

INTERNATIONAL
STANDARD

ISO/IEC
10967-2

First edition
2001-08-15

**Information technology — Language
independent arithmetic —**

**Part 2:
Elementary numerical functions**

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



Reference number
ISO/IEC 10967-2:2001(E)

© ISO/IEC 2001

PDF disclaimer

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

Adobe is a trademark of Adobe Systems Incorporated.

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

© ISO/IEC 2001

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

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

Printed in Switzerland

Contents

Foreword	viii
Introduction	ix
1 Scope	1
1.1 Inclusions	1
1.2 Exclusions	2
2 Conformity	2
3 Normative references	3
4 Symbols and definitions	4
4.1 Symbols	4
4.1.1 Sets and intervals	4
4.1.2 Operators and relations	4
4.1.3 Mathematical functions	5
4.1.4 Exceptional values	5
4.1.5 Datatypes	6
4.2 Definitions of terms	7
5 Specifications for integer and floating point operations	10
5.1 Basic integer operations	10
5.1.1 The integer <i>result</i> and <i>wrap</i> helper functions	10
5.1.2 Integer maximum and minimum	11
5.1.3 Integer diminish	11
5.1.4 Integer power and arithmetic shift	12
5.1.5 Integer square root	12
5.1.6 Divisibility tests	12
5.1.7 Integer division (with floor, round, or ceiling) and remainder	13
5.1.8 Greatest common divisor and least common positive multiple	13
5.1.9 Support operations for extended integer range	14
5.2 Basic floating point operations	15
5.2.1 The rounding and floating point <i>result</i> helper functions	15
5.2.2 Floating point maximum and minimum	17
5.2.3 Floating point diminish	18
5.2.4 Floor, round, and ceiling	19
5.2.5 Remainder after division with round to integer	20
5.2.6 Square root and reciprocal square root	20
5.2.7 Multiplication to higher precision floating point datatype	20
5.2.8 Support operations for extended floating point precision	21
5.3 Elementary transcendental floating point operations	22
5.3.1 Maximum error requirements	22
5.3.2 Sign requirements	23
5.3.3 Monotonicity requirements	23
5.3.4 The <i>result</i> * helper function	23
5.3.5 Hypotenuse	24
5.3.6 Operations for exponentiations and logarithms	24

5.3.6.1	Integer power of argument base	24
5.3.6.2	Natural exponentiation	25
5.3.6.3	Natural exponentiation, minus one	26
5.3.6.4	Exponentiation of 2	27
5.3.6.5	Exponentiation of 10	27
5.3.6.6	Exponentiation of argument base	28
5.3.6.7	Exponentiation of one plus the argument base, minus one	29
5.3.6.8	Natural logarithm	29
5.3.6.9	Natural logarithm of one plus the argument	30
5.3.6.10	2-logarithm	30
5.3.6.11	10-logarithm	31
5.3.6.12	Argument base logarithm	31
5.3.6.13	Argument base logarithm of one plus each argument	32
5.3.7	Introduction to operations for trigonometric elementary functions	32
5.3.8	Operations for radian trigonometric elementary functions	33
5.3.8.1	Radian angle normalisation	34
5.3.8.2	Radian sine	35
5.3.8.3	Radian cosine	35
5.3.8.4	Radian tangent	36
5.3.8.5	Radian cotangent	36
5.3.8.6	Radian secant	37
5.3.8.7	Radian cosecant	37
5.3.8.8	Radian cosine with sine	38
5.3.8.9	Radian arc sine	38
5.3.8.10	Radian arc cosine	38
5.3.8.11	Radian arc tangent	39
5.3.8.12	Radian arc cotangent	40
5.3.8.13	Radian arc secant	41
5.3.8.14	Radian arc cosecant	41
5.3.8.15	Radian angle from Cartesian co-ordinates	42
5.3.9	Operations for trigonometrics with given angular unit	43
5.3.9.1	Argument angular-unit angle normalisation	43
5.3.9.2	Argument angular-unit sine	44
5.3.9.3	Argument angular-unit cosine	45
5.3.9.4	Argument angular-unit tangent	45
5.3.9.5	Argument angular-unit cotangent	46
5.3.9.6	Argument angular-unit secant	47
5.3.9.7	Argument angular-unit cosecant	47
5.3.9.8	Argument angular-unit cosine with sine	48
5.3.9.9	Argument angular-unit arc sine	48
5.3.9.10	Argument angular-unit arc cosine	48
5.3.9.11	Argument angular-unit arc tangent	49
5.3.9.12	Argument angular-unit arc cotangent	50
5.3.9.13	Argument angular-unit arc secant	51
5.3.9.14	Argument angular-unit arc cosecant	51
5.3.9.15	Argument angular-unit angle from Cartesian co-ordinates	52
5.3.10	Operations for angular-unit conversions	53
5.3.10.1	Converting radian angle to argument angular-unit angle	53

