Table 9.5. The change is always in the form of added alternative(s) on the right hand side—these are underlined.

A	infNumLit	{sign}? ("inf" "infinity")
C	numberLiteral	{decNumLit} {hexNumLit} {ratNumLit} <u>{infNumLit}</u>
C	radius	{natural} <u>"?"</u>
A	emptyIntvl	"[" {sp} "]" "[" {sp} "empty" {sp} "]"
A	entireIntvl	"[" {sp} "entire" {sp} "]"
A	specialIntvl	{emptyIntvl} {entireIntvl}
C	bareIntvlLiteral	{pointIntvl} {infSupIntvl} {uncertIntvl} <u>{specialIntvl}</u>
A	Nal	"[" {sp} "nai" {sp} "]"
C	decorationLit	<u>"trv" "def" "dac" "com"</u>
C	intervalLiteral	{bareIntvlLiteral} "-" {decorationLit} <u>{Nal}</u>

The constructor `textToInterval` (12.12.7, 13.2) of any implementation shall accept any portable interval literal. An implementation may restrict support of some input strings (too long strings or strings with a rational number literal). Nevertheless, the constructor shall always return a Level 2 interval (possibly Entire in this case) that contains the Level 1 interval.

An implementation may support interval literals of more general syntax (for example, with underscores in significand). In this case there shall be a value of the conversion specifier *cs* of `intervalToText` in 13.3, that restricts output strings to the portable syntax.

12.12 Required operations on bare and decorated intervals

An implementation shall provide a \mathbb{T} -version, see 12.9, of each operation listed in 12.12.1 to 12.12.10, for each supported type \mathbb{T} . That is, those of the \mathbb{T} -version’s inputs and outputs that are intervals are of type \mathbb{T} (or the corresponding decorated type), except for the conversion operation of 12.12.10 whose output is of type \mathbb{T} and whose input is of any type.

The implementation shall provide the type-independent decoration operations of 12.12.11. It shall provide the reduction operations of 12.12.12 for the parent formats of supported 754-conforming types.

Operations in this subclause are described as functions with zero or more input arguments and one return value. It is language- or implementation-defined whether they are implemented in this way: for instance, two-output division, described in 12.12.3 as a function returning an ordered pair of intervals, might be implemented as a procedure `mulRevToPair(x, y, z1, z2)` with input arguments *x* and *y* and output arguments *z*₁ and *z*₂.

An implementation, or a part thereof, that is 754-conforming shall provide mixed-type operations, as specified in 12.6.2, for the following operations, which correspond to those that IEEE Std 754-2008 requires to be provided as *formatOf* operations.

`add, sub, mul, div, recip, sqrt, sqr, fma.`

An implementation may provide more than one version of some operations for a given type. For instance, it may provide an “accurate” version in addition to a required “tightest” one, to offer a trade-off of accuracy versus speed or code size. How such a facility is provided is language- or implementation-defined.

12.12.1 Interval constants

For each supported bare interval type \mathbb{T} there shall be a \mathbb{T} -version of each constant function `empty()` and `entire()` of 10.5.2, returning a \mathbb{T} -interval with value Empty and Entire, respectively. There shall also be a decorated version of each, returning `newDec(Empty) = Emptytrv` and `newDec(Entire) = Entiredac`, respectively, of the derived decorated type.

12.12.2 Forward-mode elementary functions

Let \mathbb{T} be a supported bare interval type and \mathbb{DT} the derived decorated type. An implementation shall provide a \mathbb{T} -version of each forward arithmetic operation in 10.5.3. Its inputs and output are \mathbb{T} -intervals, and it shall be a Level 2 interval extension of the corresponding point function. Recommended accuracies are given in 12.10.

NOTE—For operations, some of whose arguments are of integer type, such as integer power $\text{pown}(x, p)$, only the real arguments are replaced by intervals.

Each such operation shall have a decorated version with corresponding arguments of type \mathbb{DT} . It shall be a decorated interval extension as defined in 11.6—thus the interval part of its output is the same as if the bare interval operation were applied to the interval parts of its inputs.

The only freedom of choice in the decorated version is how the local decoration, denoted dv_0 in (17) of 11.6, is computed. dv_0 shall be the strongest possible (and is thus uniquely defined) if the accuracy mode of the corresponding bare interval operation is “tightest,” but otherwise is only required to obey (17).

12.12.3 Two-output division

There shall be a \mathbb{T} -version of the two-output division $\text{mulRevToPair}(\mathbf{b}, \mathbf{c})$ of 10.5.5, for each supported bare interval type \mathbb{T} , namely

$$(\mathbf{u}, \mathbf{v}) = \text{mulRevToPair}(\mathbf{b}, \mathbf{c}),$$

where \mathbf{b} and \mathbf{c} are \mathbb{T} -intervals. Each of the outputs \mathbf{u} and \mathbf{v} is a \mathbb{T} -interval that encloses the corresponding Level 1 value.

There shall be a decorated version where each of \mathbf{b} , \mathbf{c} , \mathbf{u} and \mathbf{v} is of the corresponding decorated type. If either input is NaN then both outputs are NaN . Otherwise, if \mathbf{b} and \mathbf{c} are nonempty and $0 \notin \mathbf{b}$, then \mathbf{u} is the same as the result of normal division \mathbf{c}/\mathbf{b} and shall be decorated the same way; while \mathbf{v} is empty and shall be decorated trv . In all other cases each output, empty or not, shall be decorated trv .

12.12.4 Reverse-mode elementary functions

An implementation shall provide a \mathbb{T} -version of each reverse arithmetic operation in 10.5.4, for each supported bare interval type \mathbb{T} . Its inputs and output are \mathbb{T} -intervals.

These operations shall have “trivial” decorated versions, as described in 11.7.

12.12.5 Cancellative addition and subtraction

An implementation shall provide a \mathbb{T} -version of each of the operations cancelMinus and cancelPlus in 10.5.6 for each supported bare interval type \mathbb{T} . Its inputs and output are \mathbb{T} -intervals.

The \mathbb{T} -version shall return an enclosure of the Level 1 value if the latter is defined, and Entire otherwise. It shall return Empty if the Level 1 value is Empty . Thus, for the case of $\text{cancelMinus}(\mathbf{x}, \mathbf{y})$, it returns Entire in these cases:

- either of \mathbf{x} and \mathbf{y} is unbounded;
- $\mathbf{x} \neq \emptyset$ and $\mathbf{y} = \emptyset$;
- \mathbf{x} and \mathbf{y} are nonempty bounded intervals with $\text{width}(\mathbf{x}) < \text{width}(\mathbf{y})$.

$\text{cancelPlus}(\mathbf{x}, \mathbf{y})$ shall be equivalent to $\text{cancelMinus}(\mathbf{x}, -\mathbf{y})$.

If \mathbb{T} is a 754-conforming type, the result shall be the \mathbb{T} -hull of the Level 1 result when this is defined.

These operations shall have “trivial” decorated versions, as described in 11.7.

NOTE—Two cases need care. Examples are given where \mathbb{T} is the inf-sup type of a format \mathbb{F} .

- Determining whether the Level 1 value is defined needs care when \mathbf{x} and \mathbf{y} have equal widths in finite precision. An example is `cancelMinus`($[-1, a], [-b, 1]$) where a and b are positive \mathbb{F} -numbers less than \mathbb{F} 's roundoff unit, where the latter is the smallest positive \mathbb{F} -number u for which the computed value of $1 + u$ is greater than 1.
- For any \mathbf{x}, \mathbf{y} , the Level 1 result is always bounded when it is defined, but might overflow at Level 2. An example is `cancelMinus`($[M, M], [-M, -M]$) where M is the largest finite \mathbb{F} -number. The exact value is $[2M, 2M]$, whose \mathbb{T} -hull is $[M, +\infty]$. An implementation should not return Entire in case of overflow. For any inf-sup type it can avoid doing so, since the Level 1 result has width not exceeding that of \mathbf{x} , and therefore at Level 2 cannot overflow to both $-\infty$ and $+\infty$. Not returning Entire in case of overflow makes the result Entire diagnostic: it occurs if, and only if, there is no value at Level 1.

12.12.6 Set operations

An implementation shall provide a \mathbb{T} -version of each of the operations `intersection` and `convexHull` in 10.5.7 for each supported bare interval type \mathbb{T} . Its inputs and output are \mathbb{T} -intervals.

