

IEEE Working Group P3109 Interim Report on Binary Floating-point Formats for Machine Learning

Initial release: 18 September 2023
Second interim release: 29 October 2024
Version 0.9.0: 2024-10-29
(compiled 2024-10-29)

DRAFT DOCUMENT

To cite this report please see the CITATION.cff file
found at <https://github.com/P3109/Public>.

Copyright © 2024 by The Institute of Electrical and Electronics Engineers, Inc.
Three Park Avenue
New York, New York 10016-5997, USA
All rights reserved.

This document is subject to change. USE AT YOUR OWN RISK! IEEE copyright statements SHALL NOT BE REMOVED from this draft, or modified in any way. Because this is an unapproved draft, this document must not be utilized for conformance / compliance purposes.

Contents

1	Contributors	4
2	Introduction	5
2.1	Typographical conventions and notation	5
2.2	Scope	5
3	Formats	6
3.1	Format parameters and value sets	6
3.1.1	Computation of bias	6
3.1.2	Computation of bias, $SE = 1$	7
3.1.3	Computation of bias, $SE = 0$	7
3.1.4	Observations on P3109 value sets	7
3.2	Encodings and special values	8
3.3	Subnormals	8
3.4	Not a number (NaN)	9
3.5	Zero	9
3.6	Infinities	10
3.7	Extremal values for $K = 8$	10
3.8	Equivalences	10
4	Operations	11
4.1	Operation variants and superset specifications	12
4.2	Notation and definitions	12
4.2.1	Mathematical notations	12
4.2.2	The set of extended reals	12
4.3	Projection specifications	13
4.4	Non-requirement for operator side effects	13
4.5	Operator definition template	14
4.6	Functions	15
4.6.1	Decode	15
4.6.2	Project	16
4.6.3	RoundToPrecision	17
4.6.4	Saturate	18
4.6.5	Encode	19
4.7	Conversion operations	20
4.7.1	Specification of IEEE Std 754 formats	20
4.7.2	Conversion from IEEE Std 754 formats to P3109	21
4.7.3	Conversion from P3109 to IEEE Std 754	22
4.7.4	Conversion from P3109 to P3109	23
4.8	Arithmetic operations	24
4.8.1	Unary sign operations	24
4.8.2	Binary sign operations	25
4.8.3	Binary arithmetic operations	26
4.8.4	Unary mathematical operations	27

4.8.5	Scaled addition	28
4.8.6	Scaled multiplication	29
4.9	Mixed IEEE Std 754 and P3109 operations	30
4.9.1	Scaled Fused Multiply Add	31
4.10	Comparisons, predicates, and classification	32
4.10.1	Minimum and Maximum	32
4.10.2	Comparisons	33
4.10.3	Predicates and classification	34
4.10.4	Classifier operation	35
4.10.5	Total order predicate	36
4.10.6	Comparison predicates	37
Appendices		38
A	Rationales	38
A.1	Use of infinity in computation of attention masks	38
A.2	Eight Bit Formats	38
A.3	Essential Formats	39
A.3.1	binary8p3	39
A.3.2	binary8p4	39
A.3.3	binary8p5	39
A.3.4	binary8p6	39
B	External Formats	40
C	8-bit Value Tables	40
C.1	Value Table: P1, $e_{min} = -62, e_{max} = 63$ (nonsymmetric)	41
C.2	Value Table: P2, $e_{min} = -31, e_{max} = 31$	42
C.3	Value Table: P3, $e_{min} = -15, e_{max} = 15$	43
C.4	Value Table: P4, $e_{min} = -7, e_{max} = 7$	44
C.5	Value Table: P5, $e_{min} = -3, e_{max} = 3$	45
C.6	Value Table: P6, $e_{min} = -1, e_{max} = 1$	46
C.7	Value Table: P7, $e_{min} = 0, e_{max} = 0$ (linear)	47

1 Contributors

Here are the active members of the Arithmetic Formats for Machine Learning working group.

Kiran Gunnam, *Chair*
Leonard Tsai, *Vice Chair*
Jeffrey Sarnoff, *Secretary*

Tom Thompson, *IEEE Standards Board Liaison*

Editors

Jeffrey Sarnoff, Andrew Fitzgibbon, Amos Omondi, Guy Lemieux

Paul Balanca	Michel Hack	Nathalie Revol
Luca Bertaccini	Simon Knowles	Jason Riedy
Javier Diaz Bruguera	Seokbum Ko	Ali Sazegari
David H. C. Chen	Carlo Luschi	Eric Schwarz
Marco Cococcioni	David Lutz	Oliver Sentieys
Mike Cowlshaw	Tue Ly	Michael Siu
Marius Cornea	Al Martin	Gil Tabak
Debjit Das Sarma	Paulis Micikevicius	Julio Villalba-Moreno
James Davenport	Mantas Mikaitis	Christop Wintersteiger
Jim Demmel	Aaftab Munshi	Kristopher Wong
Ken Dockser	Santosh Nagarakatte	Thomas Yeh
Massimiliano Fasi	Badreddine Noune	Chao Yu
Silviu Filip	Stuart Oberman	Aleksandr Zakharchenko
Jeff Gonion	Michael Overton	Hao Zhang
John Gustafson	Valentina Popescu	

Contacting The Working Group

To reach us, email FP-FOR-ML-STUDY-GROUP@listserv.ieee.org.

2 Introduction

This document represents ongoing discussions and current matters of consensus from IEEE Working Group P3109, “Standard for Arithmetic Formats for Machine Learning”. The Project Authorization Request (PAR) for P3109 defines the scope, need, and stakeholders as follows:

Scope of proposed standard: This standard defines a binary arithmetic and data format for machine learning-optimized domains. It also specifies the default handling of exceptions that occur in this arithmetic. This standard provides a consistent and flexible arithmetic framework optimized for Machine Learning Systems (MLS) in hardware and/or software implementations to minimize the work required to make MLS interoperable with each other, as well as other dependent systems. This standard is aligned with IEEE Std 754-2019 for Floating-Point Arithmetic.

Need for this Work: Machine Learning Systems have different arithmetic requirements from most other domains. Precisions tend to be lower, and accuracy is measured in dimensions other than just numerical (e.g. inference accuracy). Furthermore, machine learning systems are often integrated into mission-critical and safety-critical systems. With no standards specifically addressing these needs, Machine Learning Systems are built with inconsistent expectations and assumptions that hinder the compatibility and reuse of machine learning hardware, software, and training data.

Stakeholders for the Standard: System developers, vendors, and users of machine learning applications across many industries and interests including but not limited to computation, storage, medical, telecommunications, e-commerce, fleet management, automotive, robotics, and security.

