

## SLOVENSKI STANDARD SIST HD 592 S1:1997

01-avgust-1997

Binarna ari 60559:1989	Binarna aritmetika s plavajočo vejico za mikroprocesorske sisteme (IEC 60559:1989)						
Binary float	Binary floating-point arithmetic for microprocessor systems						
Binäre Gleit	Binäre Gleitpunkt-Arithmetik für Mikroprozessor-Systeme						
Arithmétiqu	Arithmétique binaire en virgule flottante pour systèmes à microprocesseurs (standards.iteh.ai)						
Ta slovenski standard je istoveten <u>z:</u> <u>HD 592 S1:1991</u>							
3590540f2fd1/sist-hd-592-s1-1997							
<u>ICS:</u> 35.160	Mikroprocesorski sistemi	Microprocessor systems					

SIST HD 592 S1:1997

en



## iTeh STANDARD PREVIEW (standards.iteh.ai)

<u>SIST HD 592 S1:1997</u> https://standards.iteh.ai/catalog/standards/sist/c49dc8af-8c2b-4e24-a213-3590540f2fd1/sist-hd-592-s1-1997

#### HARMONIZATION DOCUMENT

HD 592 S1

DOCUMENT D'HARMONISATION

HARMONISIERUNGSDOKUMENT

May 1991

UDC 621.382.049.77.037.372

Descriptors: Binary floating point arithmetic, microprocessor systems

#### ENGLISH VERSION

#### BINARY FLOATING-POINT ARITHMETIC FOR MICROPROCESSOR SYSTEMS (IEC 559:1989)

Arithmétique binaire en virgule flottante pour systèmes à microprocesseurs (CEI 559:1989) Binäre Gleitpunkt-Arithmetik für Mikroprozessor-Systeme (IEC 559:1989)

This Harmonization Document was approved by CENELEC on 1991-03-15. CENELEC members are bound to comply with the CEN/CENELEC Internal Regulations which stipulate the conditions for implementation of this Harmonization Document on a national level.

Up-to-date lists and bibliographical references concerning national implementation may be obtained on application to the Central Secretariat or to any CENELEC member.

This Harmonization Document exists in three official versions (English, French, German).

CENELEC members are the national electrotechnical committees of Austria, Belgium, Denmark, Finland, France, Germany, Greece, Iceland, Ireland, Italy, Luxembourg, Netherlands, Norway, Portugal, Spain, Sweden, Switzerland and United Kingdom.

## iTeh SPEANDARD PREVIEW

European Committee for Electrotechnical Standardization Comité Européen de Normalisation Electrotechnique Europäisches Komitee für Elektrotechnische Normung <u>SISTHD 592 S1:1997</u>

Central Secretarilardrich. deatstafts/35t/c48d9050cBrussels3-3590540f2fd1/sist-hd-592-s1-1997

© 1991 Copyright reserved to CENELEC members

Ref. No. HD 592 S1:1991 E

Page 2 HD 592 S1:1991

#### FOREWORD

The CENELEC questionnaire procedure, performed for finding out whether or not the International Standard IEC 559:1989 could be accepted without textual changes, has shown that no CENELEC common modifications were necessary for the acceptance as Harmonization Document. The reference document was submitted to the CENELEC members for formal vote and was approved by CENELEC as HD 592 S1 on 15 March 1991.

The following dates were fixed:

-	latest date of announcement of the HD at national level	(doa)	1991-09-01
-	latest date of publication of a harmonized national standard	(dop)	1992-03-01
-	latest date of withdrawal of conflicting national standards	(dow)	1992-03-01

#### ENDORSEMENT NOTICE

The text of the International Standard IEC 559:1989 was approved by CENELEC as a Harmonization Document without any modification.

\_\_\_\_\_\_

## iTeh STANDARD PREVIEW (standards.iteh.ai)

<u>SIST HD 592 S1:1997</u> https://standards.iteh.ai/catalog/standards/sist/c49dc8af-8c2b-4e24-a213-3590540f2fd1/sist-hd-592-s1-1997

# NORME INTERNATIONALE INTERNATIONAL STANDARD

# CEI IEC 60559

Deuxième édition Second edition 1989-01

## Arithmétique binaire en virgule flottante pour systèmes à microprocesseurs

## i Binary floating-point arithmetic for microprocessor systems

<u>SIST HD 592 S1:1997</u> https://standards.iteh.ai/catalog/standards/sist/c49dc8af-8c2b-4e24-a213-3590540f2fd1/sist-hd-592-s1-1997

© IEC 1989 Droits de reproduction réservés - Copyright - all rights reserved

Aucune partie de cette publication ne peut être reproduite ni utilisée sous quelque forme que ce soit et par aucun procédé, électronique ou mécanique, y compris la photocopie et les microfilms, sans l'accord écrit de l'éditeur. 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 the publisher.