These operations should return the \mathbb{T} -hull of the exact result. If either input to `intersection` is Empty, or both inputs to `convexHull` are Empty, the result shall be Empty.

NOTE—If \mathbb{T} is an inf-sup type, the operation can always return the exact result.

These operations shall have “trivial” decorated versions, as described in 11.7.

12.12.7 Constructors

For each supported bare or decorated interval type \mathbb{T} , there shall be a \mathbb{T} -version of each constructor in 10.5.8. It returns a \mathbb{T} -datum.

Difficulties in implementation. Both `textToInterval` when its input is a literal of inf-sup form, and `numsToInterval`, involve testing if a boolean value $b = (l \leq u)$ is 0 (false) or 1 (true), to determine whether the interval is empty or nonempty. Here l and u are extended-real numbers. For `textToInterval` they are values of number literals; for `numsToInterval`, they are values of datums of supported number formats.

Evaluating b as 0 when the true value is 1 (a “false negative”) leads to falsely returning Empty as an enclosure of the true nonempty interval—a containment failure. Evaluating b as 1 when the true value is 0 (a “false positive”) is undesirable, but permissible since it returns a nonempty interval as an enclosure for Empty. Implementations shall ensure that false negatives cannot occur and should ensure that false positives cannot occur.

NOTE—Evaluating b correctly can be hard, if finite l and u have values very close in a relative sense and are represented in different ways—e.g., if an implementation allows them to be floating-point values of different radices. It could be especially challenging in an extendable-precision context.

Language rules might cause such errors even when l and u are of the same number format. E.g., in C, if `long double` is supported and has more precision than `double`, default behavior might be to round `long double` inputs l and u to `double`, on entry to a `numsToInterval` call. This would be non-conforming—the comparison $l \leq u$ requires the exact values to be used, which requires use of a version of `numsToInterval` with `long double` arguments.

Bare interval constructors. A bare interval constructor call either **succeeds** or **fails**. This notion is used to determine the value returned by the corresponding decorated interval constructor.

For the constructor `numsToInterval`(l, u), the inputs l and u are datums of supported number formats, where the format of l may differ from that of u . For a given bare interval type \mathbb{T} :

- There shall be a version of `numsToInterval` where l and u are both of the same format, compatible with \mathbb{T} , see 12.4.
- If \mathbb{T} is 754-conforming, there shall be a *formatOf* version of `numsToInterval` that accepts l and u having any combination of supported IEEE 754 formats of the same radix as \mathbb{T} .

Apart from these two requirements, what (l, u) format combinations are provided is language- or implementation-defined.

The constructor call succeeds if (a) the implementation determines that the call has a Level 1 value x , see 10.5.8, or (b) it cannot determine whether a Level 1 value exists, see the discussion at the start of this subclause. The conditions under which case (b) might occur shall be documented.

Case (a) is that neither l nor u is NaN, and the exact extended-real values they denote—also called l and u for brevity—are known to satisfy $l \leq u$, $l < +\infty$, $u > -\infty$. In this case the result shall be a \mathbb{T} -interval containing x . In case (b) the result shall be a \mathbb{T} -interval containing l and u .

If \mathbb{T} is a 754-conforming type and l and u are of IEEE 754 formats with the same radix as \mathbb{T} , case (b) cannot arise, and the result shall be the \mathbb{T} -hull of x ; in particular if x is an exact \mathbb{T} -interval, the result is x . For other cases, the tightness of the result is implementation-defined.

Otherwise the call fails, and the result is Empty.

For the constructor `textToInterval(s)`, the input s is a string. The constructor call succeeds if: (a) the implementation determines that s is a bare interval literal with value x , see 12.11; or (b) s is of inf-sup form $[l, u]$ with finite l and u , but the implementation cannot determine whether a Level 1 value exists, i.e. whether $l \leq u$. The conditions under which case (b) might occur shall be documented.

In case (a) the result shall be a \mathbb{T} -interval containing x . If \mathbb{T} is a 754-conforming type, this shall be the \mathbb{T} -hull of x ; in particular if x is an exact \mathbb{T} -interval, the result is x . For other types \mathbb{T} , the tightness of the result is implementation-defined. In case (b) the result shall be an interval containing l and u .

Otherwise the call fails, and the result is Empty.

Decorated interval constructors. Let the prefix **b-** or **d-** denote the bare or decorated version of a constructor. If `b-numsToInterval(l, u)` or `b-textToInterval(s)` succeeds with result y , then `d-numsToInterval(l, u)` or `d-textToInterval(s)`, respectively, succeeds with result `newDec(y)`, see 11.5.

If s is a decorated interval literal sx_sd with Level 1 value x_{dx} , see 12.11.4, and `b-textToInterval(sx)` succeeds with result y , then `d-textToInterval(s)` succeeds with result y_{dy} , where $dy = dx$ except when $dx = \text{com}$ and overflow occurred, that is, x is bounded and y is unbounded. Then dy shall equal `dac`.

Otherwise the call fails, and the result is NaN.

Exception behavior. `UndefinedOperation` is signaled when a constructor call fails. NOTE—When signaled by the decorated constructor it will normally be ignored since returning NaN gives sufficient information.

`PossiblyUndefinedOperation` is signaled when the implementation cannot determine whether a Level 1 value exists (the two cases (b) above).

12.12.8 Numeric functions of intervals

An implementation shall provide a \mathbb{T} -version of each numeric function in Table 10.2 of 10.5.9 for each supported bare interval type \mathbb{T} , giving a result in a supported number format \mathbb{F} that should be compatible with \mathbb{T} , see 12.4. An implementation may provide several versions, returning results in different formats. If \mathbb{T} is a 754-conforming type, versions shall be provided giving a result in any supported IEEE 754 format of the same radix as \mathbb{T} .

The mapping of a Level 1 value to an \mathbb{F} -number is defined in terms of the following rounding methods, which are functions from $\overline{\mathbb{R}}$ to \mathbb{F} . NOTE—These functions help define operations of the standard but are not themselves operations of the standard.

Round toward positive: x maps to the smallest \mathbb{F} -number $\geq x$. If \mathbb{F} has signed zeros, 0 maps to $+0$.

Round toward negative: x maps to the largest \mathbb{F} -number $\leq x$. If \mathbb{F} has signed zeros, 0 maps to -0 .

Round to nearest: x maps to the \mathbb{F} -number (possibly $\pm\infty$) closest to x , with an implementation-defined rule for the distance to an infinity, and for the method of tie-breaking when more than one member of \mathbb{F} has the “closest” property. If \mathbb{F} has signed zeros, 0 maps to +0.

The implementation shall document how it handles cases not covered by the above rules, e.g., the distance to an infinity and the method of tie-breaking. If \mathbb{F} is a 754-conforming format, the tie-breaking method shall follow 4.3.1 and 4.3.3 of IEEE Std 754-2008; otherwise it is language- or implementation-defined.

In formats that have a signed zero, a Level 1 value of 0 shall be returned as -0 by **inf**, and $+0$ by all other functions in this subclause.

inf(x) returns the Level 1 value rounded toward negative.

sup(x) returns the Level 1 value rounded toward positive.

mid(x): the result m is defined by the following cases, where \underline{x}, \bar{x} are the exact (Level 1) lower and upper bounds of x .

$x = \text{Empty}$	NaN.
$x = \text{Entire}$	0.
$\underline{x} = -\infty, \bar{x}$ finite	The finite negative \mathbb{F} -number of largest magnitude.
\underline{x} finite, $\bar{x} = +\infty$	The finite positive \mathbb{F} -number of largest magnitude.
\underline{x}, \bar{x} both finite	m should be, and if \mathbb{T} is a 754-conforming type shall be, the Level 1 value rounded to nearest.

The implementation shall document how it handles the last case.

NOTE—If \mathbb{F} is compatible with \mathbb{T} , a nonempty x always contains m . If \mathbb{F} is not compatible with \mathbb{T} this might be impossible. E.g., let \mathbb{T} be **inf-sup binary64**, and let $x = [1 + u, 1 + 2u]$ where u is the **binary64** roundoff unit. If \mathbb{F} is **binary64** (compatible) then m is $1 + u$ or $1 + 2u$ depending on the tie-breaking method. If \mathbb{F} is **binary32** (incompatible) then no \mathbb{F} -numbers lie in x , and $m = 1$.

rad(x) returns NaN if x is empty, and otherwise the smallest \mathbb{F} -number r such that x is contained in the exact interval $[m - r, m + r]$, where m is the value returned by **mid**(x).