2.1 Typographical conventions and notation

Bold text describes the decisions and specifications of this document.

Text that is not in boldface is background material, typically providing rationale and arguments that represent discussions of the working group leading to a decision and specification.

2.2 Scope

This document specifies compact floating-point interchange formats (binary formats) and associated operations.

Binary formats are parameterized by their width in bits and their precision—the number of bits in the true significand (which is one more than the number of bits in the trailing significand).

The 8-bit formats defined herein shall be referred to as “binary8” formats, and further qualified by precision yielding names “binary8pP” for values $1 \leq P \leq 7$.

For example, “binary8p3” is an 8-bit format with 1 sign bit, 5 exponent bits and 3 bits of precision (of which only two need to be explicitly represented).

This version of the interim report covers interchange formats and scalar operations including rounding with saturation modes and conversion between P3109 and IEEE Std 754. vector operations, block formats, stochastic rounding, and conversion between P3109 and BFloat16 are not discussed here. These topics may be included in a later version.

3 Formats

This section describes P3109 formats, and the set of values that such formats shall represent.

The universe of values in existing floating-point usage is the finite reals, the non-finite values positive and negative infinity (Inf, -Inf), the value negative zero (-0), and not-a-number values (NaN, NaN₁, ...), subscripted by *NaN payload*.

A K-bit binary floating-point format \mathcal{F} comprises: the *value set*: a subset $\mathcal{V}_{\mathcal{F}}$ of the universe of values; and an *encoding*: a mapping from integers $0 \dots 2^K - 1$ to $\mathcal{V}_{\mathcal{F}}$.

3.1 Format parameters and value sets

The finite floating-point numbers representable with a binary format are determined by two *format-defining* parameters:

- Storage width K, the total size of the format in bits
- Precision P, the number of bits in the significand including the implicit leading bit.

All other parameters, such as the exponent of the largest finite value e_{\max} , are derived from the format-defining parameters.

IEEE Std 754-2019 includes the radix B and the minimum exponent e_{\min} in a list of format-defining parameters, this document excludes both of them for these two reasons:

- This document covers binary (radix 2) formats only, so B is not a format parameter.
- The quantity e_{\min} is determined by P and e_{\max} ; it cannot be varied independently, so it cannot be a format-defining parameter.

P3109 formats shall be defined by the parameters of storage width K, and precision P, where $1 \leq P < K$

In IEEE Std 754, e_{\max} was set for all defined formats to be $2^{W-1} - 1$, where W is the exponent field width in bits. In this document, this convention is formalized: e_{\max} is a fixed function of P, written $e_{\max}(K, P)$, with the formula following IEEE Std 754, noting that $W = K - P$.

P3109 formats shall define $e_{\max}(K, P)$ to be $2^{K-P-1} - 1$

3.1.1 Computation of bias

The choice of e_{\max} for a given format then determines the exponent bias for that format. The bias is chosen so that the exponent of the largest finite value is e_{\max} .

For IEEE Std 754 formats, the largest finite value corresponds to an exponent field which has all but the zeroth bit set (e.g. 11110 for Binary16), because all of the values with all-bits-one exponents are occupied by non-finite values (Not-a-Numbers or Infinities). This derived parameter of a format is termed *all-special exponent*, and denoted by the symbol SE, where $SE = 1$ indicates that all-bits-one exponent fields are entirely occupied by special values. For P3109 values, the all-bits-one exponent contains only one special value ($\pm\text{Inf}$) and hence $SE = 0$ unless $P = 1$. Hence, for P3109, $SE = (P = 1)$.

Table 1: Parameters for binary formats

Symbol	Parameter Description	Derived Value	binary8pP, K = 8, P =							IEEE754-2019, K =		
			7	6	5	4	3	2	1	16	32	64
K	storage (bits)	K	8	8	8	8	8	8	8	16	32	64
P	precision (bits)	P	7	6	5	4	3	2	1	11	24	53
S	sign (bits)	1	1	1	1	1	1	1	1	1	1	1
W	exponent (bits)	K − P	1	2	3	4	5	6	7	5	8	11
T	trailing significand (bits)	P − 1	6	5	4	3	2	1	0	10	23	52
SE	all-special exponent	SE	0	0	0	0	0	0	1	1	1	1
emax	maximum exponent	$2^{W-1} - 1$	0	1	3	7	15	31	63	15	127	1023
emin	minimum exponent	SE − emax	0	-1	-3	-7	-15	-31	-62	-14	-126	-1022
bias	exponent bias	1 − emin	1	2	4	8	16	32	63	15	127	1023

Format-defining parameters in bold, derived parameters in normal font.
Adapted from Table 3.5 of IEEE Std 754-2019, and extended to include the binary8pP formats. Concepts are explained in detail in this section.

3.1.2 Computation of bias, SE = 1

With P3109 formats where P = 1, in common with IEEE Std 754, the biased exponent of the largest finite value is $2^W - 2$, from which bias should be defined so that

$$(2^W - 2) - \text{bias} = \text{emax}$$

Rearranging, we obtain the following

$$\text{bias} = (2^W - 2) - (2^{W-1} - 1) = 2 \cdot (2^{W-1} - 1) - (2^{W-1} - 1) = 2^{W-1} - 1 = \text{emax}$$

Hence bias = emax.

3.1.3 Computation of bias, SE = 0

For the binary8 formats in this document where P > 1, only one of the values that has exponents with all-bits-one is non-finite ($\pm\text{Inf}$), so the biased exponent of the largest finite value is $2^W - 1$. Hence the bias calculation becomes

$$\text{bias} = (2^W - 1) - (2^{W-1} - 1) = 2^{W-1} - 1 + 1 = \text{emax} + 1$$

As, by formula, emax is odd, the bias term in the P > 1 formats is typically even, yielding a more symmetrical range, where emin = −emax.

3.1.4 Observations on P3109 value sets

Table 1 shows some properties of the P3109 value sets for K = 8.

For P ≥ 1, the value sets are subsets of the IEEE Std 754 Binary32 value set. For P ≥ 3, the binary8pP value sets are also subsets of the IEEE Std 754 Binary16 value set.

The precision P is to be strictly less than K , hence formats where $P = K$ are not defined in this document. Rationale: strictly following Table 1 would yield $e_{\max} = -\frac{1}{2}$ which means all ordinary values are irrational. Rounding this computation upward yields $e_{\max} = 0$ and $\text{bias} = 1$, where all representations within the format are subnormal, with the consequence that the value sets and encodings for $P = K - 1$ and $P = K$ are identical except for a scaling factor of $1/2$.