International Electrotechnical Commission3, rue de Varembé Geneva, SwitzerlandTelefax: +41 22 919 0300e-mail: inmail@iec.chIEC web site http: //www.iec.ch

ΞO



Commission Electrotechnique Internationale International Electrotechnical Commission Международная Электротехническая Комиссия



S

Pour prix, voir catalogue en vigueur For price, see current catalogue

-

## - 3 -

### CONTENTS

			Page		
FOREWORD					
PREFACE					
Cla					
	<b>4</b> 36				
1.	. Scope				
	1.1	Implementation objectives	7.		
	1.2	Exclusions	7		
2.	Defir	nitions	7		
3.	Formats				
-	3.1	Sets of values	13 ·		
	3.2	Basic formats	15		
	3.3	Extended formats	17		
	3.4		17		
4.	Rour	nding	19		
	4.1	Round to nearest	19		
	4.2 4.3	Rounding precision TANDARD PREVIEW	19 19		
5.	Operations				
	5.1	Arithmetic	21		
	5.2	Square root	23		
	5.3 5.4	Floating-pointatormathconversions/sist/c49dc8af-8c2b-4c24-a213	23		
	5.5	Round floating-point number to integral value	23		
	5 C	Pinany ( ) desimal convension	22		
	5.7	Comparison	23		
6.	Infir	ity, NaNs and signed zero	31		
	61	Infinity arithmetic	31		
	6.2	Operations with NaNs	31		
	6.3	The sign bit	33		
7.	Exce	ptions	33		
	7.1	Invalid operations	33		
	7.2	Division by zero	35		
	7.4	Underflow	37		
	7.5	Inexact	39		
8.	Trap	>s	39		
	8.1	Trap handler	41		
	8.2	Precedence	41		
AP	PEND	IX A - Recommended functions and predicates	43		

#### INTERNATIONAL ELECTROTECHNICAL COMMISSION

#### BINARY FLOATING-POINT ARITHMETIC FOR MICROPROCESSOR SYSTEMS

#### FOREWORD

- The formal decisions or agreements of the IEC on technical matters, prepared by Technical Committees on which all the National Committees having a special interest therein are represented, express, as nearly as possible, an international consensus of opinion on the subjects dealt with.
- 2) They have the form of recommendations for international use and they are accepted by the National Committees in that sense.
- 3) In order to promote international unification, the IEC expresses the wish that all National Committees should adopt the text of the IEC recommendation for their national rules in so far as national conditions will permit. Any divergence between the IEC recommendation and the corresponding national rules should, as far as possible, be clearly indicated in the latter arcs.itch.al)

#### SPREPACES1:1997

#### https://standards.iteh.ai/catalog/standards/sist/c49dc8af-8c2b-4e24-a213-

This standard has been prepared by Sub-Committee 47B: Microprocessor systems, of IEC Technical Committee No. 47: Semiconductor devices. (This Sub-Committee has been taken over by ISO/IEC JTC 1.)

This second edition of IEC Publication 559 replaces the first edition issued in 1982.

The text of this standard is based on the following documents:

Six Months' Rule	Report on Voting
47B(CO)19	47B(CO)26

Full information on the voting for the approval of this standard can be found in the Voting Report indicated in the above table.

### BINARY FLOATING-POINT ARITHMETIC FOR MICROPROCESSOR SYSTEMS

#### 1. Scope

#### 1.1 Implementation objectives

It is intended that an implementation of a floating-point system conforming to this standard can be realized entirely in software, entirely in hardware, or in any combination of software and hardware. It is the environment that the programmer or user of the system sees that conforms or fails to conform to this standard. Hardware components that require software support to conform shall not be said to conform apart from such software.

#### 1.2 Inclusions

This standard specifies:

- 1) basic and extended floating-point number formats;
- 2) add, subtract, multiply divide, square vroot, vremainder and compare operations; (standards.iteh.ai)
- 3) conversions between integer and floating-point numbers; <u>SIST HD 592 S1:1997</u>
- 4) conversions between different/sloading\_ippint formats;
- 5) conversions between basic format floating-point numbers and decimal strings, and
- 6) floating-point exceptions and their handling, including nonnumbers (NaNs).

#### 1.3 Exclusions

This standard does not specify:

- 1) formats of decimal strings and integers;
- 2) interpretation of the signs and significant fields of NaNs, or
- 3) binary  $\leftrightarrow$  decimal conversions to and from extended formats.

#### 2. Definitions

#### Biased exponent

The sum of the exponent and a constant (bias) chosen to make the biased exponent's range non-negative.

#### Binary floating-point number

A bit-string characterized by three components: a sign, a signed exponent, and a significand. Its numerical value, if any, is the signed product of its significand and two raised to the power of its exponent. In this standard a bit-string is not always distinguished from a number it may represent.