NOTE—**rad**(x) might be $+\infty$ even though x is bounded, if \mathbb{F} has insufficient range. However, if \mathbb{F} is an IEEE 754 format and \mathbb{T} is the derived inf-sup type, **rad**(x) is finite for all bounded nonempty intervals.

wid(x) returns NaN if x is empty. Otherwise it returns the Level 1 value rounded toward $+\infty$.

NOTE—For nonempty bounded x the ratio **wid**(x)/**rad**(x), which is always 2 at Level 1, ranges from 1 to $+\infty$ at Level 2. When \mathbb{F} is an IEEE 754 format and \mathbb{T} is the derived inf-sup type, it takes the value 1 if the bounds of x are adjacent subnormal numbers, and the value $+\infty$ if $x = [-\text{realmax}, \text{realmax}]$.

mag(x) returns NaN if x is empty. Otherwise it returns the Level 1 value rounded toward positive.

mig(x) returns NaN if x is empty. Otherwise it returns the Level 1 value rounded toward negative, except that 0 maps to +0 if \mathbb{F} has signed zeros.

Each bare interval operation in this subclause shall have a decorated version, where each input of bare interval type is replaced by an input having the corresponding decorated interval type, and the result format is that of the bare operation. Following 11.7, if any input is NaI, the result is NaN. Otherwise the result is obtained by discarding the decoration and applying the corresponding bare interval operation.

12.12.9 Boolean functions of intervals

An implementation shall provide a \mathbb{T} -version of the function **isEmpty**(x) and the function **isEntire**(x) in 10.5.10, for each supported bare interval type \mathbb{T} .

There shall be a function **isNaI**(x) for input x of any decorated type, that returns **true** if x is NaI, else **false**.

There shall be a \mathbb{T} -version of each of the comparison relations in Table 10.3 of 10.5.10, for each supported bare interval type \mathbb{T} . Its inputs are \mathbb{T} -intervals.

For a 754-conforming part of an implementation, mixed-type versions of these relations shall be provided, where the inputs have arbitrary types of the same radix.

These comparisons shall return, in all cases, the correct value of the comparison applied to the intervals denoted by the inputs as if in infinite precision. In particular `equal(\mathbf{x} , \mathbf{y})`, for those bare interval inputs \mathbf{x} and \mathbf{y} for which it is defined, shall return `true` if and only if \mathbf{x} and \mathbf{y} (ignoring their type) are the same mathematical interval, see 12.3.

Each bare interval operation in this subclause shall have a decorated version, where each input of bare interval type is replaced by an input having the corresponding decorated interval type. Following 11.7, if any input is `NaN`, the result is `false` (in particular `equal(NaN, NaN)` is `false`). Otherwise the result is obtained by discarding the decoration and applying the corresponding bare interval operation.

12.12.10 Interval type conversion

An implementation shall provide, for each supported bare interval type \mathbb{T} , the operation `\mathbb{T} -convertType` to convert an interval of any supported bare interval type to type \mathbb{T} , and the operation `\mathbb{DT} -convertType` to convert an interval of any supported decorated interval type to the corresponding decorated type \mathbb{DT} . Conversion is done by applying the \mathbb{T} -hull operation, see 12.8.1. For a bare interval \mathbf{x} :

$$\mathbb{T}\text{-convertType}(\mathbf{x}) = \text{hull}_{\mathbb{T}}(\mathbf{x}).$$

Thus if \mathbb{T} is an explicit type, see 12.8.1, the result is the unique tightest \mathbb{T} -interval containing \mathbf{x} .

Conversion of a decorated interval is done by converting the interval part, except that if the decoration is `com` and the conversion overflows (produces an unbounded interval) the decoration becomes `dac`. That is, `\mathbb{DT} -convertType(\mathbf{x}_{dx}) = \mathbf{y}_{dy}` where

$$\begin{aligned} \mathbf{y} &= \mathbb{T}\text{-convertType}(\mathbf{x}); \\ dy &= \begin{cases} \text{dac} & \text{if } dx = \text{com} \text{ and } \mathbf{y} \text{ is unbounded,} \\ dx & \text{otherwise.} \end{cases} \end{aligned}$$

12.12.11 Operations on/with decorations

An implementation shall provide the operations of 11.5. These comprise the comparison operations `=`, `≠`, `>`, `<`, `≥`, `≤` for decorations; and, for each supported bare interval type and corresponding decorated type, the operations `newDec`, `intervalPart`, `decorationPart` and `setDec`.

A call `intervalPart(NaN)`, whose value is undefined at Level 1, shall return `Empty` at Level 2, and shall signal the `IntvlPartOfNaN` exception to indicate that a valid interval has been created from the ill-formed interval.

12.12.12 Reduction operations

For each supported 754-conforming interval type, an implementation shall provide, for the parent format of that type, the four reduction operations `sum`, `dot`, `sumSquare` and `sumAbs` in 9.4 of IEEE Std 754-2008, correctly rounded.

Correctly rounded means that the returned result is defined as follows.

- If the exact result is defined as an extended-real number, return this after rounding to the relevant format according to the current rounding direction. An exact zero shall be returned as `+0` in all rounding directions, except for `roundTowardNegative`, where `−0` shall be returned.
- For `dot` and `sum`, if a `NaN` is encountered, or if infinities of both signs were encountered in the sum, `NaN` shall be returned. (“NaN encountered” includes the case $\infty \times 0$ for `dot`.)
- For `sumAbs` and `sumSquare`, if an `Infinity` is encountered, `+∞` shall be returned. Otherwise, if a `NaN` is encountered, `NaN` shall be returned.

(These rules allow for short-circuit evaluation in certain cases.)

All other behavior, such as overflow, underflow, setting of IEEE 754 flags, signaling exceptions, and behavior on vectors whose length is given as non-integral, zero or negative, shall be as specified in 9.4 of IEEE Std 754-2008. In particular, evaluation is as if in exact arithmetic up to the final rounding, with no possibility of intermediate overflow or underflow.

Also, since correct rounding applies, the Inexact flag shall be set unless an exact extended-numeric result is returned. (If a final overflow or underflow is indicated, the result is inexact.)

It is recommended that these operations for the parent format \mathbb{F} be based on an accumulator format datatype $A(\mathbb{F})$ 12.13.5, able to represent exact dot products of long \mathbb{F} vectors.

12.13 Recommended operations

12.13.1 Forward-mode elementary functions

The functions listed in 10.6.1 are arithmetic operations. If any of them is provided, it shall have a version for each bare and decorated interval type, specified as is done in 12.12.2 for the required operations.

12.13.2 Slope functions

The functions listed in 10.6.2 shall be handled in the same way as those in 12.13.1.

12.13.3 Boolean functions of intervals

The functions listed in 10.6.3 shall be handled in the same way as those in 12.13.1.

If \mathbb{T} -version of the function `isMember(m, x)` is provided, the number format \mathbb{F} of the argument m should be compatible with \mathbb{T} . An implementation may provide several \mathbb{T} -versions, with different formats of m . If \mathbb{T} is a 754-conforming type, versions should be provided with the argument m of any supported IEEE 754 format of the same radix as \mathbb{T} .

12.13.4 Extended interval comparisons

How the operations in 10.6.4 are handled at Level 2 is implementation-defined.

12.13.5 Exact reduction operations

An implementation that provides a 754-conforming type for the parent format \mathbb{F} should provide an *accumulator format* datatype $A(\mathbb{F})$ associated with the \mathbb{F} , and associated operations. An $A(\mathbb{F})$ datum z is capable of exactly representing dot products of vectors of any reasonable length of arbitrary finite \mathbb{F} -numbers.

The following operations should be provided.

- `convert` converts from an accumulator format to a floating-point format, or vice versa, or from one accumulator format to another.
- `exactAdd` and `exactSub` adds or subtracts two accumulator or floating-point format operands, of which at least one is an accumulator, giving an accumulator format result.
- `exactFma` computes $z + x * y$ where z has an accumulator format and x, y are of floating-point format, giving an accumulator format result.
- `exactDotProduct`. Let a and b be vectors of length n holding floating-point numbers of format \mathbb{F} . Then `exactDotProduct(a, b)` computes $a \cdot b = \sum_{k=1}^n a_k b_k$ exactly, giving an accumulator format result.