The following notes from IEEE Std 754-2019 have not been found to apply to the definitions in this document:

- “For binary formats, the precision $[P]$ should be at least 3, as some numerical properties do not hold for lower precisions.”
- “Similarly, e_{\max} should be at least 2 to support the operations listed in 9.2. $[\sin/\cos/\exp/\log/\text{etc}]$ ”

3.2 Encodings and special values

Encodings are mappings from integers $0 \dots 2^K - 1$ to the value set, and are defined in §4.6.5 below.

Values are considered either “special” or “ordinary”. Encodings of the special values, shared by all P3109 formats, are shown in Table 2. P3109 formats have only a single NaN, located at the position that -0 would occupy in an analogous IEEE-754 encoding, providing an increased range. The ordinary values consist of the normal and subnormal values and zero.

Table 2: Special value encodings

Special Value	Symbol	Code point	Code point ($K = 8$)
Positive Infinity	Inf	$2^{K-1} - 1$	0x7F
Negative Infinity	-Inf	$2^K - 1$	0xFF
Not a Number	NaN	2^{K-1}	0x80

These mappings are shared by all P3109 formats with bit-width K .

3.3 Subnormals

P3109 value sets shall include subnormals.

IEEE Std 754 value sets include subnormals. A value with trailing significand field T and exponent field E is interpreted as $1.T \times 2^{E-\text{bias}}$ except when all bits of the exponent bitfield are 0, in which case the value is $0.T \times 2^{e_{\min}}$.

Subnormal numbers extend the dynamic range of floating-point values and induce equal quantization steps close to zero. They can be useful when training models, where it is common to represent near-zero values for gradients. Subnormals can also be useful to represent random values drawn from certain distributions. For example, model weights are initialized to small random values at training. Subnormals are uniformly spaced around zero, and values near zero are more probable values drawn from Gaussian-like distributions. Finally, formats with narrow exponent widths necessarily have a limited range; subnormals extend this range by a power of 2 for every bit in the trailing significand.

3.4 Not a number (NaN)

P3109 value sets shall include exactly one NaN, encoded as 2^{K-1} , which shall not signal.

Many other floating-point formats define several NaN values which are returned from operations with results outside the set of values, e.g., $\text{DIV}(0, 0)$, or $\text{ADD}(+\text{Inf}, -\text{Inf})$. Multiple NaN encodings are used in other formats to allow different exceptional conditions to be distinguished.

In the context of machine learning systems, uses of NaN include:

- Debugging of code running on accelerator hardware. In A.I. accelerators, exceptions may be difficult or expensive to convey back to user code, so it is common practice to allow NaN values to propagate through calculations to indicate that an error has occurred.
- Use as a sentinel value. In some datasets, for example, where individual element values may be missing or out of range, a sentinel may be used to record the position of these values. In many cases, this will require less memory than storing such information out-of-band, such as in a coordinate-list (COO) format array. In some cases, $\pm\text{Inf}$ can be used as a missing value, but given the restricted range of P3109 formats, it is likely that infinity shall be used as a separate indicator of rounding from values outside of the finite range.
- The use of multiple NaN payloads is known in statistical code (e.g. the R system has NaN and N/A), but it is not widely used. In the context of P3109, supporting multiple NaN would reduce the already limited encoding space (e.g., occupying all code points where the exponent field is all ones, thereby reducing dynamic range) and would likely add additional hardware complexity.

3.5 Zero

P3109 formats shall have exactly one zero, encoded as the integer 0. This zero value is non-negative.

The inclusion of negative zero (-0) would incur the cost of an additional code point. Given the decision to encode only a single NaN, placing that NaN at the negative zero code point enables the strictly positive and strictly negative number ranges to be symmetric.

A key rationale for including -0 in IEEE Std 754 was the consistent implementation of branch cuts in the atan2 function [1, 2]. Although the atan function is common in deep learning, it is generally used as an activation function, rather than a trigonometric operation, and the atan2 function is rare, if not unknown, in deep learning applications. Hence, it is not expected that this standard shall define either atan or atan2 .

A secondary reason for providing -0 is the hardware simplification offered by its presence in the implementation of sign/magnitude arithmetic. However, the existence of in-market implementations is evidence that the small hardware simplification has not been sufficient to balance the loss of one code point.

It might be considered that the use of integer comparisons in sorting would argue against placing NaN at the negative zero code point. For example, the JAX machine learning framework is known to sort using integer comparison [3]. However, such sorting still requires $O(n)$ preprocessing and postprocessing steps to enable the use of two's-complement integer comparison, and already has special treatment of NaN and -0 , so eliminating -0 and placing NaN in the -0 position imposes negligible additional burden. Although sorting using comparison operations is undefined in the presence of NaNs, existing practice is to sort NaN using `totalOrder`.

3.6 Infinities

P3109 formats shall include positive and negative infinities, encoded as $2^{K-1} - 1$ and $2^K - 1$, respectively.

This decision causes a reduction in dynamic range (252 values rather than 254 for binary8, for example), while offering improved numerical robustness in important machine learning use cases.

Examples of such usage are:

- Mask values, for example, in Transformer models in machine learning. See §A.1 for more detailed discussion on this usage.
- Representation of overflow, for example, to adjust dynamic loss scaling factors [4]. Existing implementations offer several behaviors on overflow: overflow to infinity, saturation to MaxFloat, and overflow to NaN. The existence of a code point for infinity allows any of these options to be implemented in a given instantiation, while removing the code point removes the possibility of implementing the first.

3.7 Extremal values for $K = 8$

Table 3: Extremal values

Format	minSubnormal	maxSubnormal	minNormal	maxNormal
binary8p1	N/A	N/A	1×2^{-62}	1×2^{63}
binary8p2	1×2^{-32}	1×2^{-32}	1×2^{-31}	1×2^{31}
binary8p3	1×2^{-17}	$3/2 \times 2^{-16}$	1×2^{-15}	$3/2 \times 2^{15}$
binary8p4	1×2^{-10}	$7/4 \times 2^{-8}$	1×2^{-7}	$7/4 \times 2^7$
binary8p5	1×2^{-7}	$15/8 \times 2^{-4}$	1×2^{-3}	$15/8 \times 2^3$
binary8p6	1×2^{-6}	$31/16 \times 2^{-2}$	1×2^{-1}	$31/16 \times 2^1$
binary8p7	1×2^{-6}	$63/32 \times 2^{-1}$	1×2^0	$63/32 \times 2^0$

It is practical to list extremal finite values defined by the binary8 formats. Following IEEE Std 754-2019 naming patterns, we adopt: $\text{maxNormal}(\tau)$, $\text{minNormal}(\tau)$, $\text{minSubnormal}(\tau)$ where τ is a binary8 format. For example: $\text{maxNormal}(\text{binary8p4}) = 7/4 \times 2^7$ and $\text{minNormal}(\text{binary8p5}) = 1 \times 2^{-3}$.