5.3.10.2	Converting argument angular-unit angle to radian angle	54
5.3.10.3	Converting argument angular-unit angle to (another) argument angular-unit angle	55
5.3.11	Operations for hyperbolic elementary functions	56
5.3.11.1	Hyperbolic sine	56
5.3.11.2	Hyperbolic cosine	56
5.3.11.3	Hyperbolic tangent	57
5.3.11.4	Hyperbolic cotangent	58
5.3.11.5	Hyperbolic secant	58
5.3.11.6	Hyperbolic cosecant	59
5.3.11.7	Inverse hyperbolic sine	59
5.3.11.8	Inverse hyperbolic cosine	60
5.3.11.9	Inverse hyperbolic tangent	60
5.3.11.10	Inverse hyperbolic cotangent	60
5.3.11.11	Inverse hyperbolic secant	61
5.3.11.12	Inverse hyperbolic cosecant	61
5.4	Operations for conversion between numeric datatypes	62
5.4.1	Integer to integer conversions	63
5.4.2	Floating point to integer conversions	63
5.4.3	Integer to floating point conversions	64
5.4.4	Floating point to floating point conversions	64
5.4.5	Floating point to fixed point conversions	65
5.4.6	Fixed point to floating point conversions	66
5.5	Numerals as operations in a programming language	67
5.5.1	Numerals for integer datatypes	67
5.5.2	Numerals for floating point datatypes	68
6	Notification	68
6.1	Continuation values	69
7	Relationship with language standards	69
8	Documentation requirements	70
Annex A	(normative) Partial conformity	73
A.1	Maximum error relaxation	73
A.2	Extra accuracy requirements relaxation	74
A.3	Relationships to other operations relaxation	74
A.4	Very-close-to-axis angular normalisation relaxation	74
A.5	Part 1 requirements relaxation	75
Annex B	(informative) Rationale	77
B.1	Scope	77
B.1.1	Inclusions	77
B.1.2	Exclusions	78
B.2	Conformity	78
B.2.1	Validation	79
B.3	Normative references	79
B.4	Symbols and definitions	79