The result of all operations may be converted if necessary to a specified result format by application of the `convert` operation.

[Example. The Complete Arithmetic, specified by Kulisch and Snyder [B7],[B6], is an example of implementation of accumulator format and exact reduction operations. The recommended accumulator format for the binary64 format in

the Complete Arithmetic has 4 bits for sign and status, 2134 bits before the point, and 2150 after the point, for a total of 4288 bits or 536 bytes; this allows for at least 2^{88} multiply-adds before overflow can occur.]

13. Input and output (I/O) of intervals

13.1 Overview

This clause of the standard specifies conversion from a text string that holds an interval literal to an interval internal to a program (input), and the reverse (output). The methods by which strings are read from, or written to, a character stream are language- or implementation-defined, as are variations in some locales (such as specific character case matching).

Containment is preserved on input and output so that, when a program computes an enclosure of some quantity given an enclosure of the data, it can ensure this holds all the way from text data to text results.

In addition to the above I/O, which might incur rounding errors on output and/or input, each interval type \mathbb{T} has an *exact text representation*, via operations that convert any internal \mathbb{T} -interval x to a string s , and back again to recover x exactly.

13.2 Input

Input is provided for each supported bare or decorated interval type \mathbb{T} by the \mathbb{T} -version of `textToInterval(s)`, where s is a string, as specified in 12.12.7. It accepts an arbitrary interval literal s and returns a \mathbb{T} -interval enclosing the Level 1 value of s .

For 754-conforming types \mathbb{T} the required tightness is specified in 12.12.7. For other types, the tightness is implementation-defined.

NOTE—This provides the basis for free-format input of interval literals from a text stream, as might be provided by overloading the `>>` operator in C++.

13.3 Output

An implementation shall provide an operation

$$\text{intervalToText}(\mathbf{X}, cs),$$

where cs is optional. Here \mathbf{X} is a bare or decorated interval datum of any supported interval type \mathbb{T} , and cs is a string, the conversion specifier. The operation converts \mathbf{X} to an interval literal string s , see 12.11, which shall be related to \mathbf{X} as follows, where \mathbf{Y} is the Level 1 value of s .

- a) Let \mathbb{T} be a bare type. Then \mathbf{Y} shall contain \mathbf{X} and shall be empty if \mathbf{X} is empty.
- b) Let \mathbb{T} be a decorated type. If \mathbf{X} is NaI, then \mathbf{Y} shall be NaI. Otherwise, write $\mathbf{X} = x_{dx}$, $\mathbf{Y} = y_{dy}$. Then
 - 1) y shall contain x and shall be empty if x is empty.
 - 2) dy shall equal dx , except in the case that $dx = \text{com}$ and overflow occurred, that is, x is bounded and y is unbounded. Then dy shall equal dac .

NOTE— \mathbf{Y} being a Level 1 value is significant. E.g., for a bare type \mathbb{T} , it is not allowed to convert $\mathbf{X} = \emptyset$ to the string `garbage`, even though converting `garbage` back to a bare interval at Level 2 by \mathbb{T} -`textToInterval` gives \emptyset , because `garbage` has no Level 1 value as a bare interval literal.

The tightness of enclosure of \mathbf{X} by \mathbf{Y} is language- or implementation-defined.

If present, cs lets the user control the layout of the string s in a language- or implementation-defined way. The implementation shall document the recognized values of cs and their effect; other values are called *invalid*.

If cs is invalid, or makes an unsatisfiable request for a given input \mathbf{X} , the output shall still be an interval literal whose value encloses \mathbf{X} . A language- or implementation-defined extension to interval literal syntax may be used to make it obvious that this has occurred. [Example. Suppose, for uncertain form, that m is undefined or r is “unreasonably large.” Then a string such as `[Entire!uncertain form conversion error]` might be produced.

The implementation of `textToInterval` would need to accept this string as meaning the same as `[Entire]`. Such a string is not a portable literal, see 12.11.5.]

Among the user-controllable features should be the following, where l , u are the interval bounds for inf-sup form, and m , r are the base point and radius for uncertain form, as defined in 12.11.

- a) It should be possible to specify the preferred overall field width (the length of s), and whether output is in inf-sup or uncertain form.
- b) It should be possible to specify how Empty, Entire and NaI are output, e.g., whether lower or upper case, and whether Entire becomes `[Entire]` or `[-Inf, Inf]`.
- c) For l , u and m , it should be possible to specify the field width, and the number of digits after the point or the number of significant digits. For r , which is a non-negative integer ulp-count, it should be possible to specify the field width. There should be a choice of radix, at least between decimal and hexadecimal.
- d) For uncertain form, it should be possible to select the default symmetric form, or the one sided (u or d) forms. It should be possible to choose whether an exponent field is absent (and m is output to a given number of digits after the point) or present (and m is output to a given number of significant digits). Despite the normalization rules in 13.4.1, trailing zeros may be added to m as needed. E.g., if $X \approx [2.1995, 2.2007]$, s might be `2.200?7`, `2.20?1` or `2.2?1` depending on the user-requested tightness.

It is implementation-defined how large r can be in the $m ? r$ form before switching to one of the $m??$ forms denoting an unbounded interval. In $m ? r$ form, m and r should be chosen to give the tightest enclosure of X subject to m 's specified number of digits after the point, or significant digits. For example, to convert `[0.9999, 1.0001]` to this form with 2 significant digits, `9.9?2e-1`, with exact value `[0.97, 1.01]`, might be considered preferable to `1.0?1e0`, with exact value `[0.9, 1.1]`.

- e) It should be possible to output the bounds of an interval without punctuation, e.g., `1.234 2.345` instead of `[1.234, 2.345]`. For instance, this might be a convenient way to write intervals to a file for use by another application.

If cs is absent, output should be in a general-purpose layout (analogous, e.g., to the `%g` specifier of `fprintf` in C). There should be a value of cs that selects this layout explicitly.

NOTE—This provides the basis for free-format output of intervals to a text stream, as might be provided by overloading the `<<` operator in C++.

If an implementation supports more general syntax of interval literals than the portable syntax defined in 12.11.5, there shall be a value of cs that restricts output strings to the portable syntax.

If T is a 754-conforming bare type, there shall be a value of cs that produces behavior identical with that of `intervalToExact`, below. That is, the output is an interval literal that, when read back by `T-textToInterval`, recovers the original datum exactly.

13.4 Exact text representation

For any supported bare interval type T , an implementation shall provide operations `intervalToExact` and `exactToInterval`. Their purpose is to provide a portable exact representation of every bare interval datum as a string.

These operations shall obey the **recovery requirement**:

For any T -datum x , the value $s = T\text{-intervalToExact}(x)$ is a string,
such that $y = T\text{-exactToInterval}(s)$ is defined and equals x .

NOTE—From 12.3, this is datum identity: x and y have the same Level 1 value and the same type. They might differ at Level 3, e.g., a zero bound might be stored as `-0` in one and `+0` in the other.

If T is a 754-conforming type, the string s shall be an interval literal which, for nonempty x , is of inf-sup type, with the lower and upper bounds of x converted as described in 13.4.1. For such s , the operation `exactToInterval` is functionally equivalent to `textToInterval`.

If \mathbb{T} is not 754-conforming, there are no restrictions on the form of the string s apart from the above recovery requirement. However, the representation should display, in a human-readable way, the exact values of the parameters of the type's mathematical description (e.g., m , \underline{r} and \bar{r} for a type that represents an interval in triplex form $[m + \underline{r}, m + \bar{r}]$).

The algorithm by which `intervalToExact` converts x to s is regarded as part of the definition of the type and shall be documented by the implementation.

[Example. Writing a binary64 floating-point datum exactly in hexadecimal-significand form passes the "readability" test since it displays the parameters sign, exponent and significand. Dumping its 64 bits as 16 hex characters does not.]

Since `exactToInterval` creates an interval from non-interval data, it is a constructor similar to `textToInterval`. When its input is invalid, it shall return `Empty` and signal one of the exceptions `UndefinedOperation` or `PossiblyUndefinedOperation` as specified in 12.12.7.