Table 3 shows the extremal values for $1 \leq P \leq 7$. For reference, section C provides complete tables of 8-bit values.

3.8 Equivalences

The $P = K - 1$ formats are equivalent to a signed-magnitude fixed-point integer, except there are code points reserved for NaN and $\pm\text{Inf}$.

The $P = 1$ formats are “purely exponential” formats, where all values are of the form 2^n , excepting 0, NaN, $\pm\text{Inf}$.

4 Operations

This section defines the *behavior* (with no constraint as to the implementation) of operations which P3109 systems may provide. Such operations include:

- conversion between P3109 formats and between P3109 formats and IEEE-754 formats.
- comparison, classification, and informational operations
- addition, multiplication and fused multiply add

In the definitions of the above operations, certain mathematical *functions* are also defined. We emphasize that the definitions in this document are specifications of behavior, not implementation. All arithmetic represented as operating on extended real values is to be interpreted in the usual mathematical sense.

An implementation of any of the operations defined herein is conforming if it computes the same output as the defined operation, for all possible inputs. This may be attested by any proof method, including direct computation.

As defined above, a K-bit P3109 floating-point value (referred to in general as a *P3109 value*) is encoded by an integer in the range 0 to $2^K - 1$. Operation specifications operate on such values using integer arithmetic (e.g., divide and modulo) operations. Implementations may perform these operations in any equivalent manner, for example bit shifting and masking.

For arithmetic operations, an implementation may provide an α -*approximate* implementation. Such an implementation will compute values whose maximum difference from the defined outputs, over all inputs producing finite outputs, does not exceed α units in the last place. This may be attested by any proof method, including direct computation.

4.1 Operation variants and superset specifications

Each operation definition in this document is *parameterized*. Example parameters might include input and output formats, rounding modes and saturation modes. In addition operand types may be specified as if to infinite precision (e.g. s in the template definition below is defined as an integer). The set of variants so defined for a given operation is called a *superset specification*.

No implementation is required to provide all variants of a given operation, but rather should clearly define the implemented variants.

It is required that the use of an operation name defined in this document, in the context of declaration of P3109 compliance, specify precisely the variants implemented, and that the name be used only for variants which exactly follow these specifications.

For example, an implementation might declare as follows. (See later sections for definitions of terms in this example.)

“This implementation conforms to P3109, providing the following operation variants:

ConvertToP3109, ϕ in {Binary16, Binary32}, f in {binary8p3, binary8p4, binary8p5}

AddScaled, with $s_x = 0$, $s_y \in \{-128 \dots 127\}$, $\{f_x, f_y, f_z\} \subset \{\text{binary8p3, binary8p4, binary8p5}\}$

MultiplyScaled, with $s \in \{-32 \dots 32\}$

In all cases, provided projection modes π obey $\text{Rnd}_\pi \in \{\text{NearestTiesToEven, NearestTiesToAway}\}$ and $\text{Sat}_\pi \in \{\text{SatFinite, Ovflnf}\}$ ”

This example is truncated—realistic compliance declarations might run to many pages. Such declarations may be compressed by defining common *profiles*, outside of the scope of this current document.

4.2 Notation and definitions

4.2.1 Mathematical notations

We denote by $\mathbb{1}[b]$ the value 1 if $b = \text{True}$, or 0 if $b = \text{False}$.

$\text{IsEven}(I)$ is true if integer I is even, false otherwise.

Integer division is denoted by $x \div y$.

Comments are introduced with an em-dash, and are right-justified

—An example comment

4.2.2 The set of extended reals

The set of *extended reals* is the set $\mathbb{R} \cup \{-\infty, \infty\}$ of real numbers augmented with positive and negative infinity, also known as the “affinely extended real numbers”. In common with existing mathematical treatments, there is no negative zero in this set. In the extended reals, certain operations, such as $\infty - \infty$, are *undefined*. The definitions in this specification are constructed such that no operation on extended real quantities produces an undefined result.

Operation specifications will, in general, convert floating-point operands to extended real values in order to define their behaviors.

4.3 Projection specifications

Operations on P3109 values are defined via conversion to extended real values, on which the mathematical operation is performed, before conversion back to the appropriate P3109 range. In general, operation results will not be exact P3109 values, and hence will be *projected* into the P3109 range via rounding and overflow handling. A *projection specification* is a pair (rounding mode, saturation mode). For a given projection specification π , these are written Rnd_π , Sat_π .

The defined rounding modes are as follows. The precise specifications of these modes are in the function `RoundToPrecision` (§4.6.2):

<code>NearestTiesToEven</code>	Round to nearest, ties to even
<code>NearestTiesToAway</code>	Round to nearest, ties away from zero
<code>TowardPositive</code>	Round toward positive
<code>TowardNegative</code>	Round toward negative
<code>TowardZero</code>	Round toward zero

Values are first rounded to the target precision, with exponent unbounded above. Those which are then outside the maximum value in the target format are then treated according to the saturation mode.

The defined saturation modes are as follows. The precise specifications of these modes are in the function `Saturate` (§4.6.4):

<code>SatMax</code>	All return values are clamped to the representable finite range.
<code>SatFinite</code>	Finite out-of-range values are clamped to the representable finite range, infinities are preserved.
<code>Ovflnf</code>	Out-of-range values are replaced with extreme value, positive or negative infinity as indicated by the rounding mode.

4.4 Non-requirement for operator side effects

The operator definitions herein describe no side effects, such as the setting of flags, or the triggering of interrupts, and hence do not return values other than the defined return value. Typically NaN is returned when input values are out of range (e.g. $\log(-1.0)$). When saturation is specified, there is no direct mechanism to distinguish overflowed values from values which round to the format's maximum value. An implementation which has side effects will still conform to this specification providing its return value on all inputs matches the definitions herein.

4.5 Operator definition template

Operations are defined according to the following template:

Signature

$\text{Operator}_f(x, s) \rightarrow Z$

—Naming the operator, its parameters, operands and result.

Parameters

f : format, precision P_f

—Parameters specify a family of related operations

Operands

x : P3109 value, in format f

s : Integer scale

Output

Z : extended real value or NaN

—Result value and type

Behavior

$\text{Operator}(\text{NaN}, 0) \rightarrow \text{NaN}$

$\text{Operator}(x \in \{-\text{Inf}, \text{Inf}\}, 0) \rightarrow x$

$\text{Operator}(*, 0) \rightarrow 0$

$\text{Operator}(x, s < 0) \rightarrow -\text{Operator}_f(x, -s)$

—An ordered sequence of pattern-matching declarations.