#### Denormalized number

A nonzero floating-point number, the exponent of which has a reserved value, usually the format's minimum, and the explicit or implicit leading significant bit of which is zero.

#### Destination

The location for the result of a binary or unary operation. The destination may be either explicitly designated by the user or implicitly supplied by the system (e.g. intermediate results in sub-expressions or arguments for procedures). Some languages place the results of intermediate calculations in destinations beyond the user's control. Nonetheless, this standard defines the result of an operation in terms of that destination's format as well as the operands' values.

## Exponent iTeh STANDARD PREVIEW

The component of **Spinarya floating point**) number that normally signifies the integer power to which two is raised in determining the value of the represented number Occasionally the exponent is called the signed or numbiased, exponent and add/sist/c49dc8af-8c2b-4e24-a213-

3590540f2fd1/sist-hd-592-s1-1997

#### Fraction

The field of the significand that lies to the right of its implied binary point.

#### Mode

A variable that a user may set, sense, save and restore, to control the execution of subsequent arithmetic operations. The default mode is the mode that a program can assume to be in effect unless an explicitly contrary statement is included either in the program or in its specification.

The following modes shall be implemented:

- 1) rounding, to control the direction of rounding errors, and in certain implementations.
- 2) rounding precision, to shorten the precision of results. The implementor may, at his option, implement the following modes:
- 3) traps disabled/enabled, to handle exceptions.

#### 559 © IEC

#### - 11 -

#### NaN

Not a number; a symbolic entity encoded in floating-point format. There are two types of NaNs (see 6.2). Signalling NaNs signal the invalid operation exception (see 7.1) whenever they appear as operands. Quiet NaNs propagate through almost every arithmetic operation without signalling exceptions.

#### Result

The bit-string (usually representing a number) that is delivered to the destination.

#### Significant

The component of a binary floating-point number which consists of an explicit or implicit leading bit to the left of its implied binary point and a fraction field to the right.

#### Shall

The word "shall" signifies that which is obligatory in any conforming implementation.

#### Should

**iTeh STANDARD PREVIEW** The word "should" signifies that which is strongly recommended as keeping with a the a intent left athe standard, being in although architectural or other constraints beyond the scope of this standard may, on occasion, render the recommendations impractical.

Status flag

https://standards.iteh.ai/catalog/standards/sist/c49dc8af-8c2b-4e24-a213-3590540f2fd1/sist-hd-592-s1-1997

A variable that may take two states, set and clear. A user may clear a flag, copy it, or restore it to a previous state. When set, a status flag may contain additional system-dependent information, possibly inaccessible to some users. The operations of this standard may, as a side-effect, set some of the following flags: inexact result, underflow, overflow, divide by zero and invalid operation.

#### User

Any person, hardware, or program not itself specified by this standard, having access to and controlling those operations of the programming environment specified in this standard.

#### 3. Formats

This standard defines four floating-point formats in two groups, basic and extended, each having two widths, single and double. The standard levels of implementation are distinguished by the combinations of formats supported.

559 © IEC

3.1 Sets of values

This sub-clause concerns only the numerical values representable within a format, not the encodings which are the subject of the following sub-clauses. The only values representable in a chosen format are those specified via the following three integer parameters:

P = number of significant bits (precision)

 $E_{\text{max}}$  = maximum exponent, and

 $E_{\min}$  = minimum exponent

Each format's parameters are displayed in Table 1. Within each format just the following entities shall be provided:

Numbers of the form  $(-1)^{s_2} (b_0 b_1 b_2 \dots b_{p-1})$ 

where:

```
s is 0 or 1;
E is any integer between E and E inclusive, and each b is 0
or 1.
```

```
Two infinities, to and Standard PREVIEW at least one signalling NaN, and at least one quiet NaN. (standards.iteh.ai)
```

	Format				
Parameter	Single	Single Extended	Double	Double Extended	
Р	24	≥32	53	≥64	
E max	+127	≥+1 023	+1 023	≥+16 383	
E min	-126	≤-1 022	-1 022	≤-16 382	
Exponent bias	+127	Unspeci- fied	+1 023	Unspeci- fied	
Exponent width (bits)	8	≥11	11	≥15	
Format width (bits)	32	≥43	64	≥79	

Table	1 -	Summa	ary of	forma	t:1997	ramet	ers	
https	://stan	dards.iteh.a	i/catalog/s	standards	/sist/c49	dc8af-8	c2b-4e24	-a213
		350	0.05/0.60 fd	1/gist hd	502 c1	1007		

The foregoing description enumerates some values redundantly, for example:

$$2^{0}(1.0) = 2^{1}(0.1) = 2^{2}(0.01) = \dots$$

However, the encodings of such nonzero values may be redundant only in extended formats (see 3.3). The nonzero values of the form  $\pm 2^{E}$ min  $(0 \cdot b_1 b_2 \dots b_{p-1})$  are called denormalized. Reserved exponents