13.4.1 Conversion of IEEE 754 numbers to strings

An IEEE 754 format \mathbb{F} is defined by the parameters $b, p, emax$, and $emin$, where

b is the radix, 2 or 10;

p is the number of digits in the significand (precision);

$emax$ is the maximum exponent; and

$emin = 1 - emax$ is the minimum exponent.

(See also 3.3 of IEEE Std 754-2008)

A finite \mathbb{F} -number x can be represented $(-1)^s \times b^e \times m$ where $s = 0$ or 1 , e is an integer, $emin \leq e \leq emax$, and m has a p -digit radix b expansion $d_0.d_1d_2 \dots d_{p-1}$, where d_i is an integer digit $0 \leq d_i < b$ (so $0 \leq m < b$). As used within interval literals, x denotes a real number, with no distinction between -0 and $+0$. To make the representation unique, constraints are imposed in three mutually exclusive cases:

- A *normal* number, with $|x| \geq b^{emin}$, shall have $d_0 \geq 1$ (so $1 \leq m < b$).
- A *subnormal* number, with $0 < |x| < b^{emin}$, shall have $e = emin$, which implies $d_0 = 0$ (so $0 < m < 1$).
- *Zero*, $x = 0$, shall have sign bit $s = 0$ and exponent $e = 0$ (and necessarily $m = 0$).

NOTE—For $b = 2$ the standard form used by IEEE 754 is the same as this, except for replacing zero by two signed zeros, with exponent $e = emin$. For $b = 10$, IEEE 754 normal numbers have several representations, if they need fewer than p digits in their expansion. The standard form above chooses the representation with smallest quantum, which is the unique one having $d_0 \neq 0$.

The rules given below for converting x to a string $xstr$ allow user-, language- or implementation-defined choice while ensuring the values of s , m and e are easily found from $xstr$ in each of these cases, even without knowledge of the format parameters $p, emax, emin$.

$xstr$ is the concatenation of: a sign part $sstr$; a significand part $mstr$; and an exponent part $estr$. If $b = 2$, the hex-indicator "0x" is prefixed to $mstr$.

$sstr$ is "-" or an optional "+", as appropriate.

If $b = 10$, $mstr$ is the (decimal) expansion $d_0.d_1d_2 \dots d_{p-1}$, optionally abbreviated by removing some or all trailing zeros. If this leaves no digits after the point, the point may be removed. If $b = 2$, $mstr$ is formed from the (binary) expansion $d_0.d_1d_2 \dots d_{p-1}$, abbreviated in the same way, and then converted to a hexadecimal string $D_0.D_1 \dots$ (so necessarily D_0 is 1 if x is normal, 0 if x is subnormal or zero).

$estr$ consists of "e" if $b = 10$, "p" if $b = 2$, followed by the exponent e written as a signed decimal integer, with the sign optional if $e \geq 0$.

[Examples. In any binary format, the number 2 (with $s = 0$, $m = 1$, $e = 1$) may be written as 0x1p1 or +0x1.0p+01, etc., but not as 0x2p0; while $\frac{1}{2}$ may be written as 0x1p-1 or +0x1.0p-01, etc. The number -4095 (with $s = 1$, $m = \frac{4095}{2048}$, $e = 11$) may be written as -0x1.ffep+11.

In decimal32 (see IEEE Std 754-2008 Table 3.6), which has $p = 7$, the smallest positive normal number may be written $1\text{e-}95$ or $+1.000000\text{e-}95$, etc.; and the next number below it as $0.999999\text{e-}95$. The smallest positive number may be written $0.000001\text{e-}95$.]

Above, alphabetic characters have been written in lowercase, but they may be in either case.

A shorter form for subnormal numbers may be used, normalized by requiring $d_1 \neq 0$; however, to find the canonical m and e from $xstr$ one then needs to know $emin$. For instance, the smallest positive decimal32 number $x = 0.000001\text{e-}95$ has the shorter form $0.1\text{e-}101$, but to deduce that x has $m = 0.000001$ and $e = -95$ one needs to know that $emin = -95$ for this format.

13.4.2 Exact representations of comparable types

The exact text representation of a bare interval of any type should also be an exact representation in any wider (in the sense of 12.5.1) type, which when converted back produces the mathematically same interval.

That is, let type \mathbb{T}' be wider than type \mathbb{T} . Let x be a \mathbb{T} -interval and let

$$s = \mathbb{T}\text{-intervalToExact}(x).$$

Then

$$x' = \mathbb{T}'\text{-exactToInterval}(s)$$

should be defined and equal to $\mathbb{T}'\text{-convertType}(x)$.

NOTE—If \mathbb{T} and \mathbb{T}' are 754-conforming types, this property holds automatically, because of the properties of `textToInterval` and the fact that s is an interval literal.

14. Levels 3 and 4 description

14.1 Overview

This clause is mostly about Level 3, where Level 2 datums are represented, and operations on them described. 14.4, on interchange representation and encoding, also contains some Level 4 requirements.

Level 3 entities are here called *objects*. They represent Level 2 datums and may be referred to as *concrete*, while the datums are *abstract*.

14.2 Representation

Individual datums of an abstract type or format are *represented* by individual objects of its concrete type or format. The property that defines a representation, for a given type or format, is:

(30)

Each datum shall be represented by at least one object. Each object shall represent at most one datum.

[Examples. Let \mathbb{F} be an IEEE 754 format and let \mathbb{T} the derived (bare) inf-sup type. Three possible representations are:

- **inf-sup form.** Any \mathbb{T} -interval x is represented at Level 3 by the object $(\inf(x), \sup(x))$ of two Level 2 numbers – members of \mathbb{F} . All intervals have only one Level 3 representative because operations `inf` and `sup` are uniquely defined at Level 2 12.12.8: interval $[0, 0]$ has representative $(-0, +0)$, interval `Empty` has representative $(+\infty, -\infty)$.
- **inf-sup-nan form.** The objects are defined to be pairs (l, u) where l, u are members of \mathbb{F} . A nonempty \mathbb{T} -interval $x = [\underline{x}, \bar{x}]$ is represented by an object (l, u) such that the values of l and u are \underline{x} and \bar{x} , and `Empty` is represented by (NaN, NaN) . Its valid objects are (NaN, NaN) , together with all (l, u) such that l, u are not NaN and $l \leq u$, $l < +\infty$, $u > -\infty$.
- **neginf-sup-nan form.** This is as the previous, except that for a nonempty \mathbb{T} -interval x the value of l is $-\underline{x}$.

If, in these descriptions l, u and NaN are viewed as Level 2 datums, then interval $[0, 0]$ has four representatives in **inf-sup-nan** and **neginf-sup-nan** forms: $(-0, +0)$, $(-0, -0)$, $(+0, +0)$, $(+0, -0)$. Each nonempty interval with nonzero bounds has only one representative: there are unique l and u . `Empty` has also only one representative: there is a unique NaN. However, NaN itself has representatives, and from this viewpoint `Empty` has more than one representative: there are many NaNs, quiet or signaling and with different payloads, to use in `Empty` = (NaN, NaN) .

/

14.3 Operations and representation

Each Level 2 (abstract) library operation is implemented by a corresponding Level 3 (concrete) operation, whose behavior shall be consistent with the abstract operation.

When an input Level 3 object does not represent a Level 2 datum, the result is implementation-defined. An implementation shall provide means for the user to specify that an `InvalidOperand` exception be signaled when this occurs.

To promote reproducibility (1.8, 6.5), an implementation should provide a computational mode where, provided certain programming restrictions are adhered to, the computed *values* depend only on the Level 2 values of inputs.

NOTE—Such an effect might be achieved by canonicalizing, i.e., ensuring that, at least for operations with numeric output, the *representative* of the output is independent of the representatives of the inputs. However, doing so incurs a cost, and there will be an implementation-defined trade-off between the run time overhead of a “reproducible mode,” and its scope, i.e., how few programming restrictions it imposes.