—Exact pattern: only the provided operands match.

—Match for all x values in the given set.

—The $*$ symbol matches any value.

—Operators in pattern RHS have explicit parameters.

Notes

Notes on the operation.

In a sequence of pattern-matching declarations, the first matching pattern in the order presented in this document defines the behavior for a given operand sequence.

In pattern matching, certain shorthand notations are used, as follows. Consider an operation which takes a P3109 value x in format f_x , and returns a P3109 value in format f_y . One matching pattern might be $\text{Operator}(\text{NaN}) \rightarrow \text{NaN}$, where the input and output NaN values are in different formats, although this is not explicitly marked. A more explicit presentation, such as $\text{Operator}(\text{NaN}_{f_x}) \rightarrow \text{NaN}_{f_y}$ was considered to increase difficulty of comprehension without a compensatory reduction in ambiguity.

Similarly, parameter subscripts are elided on the left-hand-side of patterns where they are common to all, but are written explicitly on the right.

Ancillary information, typically parameters, may include information such as an operand's format or rounding behaviors. An implementation may choose to provide that information using any appropriate mechanism. For example, available mechanisms in a hardware implementation might include passing the information in a hardware register, or as additional bits in an operation's opcode.

4.6 Functions

This section defines mathematical functions which are referenced in the later definitions of operations, but which themselves need not (and, in cases where the inputs or outputs are real, cannot) be provided in a conforming implementation.

Behavior that is specified through cases should be read from top to bottom.

4.6.1 Decode

Decode P3109 value to extended real or NaN.

Signature

$$\text{Decode}_f(x) \rightarrow X$$

Parameters

f : format, width K_f , precision P_f , exponent bias b_f

Operands

x : P3109 value, in format f

Output

X : extended real value or NaN

Behavior

$$\text{Decode}(2^{K_f-1}) \rightarrow \text{NaN}$$

$$\text{Decode}(2^{K_f-1} - 1) \rightarrow \infty$$

$$\text{Decode}(2^{K_f-1} < x < 2^{K_f}) \rightarrow -\text{Decode}_f(x - 2^{K_f-1})$$

$$\text{Decode}(0 \leq x < 2^{K_f-1} - 1) \rightarrow X$$

where

$$X = \begin{cases} (0 + T \times 2^{1-P_f}) \times 2^{E+1} & \text{if } E = -b_f \\ (1 + T \times 2^{1-P_f}) \times 2^E & \text{Otherwise} \end{cases} \quad \begin{array}{l} \text{---Subnormal} \\ \text{---Normal} \end{array}$$

$$T = x \bmod 2^{P_f-1}$$

$$E = x \div 2^{P_f-1} - b_f$$

4.6.2 Project

Project extended real value to P3109 format f , applying specified rounding and saturation.

Signature

$$\text{Project}_{f,\pi}(X) \rightarrow x$$

Parameters

f : target format, precision P_f , exponent bias b_f , maximum finite value M_f

π : projection specification: rounding mode Rnd_π , saturation Sat_π

Operands

X : extended real value

Output

x : P3109 value, format f

Behavior

$$\text{Project}_{f,\pi}(X) \rightarrow x$$

where

$R = \text{RoundToPrecision}_{P_f, b_f, \text{Rnd}_\pi}(X)$ —Round to precision P_f with exponent unbounded from above

$S = \text{Saturate}_{M_f}(\text{Sat}_\pi, \text{Rnd}_\pi, R)$

$x = \text{Encode}_f(S)$

4.6.3 RoundToPrecision

Convert extended real value to extended real value representable with a given precision. The exponent is bounded below by $2 - P - b$ and unbounded above.

Signature

$$\text{RoundToPrecision}_{P,b,\text{Rnd}}(X) \rightarrow Z$$

Parameters

P : integer precision

b : exponent bias

Rnd : rounding mode

Operands

X : extended real value

Output

Z : extended real value, of the form $N \times 2^E$, where $N \in \mathbb{Z}$, $0 \leq |N| < 2^P$, and $2 - P - b \leq E$.

Behavior

$$\text{RoundToPrecision}(X \in \{0, -\infty, \infty\}) \rightarrow X$$

$$\text{RoundToPrecision}(X) \rightarrow Z$$

where

$$E = \max(\lfloor \log_2(|X|) \rfloor, 1 - b) - P + 1 \quad \text{---Subnormals handled by } \max(\cdot, 1 - b)$$

$$S = |X| \times 2^{-E} \quad \text{---Real-valued significand, to be rounded to integer}$$

$$\Delta = S - \lfloor S \rfloor$$

$$\text{CodelsEven} = \begin{cases} \text{IsEven}(\lfloor S \rfloor) & \text{if } P > 1 \\ (\lfloor S \rfloor = 0) \vee \text{IsEven}(E + b) & \text{Otherwise} \end{cases}$$

$$I = \lfloor S \rfloor + \mathbb{1}[\text{RoundAway}(\text{Rnd})]$$

$$Z = \text{sign}(X) \times I \times 2^E$$

and

$$\text{RoundAway}(\text{NearestTiesToEven}) = \Delta > 0.5 \vee (\Delta = 0.5 \wedge \neg \text{CodelsEven})$$

$$\text{RoundAway}(\text{NearestTiesToAway}) = \Delta \geq 0.5$$

$$\text{RoundAway}(\text{TowardPositive}) = \Delta > 0 \wedge X > 0$$

$$\text{RoundAway}(\text{TowardNegative}) = \Delta > 0 \wedge X < 0$$

$$\text{RoundAway}(\text{TowardZero}) = \text{False}$$

Notes

The $P = 1$ logic for CodelsEven uses the observation that $I = 0 \implies Z = 0$, hence the code is 0, hence even.

Only precision and bias are needed for this definition, not format width.

Note that I may be set to 2^P , which might appear to preclude its representation in $P - 1$ bits of explicit significand, but the computed real value is then $(1 + 0) \times 2^{E+1+P}$, representable as the first number in the next binade.

4.6.4 Saturate

Saturate extended real to $\pm\infty$, or to maximum value, according to projection specification parameters Sat, Rnd.