B.4.1	Symbols	79
B.4.1.1	Sets and intervals	79
B.4.1.2	Operators and relations	80
B.4.1.3	Mathematical functions	80
B.4.1.4	Exceptional values	80
B.4.1.5	Datatypes	81
B.4.2	Definitions of terms	81
B.5	Specifications for the numerical functions	81
B.5.1	Basic integer operations	82
B.5.1.1	The integer <i>result</i> and <i>wrap</i> helper functions	82
B.5.1.2	Integer maximum and minimum	82
B.5.1.3	Integer diminish	82
B.5.1.4	Integer power and arithmetic shift	83
B.5.1.5	Integer square root	83
B.5.1.6	Divisibility tests	83
B.5.1.7	Integer division (with floor, round, or ceiling) and remainder	83
B.5.1.8	Greatest common divisor and least common positive multiple	84
B.5.1.9	Support operations for extended integer range	84
B.5.2	Basic floating point operations	84
B.5.2.1	The rounding and floating point <i>result</i> helper functions	86
B.5.2.2	Floating point maximum and minimum	86
B.5.2.3	Floating point diminish	86
B.5.2.4	Floor, round, and ceiling	86
B.5.2.5	Remainder after division and round to integer	87
B.5.2.6	Square root and reciprocal square root	87
B.5.2.7	Multiplication to higher precision floating point datatype	88
B.5.2.8	Support operations for extended floating point precision	88
B.5.3	Elementary transcendental floating point operations	89
B.5.3.1	Maximum error requirements	89
B.5.3.2	Sign requirements	90
B.5.3.3	Monotonicity requirements	90
B.5.3.4	The <i>result</i> * helper function	90
B.5.3.5	Hypotenuse	91
B.5.3.6	Operations for exponentiations and logarithms	91
B.5.3.7	Introduction to operations for trigonometric elementary functions	93
B.5.3.8	Operations for radian trigonometric elementary functions	94
B.5.3.9	Operations for trigonometrics with given angular unit	96
B.5.3.10	Operations for angular-unit conversions	97
B.5.3.11	Operations for hyperbolic elementary functions	98
B.5.4	Operations for conversion between numeric datatypes	98
B.5.5	Numerals as operations in a programming language	99
B.5.5.1	Numerals for integer datatypes	99
B.5.5.2	Numerals for floating point datatypes	99
B.6	Notification	100
B.6.1	Continuation values	100
B.7	Relationship with language standards	101
B.8	Documentation requirements	101

Annex C (informative) Example bindings for specific languages	103
C.1 Ada	104
C.2 BASIC	110
C.3 C	114
C.4 C++	120
C.5 Fortran	126
C.6 Haskell	132
C.7 Java	137
C.8 Common Lisp	142
C.9 ISLisp	147
C.10 Modula-2	152
C.11 Pascal and Extended Pascal	157
C.12 PL/I	162
C.13 SML	167
Annex D (informative) Bibliography	173
Annex E (informative) Possible changes to part 1	177

Foreword

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

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

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

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

International Standard ISO/IEC 10967-2 was prepared by Joint Technical Committee ISO/IEC JTC 1, *Information technology*, Subcommittee SC 22, *Programming languages, their environments and system software interfaces*.

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

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

Additional parts will specify other arithmetic datatypes or arithmetic operations.

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

Introduction

The aims

Portability is a key issue for scientific and numerical software in today's heterogeneous computing environment. Such software may be required to run on systems ranging from personal computers to high performance pipelined vector processors and massively parallel systems, and the source code may be ported between several programming languages. Part 1 of ISO/IEC 10967 specifies the basic properties of integer and floating point types that can be relied upon in writing portable software.

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

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

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

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

The content

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

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

The benefits

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

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

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

Programs will be able to determine (at run time) the crucial numeric properties of the implementation. They will be able to reject unsuitable implementations, and (possibly) to correctly characterize the accuracy of their own results. Programs will be able to detect (and possibly correct for) exceptions in arithmetic processing.

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

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

Information technology — Language independent arithmetic —

Part 2: Elementary numerical functions

1 Scope

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

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

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

1.1 Inclusions

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

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

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

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

This part also provides specifications for:

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

1.2 Exclusions

This part provides no specifications for

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

Furthermore, this part does not provide specifications for how the operations should be implemented or which algorithms are to be used for the various operations.

2 Conformity

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

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

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

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

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

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

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

An implementation is permitted to have modes of operation that do not conform to this part. A conforming implementation shall specify how to select the modes of operation that ensure conformity. However, a mode of operation that conforms to this part should be the default mode of operation.

NOTES

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

3 Normative references

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

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

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

NOTE – See also annex E.