As an example, let \mathbb{F} be an IEEE 754 decimal format and \mathbb{T} the derived inf-sup type. Suppose a \mathbb{T} -interval $[l, u]$ is represented at Level 3 as the pair of \mathbb{F} -numbers (l, u) . Let f be the expression

$$f(x) = \text{sameQuantum}(\text{inf}(x), 0.3)$$

(see 5.7.3 of IEEE Std 754-2008). Suppose a program reads x using `textToInterval`, first from the string `[0.3,inf]`, then from `[0.30,inf]`. A straightforward implementation might store them with the two Level 3 representations $x' = (0.3, +\infty)$ and $x'' = (0.30, +\infty)$, where 0.3 and 0.30 stand for the IEEE 754 decimal number objects $(0, -1, 3)$ and $(0, -2, 30)$ denoting the same real value $0.3 = 3 \times 10^{-1} = 30 \times 10^{-2}$. Thus $x' = x''$ at both Level 1 and Level 2, and the user might expect they give the same output, but

$$\begin{aligned} f(x') &= \text{sameQuantum}(0.3, 0.3) && = \text{true;} \\ f(x'') &= \text{sameQuantum}(0.30, 0.3) && = \text{false.} \end{aligned}$$

The standard does not say which of these results is “correct.” Rather than change the implementation, it might be better to document that such code must be avoided if reproducible behavior is required.

14.4 Interchange representations and encodings

The purpose of interchange representations and encodings is to allow the loss-free exchange of Level 2 interval data between 754-conforming implementations. This is done by imposing a standard Level 3 representation using Level 2 number datums and delegating interchange encoding of number datums to the IEEE 754 standard.

Let x be a datum of the bare interval inf-sup type \mathbb{T} derived from a supported IEEE 754 format \mathbb{F} . Its standard Level 3 representative is an ordered pair $(\text{inf}(x), \text{sup}(x))$ of two Level 2 \mathbb{F} -numbers as defined in 12.12.8. For example, the only representative of Empty is the pair $(+\infty, -\infty)$ and the only representative of $[0, 0]$ is the pair $(-0, +0)$.

Let x_{dx} be a datum of the decorated interval type \mathbb{DT} derived from \mathbb{T} . Its standard Level 3 representative is an ordered triple $(\text{inf}(x_{dx}), \text{sup}(x_{dx}), dx)$ of two Level 2 \mathbb{F} -datums and a decoration. For example, the only representative of $\text{Empty}_{\text{trv}}$ is the triple $(+\infty, -\infty, \text{trv})$ and the only representative of NaN is the triple $(\text{NaN}, \text{NaN}, \text{ill})$.

Interchange Level 4 encoding of an interval datum is a bit string that comprises, in the order defined above, the IEEE 754 interchange encodings of the two floating-datums, and, for decorated intervals, of the decoration represented as an **octet** (bit string of length 8, equivalently 8-bit byte) as follows:

```
ill  00000000
trv  00000100
def  00001000
dac  00001100
com  00010000
```

This encoding permits future refinement without disturbing the propagation order of the decorations.

NOTE—The above rules imply that an interval has a unique interchange representation if it is not NaN and in a binary format, but not generally otherwise. The reason for the rules is that the sign of a zero bound cannot convey any information relevant to intervals; but an implementation may potentially use cohort information in decimal formats (IEEE Std 754-2008 3.5.1), or a NaN payload.

When (optional) compressed intervals are supported, their interchange representation is as described above for bare intervals, if the compressed datum is a non-empty interval; if it is a decoration it is represented as a pair of floating-point datums, the first of which is NaN, and the second is a floating-point integer that encodes the decoration as follows:

ill	0.0
emp	2.0
trv	4.0
def	8.0
dac	12.0

NOTE—These encodings match the octet-encoded decorations of decorated intervals, when interpreted as small integers, with a new decoration **emp** to represent a compressed Empty interval. The **com** decoration cannot occur in a compressed interval.

Export and import of interchange formats normally occur as a sequence of octets, e.g., in a file or a network packet. There is therefore a need to define the **octet-encoding** that maps the conceptual Level 4 bit string encodings of floating-point datums (as specified by IEEE Std 754-2008) and of decorations (not specified in IEEE Std 754-2008) into an octet sequence. Octet-encoding might optionally insert padding zero octets for alignment.

Applications exchanging data need to describe the types and layout thereof; standards like this one (and IEEE Std 754-2008) only define the representation of individual datums, and then only as conceptual Level 4 entities. Environments whose primary focus is universal portability (e.g., Java) may fully define the representation at the level of an octet sequence, even when this does not match the natural in-memory layout of the platform. The standard merely requires the parameters of the octet-encoding be communicated, which may avoid conversion costs when bulk data are exchanged among similar systems. Those parameters are¹⁴:

- width in bits of the floating-point interchange format (see 3.6 of IEEE Std 754-2008);
- for decimal formats, whether BID or DPD encoding is used;
- byte order (Endianness) of the floating-point datums: Big-Endian, Little-Endian, or mixed (in the last case it might be useful to export a template);
- number of padding zero octets inserted before or after the decoration for alignment.

The standard does not define how this information is to be conveyed. This is not a limitation because applications must already agree on what it is that is being exchanged.

[Example. The interchange representation of *inf-sup* binary32 decorated interval $[-1, 3]_{\text{com}}$ is a triple

$(-0x1.000000p0, +0x1.800000p1, \text{com})$.

The Level 4 interchange encodings of its fields are these bit strings (with underscores for readability):

$-0x1.000000p0$	10111111_10000000_00000000_00000000 ,
$+0x1.800000p1$	01000000_01000000_00000000_00000000 ,
com	00010000 .

The Level 4 interchange encoding of the interval is a bit string of length 72:

10111111_10000000_00000000_00000000_01000000_01000000_00000000_00000000_00010000 .

¹⁴Those parameters describe the exported data, not necessarily what the platform's native representation is. The possibility of universally portable representations, e.g., as used by Java, is included.

The Big-Endian interchange octet-encodings of its fields are these sequences of octets:

$-0 \times 1.000000p0$	10111111	10000000	00000000	00000000
$+0 \times 1.800000p1$	01000000	01000000	00000000	00000000
com	00010000			

and the interchange octet-encoding of the interval is (assuming no padding octets are inserted) a sequence of 9 octets:

10111111 10000000 00000000 00000000 01000000 01000000 00000000 00000000 00010000 .

The Little-Endian interchange octet-encodings of its fields are sequences of octets:

$-0 \times 1.000000p0$	00000000	00000000	10000000	10111111
$+0 \times 1.800000p1$	00000000	00000000	01000000	01000000
com	00010000			

and the interchange octet-encoding of the interval is a sequence of 9 octets:

00000000 00000000 10000000 10111111 00000000 00000000 01000000 01000000 00010000 .

]

A 754-conforming implementation shall provide an interchange encoding for each supported IEEE 754 interval type. Interchange encodings for non-IEEE-754 interval types, and on non-IEEE-754 systems, are optional and implementation-defined. If an implementation provides other decoration attributes besides the standard ones, then how it maps them to an interchange encoding is implementation-defined.

ANNEX A

Bibliography (informative)

Bibliographical references are resources that provide additional or helpful material but do not need to be understood or used to implement this standard. Reference to these resources is made for informational use only.

- [B1] Alefeld, G. and Mayer, G., “Interval analysis: theory and applications,” *Journal of Computational and Applied Mathematics*, vol. 121, no. 1–2, pp. 421–464, September 2000.
- [B2] Allen, J. F., “Maintaining knowledge about temporal intervals,” *Communications of the ACM*, vol. 26, no. 11, pp. 832–843, November 1983.
- [B3] ISO/IEC 9899:1999, “1999 standard for Programming Language C (C99).”¹⁵
- [B4] Griewank, A. and Walther, A., *Evaluating Derivatives: Principles and Techniques of Algorithmic Differentiation*, Philadelphia, PA: SIAM, 2nd ed., 2008.
- [B5] Kreinovich, V., “Interval computations web site,” accessed August 28, 2014, <http://cs.utep.edu/interval-comp>.
- [B6] Kulisch, U., *Computer arithmetic and validity: Theory, implementation, and applications*, Berlin: Walter de Gruyter, 2nd ed., 2013.
- [B7] Kulisch, U. and Snyder, V., “The exact dot product as basic tool for long interval arithmetic,” *Computing*, vol. 91, no. 3, pp. 307–313, November 2011.
- [B8] Moore, R. E., *Interval analysis*, Englewood Cliffs: Prentice-Hall, 1966.
- [B9] Nehmeier, M., Siegel, S., and von Gudenberg, J. W., “Parallel detection of interval overlapping,” in *Applied Parallel and Scientific Computing*, Springer, vol. 7134 of *Lecture Notes in Computer Science*, pp. 127–136, 2012.
- [B10] Neumaier, A., “Vienna proposal for interval standardization,” December 2008, accessed August 28, 2014, <http://www.mat.univie.ac.at/~neum/ms/1788.pdf>.
- [B11] The Organization for the Advancement of Structured Information Standards (OASIS), “Conformance requirements for specifications,” accessed August 28, 2014, <https://www.oasis-open.org/committees/ioc>.
- [B12] Pryce, J. D., “The decoration system,” P1788 position paper, version 03.2, June 2011; accessed August 28, 2014, <http://grouper.ieee.org/groups/1788/PositionPapers/DecorationsANJDP.pdf>.
- [B13] Pryce, J. D. and Corliss, G. F., “Interval arithmetic with containment sets,” *Computing*, vol. 78, no. 3, pp. 251–276, November 2006.
- [B14] Rump, S. M., “Verification methods: Rigorous results using floating-point arithmetic,” *Acta Numerica*, vol. 19, pp. 287–449, May 2010.
- [B15] Rump, S. M., “Interval arithmetic over finitely many endpoints,” *BIT Numerical Mathematics*, vol. 52, no. 4, pp. 1059–1075, 2012.