Signature

$\text{Saturate}_M(\text{Sat}, \text{Rnd}, X) \rightarrow Z$

Parameters

M : real maximum value

Operands

Sat : saturation mode

Rnd : rounding mode

X : extended real value

Output

Z : extended real value

Behavior

$\text{Saturate}(*, *, X \in [-M, M]) \rightarrow X$

$\text{Saturate}(\text{SatMax}, *, X \notin [-M, M]) \rightarrow \text{sign}(X) \times M$

$\text{Saturate}(\text{SatFinite}, *, X \in \{\pm\infty\}) \rightarrow X$

$\text{Saturate}(\text{SatFinite}, *, |X| \in [M, \infty)) \rightarrow \text{sign}(X) \times M$

$\text{Saturate}(\text{Ovflnf}, *, X \in \{\pm\infty\}) \rightarrow X$

$\text{Saturate}(\text{Ovflnf}, \text{TowardZero}, |X| \in [M, \infty)) \rightarrow \text{sign}(X) \times M$

$\text{Saturate}(\text{Ovflnf}, \text{TowardPositive}, X \in (-\infty, -M)) \rightarrow -M$

$\text{Saturate}(\text{Ovflnf}, \text{TowardNegative}, X \in (M, \infty)) \rightarrow M$

$\text{Saturate}(\text{Ovflnf}, *, X) \rightarrow \text{sign}(X) \times \infty$

Note

In Saturate, all values above M are sent to either M or ∞ . Nevertheless, Project (§4.6.2) has the property that values below $M + \frac{1}{2} \text{ulp}(M)$ will round to M because RoundToPrecision precedes saturation in the definition of Project.

4.6.5 Encode

Encode extended real value or NaN to P3109 format f . Encode must be applied only to a value which is in the value set of format f . Such a value may be produced by, for example RoundToPrecision and Saturate, or the input may be known to be in the value set, for example, in the negation of a value already in the set.

Signature

$$\text{Encode}_f(X) \rightarrow z$$

Parameters

f : target format, width K_f , precision P_f , exponent bias b_f

Operands

X : extended real value or NaN, in the value set of format f

Output

z : P3109 value, format f

Behavior

$$\text{Encode}(\text{NaN}) \rightarrow 2^{K_f-1}$$

$$\text{Encode}(\infty) \rightarrow 2^{K_f-1} - 1$$

$$\text{Encode}(X < 0) \rightarrow \text{Encode}_f(-X) + 2^{K_f-1}$$

$$\text{Encode}(0) \rightarrow 0$$

$$\text{Encode}(X > 0) \rightarrow z$$

where

$$z = \begin{cases} T & \text{if } S < 2^{P_f-1} \\ T + (E + b_f) \times 2^{P_f-1} & \text{Otherwise} \end{cases} \quad \text{—Subnormals}$$

and

$$E = \max(\lfloor \log_2(X) \rfloor, 1 - b_f)$$

$$S = X \times 2^{-E} \times 2^{P_f-1} \quad \text{—}S \text{ is the significand}$$

$$T = S \bmod 2^{P_f-1}$$

Note

Because of the precondition that X is in the value set of format f , it follows that $S \in \mathbb{N}$.

4.7 Conversion operations

4.7.1 Specification of IEEE Std 754 formats

In specifying IEEE Std 754 formats, the symbol ϕ is used, from which format parameters are extracted as needed as follows:

M_ϕ : maximum finite value (e.g., 65504.0 for Binary16)

P_ϕ : precision (e.g., 11 for Binary16)

B_ϕ : exponent bias (e.g., 15 for Binary16)

We will make use of the function `Encode754`, defined as follows:

Signature

$\text{Encode754}_\phi(X) \rightarrow z$

Parameters

ϕ : target IEEE Std 754 format

Operands

X : NaN or extended real value, in the value set of format ϕ

Output

z : IEEE Std 754 value, format ϕ

Behavior

$\text{Encode754}(\text{NaN}) \rightarrow$ Any quiet IEEE Std 754 NaN

$\text{Encode754}(X) \rightarrow$ The code in ϕ that decodes to X

Notes

In this document, the `Encode754` function is only called with arguments in the value set of format ϕ , therefore encoding is unambiguous and independent of rounding mode. Conversely, it is an error in this document if this condition is not guaranteed.

An implementation may return any quiet NaN. It is recommended that the quiet NaN with zero payload is returned.

4.7.2 Conversion from IEEE Std 754 formats to P3109

Convert a value in an IEEE Std 754 format to the corresponding value in a given P3109 format, considering rounding and saturation.

Signature

$\text{ConvertToP3109}_{\phi, f, \pi}(X) \rightarrow z$

Parameters

ϕ : source IEEE Std 754 format

f : target format

π : projection specification

Operands

X : IEEE-754 value, format ϕ

Output

z : P3109 value, format f

Behavior

$\text{ConvertToP3109}(\text{Any IEEE Std 754 NaN}) \rightarrow \text{NaN}$

$\text{ConvertToP3109}(X) \rightarrow \text{Project}_{f, \pi}(X')$ **where** $X' = \text{AsExtendedReal}_{\phi}(X)$

Note

$\text{AsExtendedReal}_{\phi}$ is a function which converts a non-NaN encoded IEEE Std 754 value to a value in the extended reals.
 $\text{AsExtendedReal}(-0) \rightarrow 0$.

4.7.3 Conversion from P3109 to IEEE Std 754

Convert a P3109 value to a value in an IEEE Std 754 format. We note that for IEEE Std 754 binary $\{K\}$ formats for $K > 16$, all P3109 values are values in the target format, hence the conversion of non-NaN values is unambiguous. For Binary16 outputs, some P3109 values will be out of the target range; it is necessary to be precise about rounding and saturation.

Signature

$\text{ConvertToIEEE754}_{f,\pi,\phi}(x) \rightarrow X$

Parameters

f : input format

π : projection specification

ϕ : IEEE Std 754 result format, precision P_ϕ , bias B_ϕ , maximum value M_ϕ

Operands

x : P3109 value, format f

Output

X : IEEE-754 value, format ϕ

Behavior

$\text{ConvertToIEEE754}(\text{NaN}) \rightarrow \text{Encode754}_\phi(\text{NaN})$

$\text{ConvertToIEEE754}(x) \rightarrow \text{Encode754}_\phi(X)$

where

$Y = \text{Decode}_f(x)$

$R = \text{RoundToPrecision}_{P_\phi, B_\phi, \text{Rnd}_\pi}(Y)$

$X = \text{Saturate}_{M_\phi}(\text{Sat}_\pi, \text{Rnd}_\pi, R)$

4.7.4 Conversion from P3109 to P3109

Convert a P3109 value to another P3109 format.

Signature

$\text{ConvertP3109ToP3109}_{f_x, f_z, \pi}(x) \rightarrow z$

Parameters

f_x : input format

f_z : target format

π : projection specification

Operands

x : P3109 value, format f_x

Output

z : P3109 value, format f_z

Behavior

$\text{ConvertP3109ToP3109}(\text{NaN}) \rightarrow \text{NaN}$

$\text{ConvertP3109ToP3109}(x) \rightarrow \text{Project}_{f_z, \pi}(X)$ where $X = \text{Decode}_{f_x}(x)$

4.8 Arithmetic operations

Arithmetic operations which take one or more P3109 values as arguments, and return one or more P3109 values are defined in this section. Arithmetic operations which mix IEEE Std 754 and P3109 values are described in §4.9.

4.8.1 Unary sign operations

Operations which take P3109 values as input, and whose outputs are guaranteed to be exact values in the same P3109 format. These operations do not permit the output format to be different from the input format.

Signature

Operation_{*f*}(*x*) → *z*

Parameters

f : input and output format

Operands

x : P3109 value, format *f*

Output

z : P3109 value, format *f*

Behavior

$$\text{Abs}(x) \rightarrow \begin{cases} \text{NaN} & \text{isNaN}(x) \\ \text{Encode}_f(|X|) & \text{Otherwise} \end{cases}$$

where $X = \text{Decode}_f(x)$

$$\text{Negate}(x) \rightarrow \begin{cases} \text{NaN} & \text{isNaN}(x) \\ \text{Encode}_f(-X) & \text{Otherwise} \end{cases}$$

where $X = \text{Decode}_f(x)$

Note

The extended reals do not have negative zero, hence $-0 = 0$ and $\text{Abs}(0) = 0$.

4.8.2 Binary sign operations

Signature

$\text{Operation}_f(x, y) \rightarrow z$

Parameters

f : input and output format

Operands

x : P3109 value, format f

y : P3109 value, format f

Output

z : P3109 value, format f

Behavior

$$\text{CopySign}(x, y) \rightarrow \begin{cases} \text{NaN} & \text{isNaN}(x) \vee \text{isNaN}(y) \\ \text{Abs}(x) & Y \geq 0 \\ \text{Negate}(\text{Abs}(x)) & Y < 0 \end{cases}$$

where $Y = \text{Decode}_f(y)$

4.8.3 Binary arithmetic operations

These operations take two P3109 values as input, and return a P3109 value. The definitions allow for all input and output formats to differ, a given implementation may supply any subset of the defined operations.

Signature

Operation $_{f_x, f_y, f_z, \pi}(x, y) \rightarrow z$

Parameters

f_x : format of x

f_y : format of y

f_z : format of z

π : projection specification

Operands

x : P3109 value, format f_x

y : P3109 value, format f_y

Output

z : P3109 value, format f_z

Behavior

Add(*, NaN) \rightarrow NaN

Add(NaN, *) \rightarrow NaN

Add(-Inf, Inf) \rightarrow NaN

Add(Inf, -Inf) \rightarrow NaN

Add(x, y) \rightarrow Project $_{f_z, \pi}(X + Y)$ where $X = \text{Decode}_{f_x}(x), Y = \text{Decode}_{f_y}(y)$

Subtract(*, NaN) \rightarrow NaN

Subtract(NaN, *) \rightarrow NaN

Subtract(Inf, Inf) \rightarrow NaN

Subtract(-Inf, -Inf) \rightarrow NaN

Subtract(x, y) \rightarrow Project $_{f_z, \pi}(X - Y)$ where $X = \text{Decode}_{f_x}(x), Y = \text{Decode}_{f_y}(y)$

Multiply(*, NaN) \rightarrow NaN

Multiply(NaN, *) \rightarrow NaN

Multiply(0, \pm Inf) \rightarrow NaN

Multiply(\pm Inf, 0) \rightarrow NaN

Multiply(x, y) \rightarrow Project $_{f_z, \pi}(X \times Y)$ where $X = \text{Decode}_{f_x}(x), Y = \text{Decode}_{f_y}(y)$

Divide(*, NaN) \rightarrow NaN

Divide(NaN, *) \rightarrow NaN

Divide(*, 0) \rightarrow NaN

Divide(x, y) \rightarrow Project $_{f_z, \pi}(X/Y)$ where $X = \text{Decode}_{f_x}(x), Y = \text{Decode}_{f_y}(y)$

Notes

Divide(x, \pm Inf) where x is finite yields 0.

Divide($x, 0$) yields NaN. With x finite, it would be inconsistent to return Inf: $1/(1/-\text{Inf}) \rightarrow \text{Inf}$.

4.8.4 Unary mathematical operations

Signature

Operation _{f_x, f_z, π} (x) $\rightarrow z$

Parameters

f_x : input format

f_z : output format

π : projection specification

Operands

x : P3109 value, format f_x

Output

z : P3109 value, format f_z

Behavior

$$\text{Sqrt}(x) \rightarrow \begin{cases} \text{NaN} & \text{isNaN}(x) \\ \text{NaN} & X < 0 \\ \text{Project}_{f_z, \pi}(\sqrt{X}) & X \geq 0 \end{cases}$$

where $X = \text{Decode}_{f_x}(x)$

$$\text{Exp}(x) \rightarrow \begin{cases} \text{NaN} & \text{isNaN}(x) \\ \text{Project}_{f_z, \pi}(e^X) & \text{Otherwise} \end{cases}$$

where $X = \text{Decode}_{f_x}(x)$

$$\text{Exp2}(x) \rightarrow \begin{cases} \text{NaN} & \text{isNaN}(x) \\ \text{Project}_{f_z, \pi}(2^X) & \text{Otherwise} \end{cases}$$

where $X = \text{Decode}_{f_x}(x)$

$$\text{Log}(x) \rightarrow \begin{cases} \text{NaN} & \text{isNaN}(x) \\ \text{NaN} & X < 0 \\ \text{Project}_{f_z, \pi}(\log(X)) & X \geq 0 \end{cases}$$

where $X = \text{Decode}_{f_x}(x)$

$$\text{Log2}(x) \rightarrow \begin{cases} \text{NaN} & \text{isNaN}(x) \\ \text{NaN} & X < 0 \\ \text{Project}_{f_z, \pi}(\log_2(X)) & X \geq 0 \end{cases}$$

where $X = \text{Decode}_{f_x}(x)$

4.8.5 Scaled addition

Compute $X \times 2^{s_x} + Y \times 2^{s_y}$, and return a P3109 value. Scaling is applied in the extended reals, before projection to the target format.