¹⁵ISO/IEC publications are available from the International Organization for Standards, ISO Central Secretariat, 1, ch. de la Voie-Creuse, CP 56, CH-1211 Geneva 20, Switzerland (<http://www.iso.org/>).

- [B16] Sun Microsystems, Oracle Engineered Systems, “Fortran 95 Interval Arithmetic Programming Reference”; accessed 28 Jan 2015, <http://docs.oracle.com/cd/E19422-01/819-3695/>.
- [B17] Wikipedia, “Brouwer fixed-point theorem”; accessed 28 Jan 2015, http://en.wikipedia.org/wiki/Brouwer_fixed-point_theorem.

ANNEX B

The fundamental theorem of decorated interval arithmetic for the set-based flavor (informative)

This Annex states and proves the Fundamental Theorem 6.1 for the set-based flavor, in a computationally verifiable form where the decoration system is used to identify the different cases of that theorem. It thereby proves that a conforming implementation of this flavor can be used to conclude while evaluating a point function that it is everywhere defined, or everywhere defined and continuous, etc., on an input box.

B1. Preliminaries

The variables in an expression. We denote the set of variables occurring in expression f by $V(f)$. This set is defined as follows:

- If f is a single named variable, e.g. x , then $V(f) = \{x\}$. The x is here regarded as a symbol, chosen from some set of allowed symbols.
- If $f = h(g_1, \dots, g_k)$ then

$$V(f) = \bigcup_{i=1}^k V(g_i).$$

In particular, if $k = 0$ (so f is a constant, see below) then $V(f) = \emptyset$.

Argument association by keyword. It is helpful to use both positional and keyword notation for associating dummy arguments to actual arguments when invoking a function

Library operations φ use positional notation $\varphi(v_1, \dots, v_k)$ with the arguments in a fixed order, and the names of the dummy arguments are not important. However the theorem is simpler to state and prove if, for a function defined by expression, the variables occurring in the expression are linked *by name* to their actual arguments. E.g., expression $f = b/d + a$ is deemed to define a function f with arguments a, b, d in no specified order. Keyword notation¹⁶ is illustrated by $f(a = 5, b = 12, d = 4)$, which denotes the evaluation $12/4 + 5$ with result 8.

Formally, positional notation numbers the elements of $V(f)$ from 1 to n and indexes actual arguments x_i to match, making a map $x : i \mapsto x_i$ defined on $\{1, \dots, n\}$; while keyword notation indexes the arguments x_v directly by the variables v , making a map $x : v \mapsto x_v$ defined on $V(f)$.

Such an x is regarded as an actual-argument *vector*. When $U \subseteq V(f)$, notation x_U means the *subvector* comprising those x_v for $v \in U$, i.e., the restriction of map x to subset U of its domain.

Keyword notation helps solve two related problems, of *empty interval inputs* and of *ill-formed interval inputs*. These occur with positional notation when there is an argument that does not actually occur in the expression.

[Example. Make the positional definition $y = f(a, b, c, d) = b/d + a$, where argument c does not occur in the defining expression $f = b/d + a$.

¹⁶Similar to the keyword argument feature in Fortran, the `subs` command in Maple, etc.

- Suppose, with bare interval evaluation, we evaluate $y = f(a, b, c, d)$ where a, b, d are nonempty but $c = \emptyset$. What actually results from evaluating the expression f is $y_C = b/d + a$, which (whether at Level 1 or 2) is in general nonempty.

Using the positional definition, the input box is the product of 4 intervals, $x_P = a \times b \times c \times d = a \times b \times \emptyset \times d = \emptyset$, so that the theoretical range is $y_P = \text{Rge}(f | \emptyset) = \emptyset$.

Using the keyword definition, the input box is the product of 3 intervals x_v for $v \in V(f)$, that is $x_K = a \times b \times d \neq \emptyset$, making the theoretical range to be $y_K = \text{Rge}(f | x_K) \neq \emptyset$. In fact—if using the Level 1 interval versions of operations— y_C coincides with y_K whenever $0 \notin d$.

- A similar case arises with decorated interval evaluation, where inputs a_{da}, b_{db}, d_{dd} are well-formed but c_{dc} is NaI. Positional notation suggests the input box, and hence the result, should be NaI, but this result cannot come by evaluating the expression one operation at a time; the only way to get it is to redefine evaluation to include a preliminary scan that looks for NaI inputs.

]

This example shows that keyword notation defines a theoretical result having a more useful and common-sense relation to the computed result of bare or decorated interval evaluation of an expression, than does positional notation.

Constant functions. It is necessary to clarify the case of a *zero-argument* arithmetic operation h . A point argument of a general k -ary h is a tuple $u = (u_1, \dots, u_k) \in \mathbb{R}^k$. For $k = 0$ this is the empty tuple $()$, which is the unique element of \mathbb{R}^0 . Since \mathbb{R}^0 is a singleton set, a particular $h : \mathbb{R}^0 \rightarrow \mathbb{R}$ is either total—defined everywhere—on \mathbb{R}^0 , with a unique real value; or defined nowhere on \mathbb{R}^0 .

- The total functions are in one to one correspondence with \mathbb{R} . If c is a symbol denoting a real number, the notation $c()$ means the total function on \mathbb{R}^0 with value c .
- The nowhere defined function on \mathbb{R}^0 is by definition the “Not a Number” function NaN.

In this way, each $h : \mathbb{R}^0 \rightarrow \mathbb{R}$ is either identified with a real constant—e.g., the functions $0.1()$ and $\pi()$ are identified with 0.1 and π —or is the unique NaN function.

The Level 1 interval version (the natural interval extension) of a total $c()$ is $\mathbf{c} : \mathbb{R}^0 \rightarrow \overline{\mathbb{R}}$ with value $[c, c]$; that of the NaN function is the “Not an Interval” function NaI with value \emptyset .

The decorated version of $c()$ produces $[c, c]_{\text{com}}$; the decorated version of NaI produces \emptyset_{ill} .

At Level 2 one must consider finite precision. E.g., if \mathbb{T} is the type inf-sup binary64, and $c = 0.1$, then $[c, c]$ is not a \mathbb{T} -interval and must be enclosed in a larger interval. A general interval extension of the constant $c()$ is any interval $[\underline{c}, \bar{c}]$ that contains c . Its decorated version is $[\underline{c}, \bar{c}]_{dc}$ where

$$dc = \begin{cases} \text{any decoration} \supseteq \text{com}, & \text{if } [\underline{c}, \bar{c}] \text{ is bounded,} \\ \text{any decoration} \supseteq \text{dac}, & \text{if } [\underline{c}, \bar{c}] \text{ is unbounded.} \end{cases}$$

Initialization of intervals. Define a decorated interval $\mathbf{X} = x_{dx}$ to be **validly initialized** if either $x_{dx} = \emptyset_{\text{ill}}$ or $p_{dx}(\text{Id}, x)$ holds, the latter being equivalent to $dx \supseteq \text{Dec}(\text{Id} | x)$. An \mathbf{X} for which $dx \neq \text{ill}$ is called **tightly initialized** if $dx = \text{Dec}(\text{Id} | x)$. A decorated box is validly, or tightly, initialized if each of its decorated interval components is so.

The possible valid initializing decorations for the various kinds of interval are shown in the table below, with the tight initialization case underlined.

Interval x	Valid initializing decoration
nonempty and bounded,	<u>com</u> , dac, def, trv
unbounded,	<u>dac</u> , def, trv
empty,	<u>trv</u>
any nonempty bare interval; or <u>empty</u>	<u>ill</u>

The tightly initialized cases coincide with those obtainable from a decorated constructor. If it succeeds, the decoration is as set by the **newDec** function. If it fails, the result is \emptyset_{ill} .

B2. The theorem

THEOREM B2.1 (Fundamental Theorem of Decorated Interval Arithmetic, FTDIA). *Let f be an arithmetic expression denoting a real function f . Suppose f is evaluated, possibly in finite precision, on a validly initialized decorated box $\mathbf{X} = \mathbf{x}_{dx}$, with components indexed over $V(f)$, to give result $\mathbf{Y} = \mathbf{y}_{dy}$.*

If some component of \mathbf{X} is decorated `ill` then

$$dy = \text{ill}. \quad (31)$$

If no component of \mathbf{X} is decorated `ill`, and none of the operations φ of f is an everywhere undefined function `Un`, then $dy \neq \text{ill}$ and

$$\mathbf{y} \supseteq \text{Rge}(f | \mathbf{x}) \quad \text{and} \quad (32)$$

$$dy \supseteq \text{Dec}(f | \mathbf{x}). \quad (33)$$

In more primitive terms, (32, 33) say that $f(x) \in \mathbf{y}$ for all $x \in \mathbf{x}$, and $p_{dy}(f, \mathbf{x})$ holds.

Proof. We proceed by induction on the number of operations in f .

Base case. This is the case when f has zero operations. That is, it is a single variable, say x , and denotes the identity function $f = \text{Id} : x \mapsto x$ ($x \in \mathbb{R}$). Then the output \mathbf{y}_{dy} equals the input decorated interval $\mathbf{X} = \mathbf{x}_{dx}$. **NOTE**—In finite precision it appears this need not be the case because of e.g., type conversion. However this is to be treated as a separate unary arithmetic operation—a decorated interval extension of `Id`—which must obey the usual rules for arithmetic operations.

Thus if $dx = \text{ill}$, then $dy = \text{ill}$ and (31) holds. If $dx \neq \text{ill}$, then $\text{Rge}(f | \mathbf{x}) = \text{Rge}(\text{Id} | \mathbf{x}) = \mathbf{x} = \mathbf{y}$ so (32) holds. Also $\text{ill} \neq dy = dx \supseteq \text{Dec}(\text{Id} | \mathbf{x}) = \text{Dec}(f | \mathbf{x})$ by the definition of valid initialization, so (33) holds. Hence the Theorem is true in this case.

Induction step. Suppose the Theorem has been established for all expressions with fewer than N operations, and let f have N operations.

By the definition of an expression, f is obtained from k expressions g_i , $i = 1, \dots, k$, through an arithmetic operation h of some arity $k \geq 0$. That is, $f = h(g_1, \dots, g_k)$, where necessarily each g_i has fewer than N operations. (Note the case $k = 0$, of a constant function, is included in this argument.)

Let $V_i = V(g_i)$, that is, the set of variables occurring in g_i . Then $V(f) = \bigcup_{i=1}^k V(g_i)$. When f , h , and the g_i are regarded as defining point functions, this means $f(x) = h(g_1(x_{V_1}), \dots, g_k(x_{V_k}))$, where x_{V_i} is $V(g_i)$ regarded as a vector.

By the inductive hypothesis the theorem holds for each g_i . That is, we have computed an interval \mathbf{g}_i and a decoration dg_i , such that by (32, 33),

$$\mathbf{g}_i \supseteq \text{Rge}(g_i | \mathbf{x}_{V_i}) \quad \text{and} \quad (34)$$

$$p_{dg_i}(g_i, \mathbf{x}_{V_i}) \quad \text{holds}. \quad (35)$$

By the definition of a decorated interval version of f , \mathbf{y}_{dy} is computed using a decorated interval extension of h , hence by the definition in 11.6,

$$\mathbf{y} \supseteq \text{Rge}(h | \mathbf{g}) \quad \text{and} \quad (36)$$

$$dy = \min\{dh, dg_1, \dots, dg_k\} \quad (37)$$

for some dh such that

$$p_{dh}(h, \mathbf{g}) \quad \text{holds}. \quad (38)$$

Corresponding to the different meanings of the decorations, (31, 32, 33) are verified case by case.

Case $df = \text{ill}$. Then either some $dg_i = \text{ill}$ or $dh = \text{ill}$.

— If $dg_i = \text{ill}$, by (35) $\text{Dom } g_i$ is empty, and therefore $\text{Dom } f$ is empty.

— If $dh = \text{ill}$, then by (38) $\text{Dom } h$ is empty, and therefore $\text{Dom } f$ is empty.

In both cases, $\text{Rge}(f | \mathbf{x}) = \emptyset$ and $p_{\text{ill}}(f, \mathbf{x})$ holds. Since the computed result is \emptyset_{ill} , (32, 33) hold.

Case $df = \text{trv}$. This is always true, and nothing needs to be shown.

Case $df = \text{def}$. Then each $dg_i \geq \text{def}$. Hence, for all $i = 1, \dots, k$, \mathbf{x}_{V_i} is a nonempty subset of $\text{Dom } g_i$, and by (34), $\mathbf{g}_i \supseteq \text{Rge}(g_i | \mathbf{x}_{V_i})$. Since also $dh \geq \text{def}$, \mathbf{g} is a nonempty subset of $\text{Dom } h$. Hence f is everywhere defined on \mathbf{x} , that is $p_{\text{def}}(f, \mathbf{x})$ holds, and (33) holds.

It remains to show (32). For any $v \in \text{Rge}(f | \mathbf{x})$, there is $x \in \mathbf{x}$ such that $v = f(x)$. Then there exist $x_{V_i} \in \mathbf{x}_{V_i}$ for $i = 1, \dots, k$ such that

$$v = f(x) = h(g_1(x_{V_1}), \dots, g_k(x_{V_k})).$$

Denote $u_i = g_i(x_{V_i})$ and $u = (u_1, \dots, u_k)$. Since $u_i \in \text{Rge}(g_i | \mathbf{x}_{V_i}) \subseteq \mathbf{g}_i$ and $u \in \mathbf{g}$, we have

$$\begin{aligned} v = f(x) &= h(u_1, \dots, u_k) \\ &\in \{h(u) \mid u \in \mathbf{g}\} \\ &= \text{Rge}(h | \mathbf{g}) \subseteq \mathbf{y}. \end{aligned}$$

Since v was arbitrary, this proves (32). It holds in the next two cases since **dac** and **com** are stronger than **def**.

Case $df = \text{dac}$. This is as the **def** case with the addition that the restriction of each g_i to \mathbf{x}_{V_i} is everywhere continuous, and the restriction of h to \mathbf{g} is everywhere continuous. Hence the restriction of f to \mathbf{x} is everywhere defined and continuous, and (33) holds.

Case $df = \text{com}$. Then each $dg_i = \text{com}$. Hence, for all $i = 1, \dots, k$, \mathbf{x}_{V_i} is a bounded, nonempty subset of $\text{Dom } g_i$, g_i is continuous at each point in \mathbf{x}_{V_i} , and the computed \mathbf{g}_i is bounded. Since also $dh = \text{com}$, h is defined and continuous at each point of \mathbf{g} , f is defined and continuous at each point of \mathbf{x} , and the computed \mathbf{y} is bounded; (33) holds.

This completes the induction step and the proof. □

From the details of the proof one sees

COROLLARY B2.2. *If the result \mathbf{y}_{dy} has $dy = \text{ill}$, and the expression has at least one operation, then $\mathbf{y} = \emptyset$.*

This gives an argument in favor of assuming that if the decoration is **ill** then the interval part is *always* empty, i.e., of defining **NaI** to be \emptyset_{ill} .

Consensus

WE BUILD IT.

Connect with us on:



Facebook: <https://www.facebook.com/ieeesa>



Twitter: @ieeesa



LinkedIn: <http://www.linkedin.com/groups/IEEE-SA-Official-IEEE-Standards-Association-1791118>



IEEE-SA Standards Insight blog: <http://standardsinsight.com>



YouTube: IEEE-SA Channel