Signature

$\text{AddScaled}_{f_x, f_y, f_z, \pi}(x, s_x, y, s_y) \rightarrow z$

Parameters

f_x : format of x
 f_y : format of y
 f_z : format of z
 π : projection specification

Operands

x : P3109 value, format f_x
 s_x : integer log-scale factor for x
 y : P3109 value, format f_y
 s_y : integer log-scale factor for y

Output

z : P3109 value, format f_z

Behavior

$\text{AddScaled}(\text{NaN}, *, *, *) \rightarrow \text{NaN}$
 $\text{AddScaled}(*, *, \text{NaN}, *) \rightarrow \text{NaN}$
 $\text{AddScaled}(-\text{Inf}, *, \text{Inf}, *) \rightarrow \text{NaN}$
 $\text{AddScaled}(\text{Inf}, *, -\text{Inf}, *) \rightarrow \text{NaN}$
 $\text{AddScaled}(x, s_x, y, s_y) \rightarrow \text{Project}_{f_z, \pi}(Z)$

where

$$\begin{aligned} Z &= X \times 2^{s_x} + Y \times 2^{s_y} \\ X &= \text{Decode}_{f_x}(x) \\ Y &= \text{Decode}_{f_y}(y) \end{aligned}$$

Note

The valid range of the log-scale factors s_x and s_y is not specified. Implementers shall declare the provision of `AddScaled` with specified constraints on legal values for these operands.

Example declarations might include

“AddScaled, with log-scale factors in $\{-32 \dots 31\}$ ”
“AddScaled, with $s_x = 0$, $s_y \in \{-128 \dots 127\}$ ”
“AddScaled, with $s_x \in \{-128 \dots 127\}$, $s_y = s_x$ ”

See §4.1 for further discussion.

4.8.6 Scaled multiplication

Compute $X \times Y \times 2^s$, and return a P3109 value. Scaling is applied in the extended reals, before projection to the target format.

Signature

$\text{MultiplyScaled}_{f_x, f_y, f_z, \pi}(x, y, s) \rightarrow z$

Parameters

f_x : format of x

f_y : format of y

f_z : format of z

π : projection specification

Operands

x : P3109 value, format f_x

y : P3109 value, format f_y

s : integer log-scale factor

Output

z : P3109 value, format f_z

Behavior

$\text{MultiplyScaled}(*, \text{NaN}, *) \rightarrow \text{NaN}$

$\text{MultiplyScaled}(\text{NaN}, *, *) \rightarrow \text{NaN}$

$\text{MultiplyScaled}(0, \pm\text{Inf}, *) \rightarrow \text{NaN}$

$\text{MultiplyScaled}(\pm\text{Inf}, 0, *) \rightarrow \text{NaN}$

$\text{MultiplyScaled}(x, y, s) \rightarrow \text{Project}_{f_z, \pi}(Z)$

where

$$Z = X \times Y \times 2^s$$

$$X = \text{Decode}_{f_x}(x)$$

$$Y = \text{Decode}_{f_y}(y)$$

Note

The valid range of the log-scale factor s is not specified. Implementations shall declare the provision of `MultiplyScaled` with specified constraints on legal values for these operands.

An example such declaration might be

“`MultiplyScaled`, with log-scale factor in $\{-32 \dots 31\}$ ”

See §4.1 for further discussion.

4.9 Mixed IEEE Std 754 and P3109 operations

This section describes operations which take a mix of IEEE Std 754 and P3109 operands and return IEEE Std 754 values. An implementation is free to define additional fused operations in which P3109 operands are upconverted to IEEE Std 754 before operating.

This section addresses operations that cannot be expressed as a fused sequence of operations. For instance, it is not possible to upconvert all values from binary8p1 to Binary16. Moreover, the operations described here may include saturation and rounding modes not available in definitions based solely on upconversion.

Like the operations defined previously, these are superset specifications. An implementation might provide any subset of the parameterized set of operations, accompanied by a statement specifying which subset is implemented.

These definitions make use of the functions `AsExtendedReal` and `Encode754` defined above.

4.9.1 Scaled Fused Multiply Add

Compute $Z = A \times 2^{s_a} + X \times Y \times 2^s$, with A and Z in the same format. Scaling is applied in the extended reals, before projection to the target format.

Signature

$\text{ScaledFMA}_{\phi, f_x, f_y, \pi}(a, s_a, x, y, s) \rightarrow z$

Parameters

ϕ : IEEE Std 754 format of a and target, precision P_ϕ , bias B_ϕ , maximum value M_ϕ
 f_x : format of x
 f_y : format of y
 π : projection specification

Operands

a : IEEE Std 754 value, format ϕ
 s_a : integer log-scale factor
 x : P3109 value, format f_x
 y : P3109 value, format f_y
 s : integer log-scale factor

—or NaN?

Output

z : IEEE Std 754 value, format ϕ

Behavior

$\text{ScaledFMA}(\text{NaN}, *, *, *, *) \rightarrow \text{NaN}$
 $\text{ScaledFMA}(*, *, \text{NaN}, *, *) \rightarrow \text{NaN}$
 $\text{ScaledFMA}(*, *, *, \text{NaN}, *) \rightarrow \text{NaN}$
 $\text{ScaledFMA}(a, s_a, x, y, s) \rightarrow z$

where

$A = \text{AsExtendedReal}_\phi(a)$
 $X = \text{Decode}_{f_x}(x)$
 $Y = \text{Decode}_{f_y}(y)$
 $Z = A \times 2^{s_a} + X \times Y \times 2^s$
 $Z' = \text{RoundToPrecision}_{P_\phi, B_\phi, \text{Rnd}_\pi}(Z)$
 $Z'' = \text{Saturate}_{M_\phi}(\text{Sat}_\pi, \text{Rnd}_\pi, Z')$
 $z = \text{Encode754}_\phi(Z'')$

— Z'' guaranteed representable in format ϕ

Note

The valid range of the log-scale factors are not specified. Implementations shall declare the provision of ScaledFMA with specified constraints on legal values for these operands.

An example such declaration might be:

“ScaledFMA, with log-scale factors in $\{-32 \dots 31\}$, $f_x = f_y$ in $\{\text{binary8p3}, \text{binary8p4}\}$, ϕ in $\{\text{Binary16}, \text{Binary32}\}$ ”

4.10 Comparisons, predicates, and classification

4.10.1 Minimum and Maximum

These operations take two P3109 values as input, and return a P3109 value. Operands and return value are in the same format.

Signature

$$\text{Operation}_f(x, y) \rightarrow z$$

Parameters

$$f : \text{format}$$

Operands

$$x : \text{P3109 value, format } f$$
$$y : \text{P3109 value, format } f$$

Output

$$z : \text{P3109 value, format } f$$

Behavior

$$\text{Minimum}(*, \text{NaN}) \rightarrow \text{NaN}$$
$$\text{Minimum}(\text{NaN}, *) \rightarrow \text{NaN}$$
$$\text{Minimum}(x, y) \rightarrow \text{Encode}_f(\min(X, Y)) \quad \textbf{where} \quad X = \text{Decode}_f(x), Y = \text{Decode}_f(y)$$
$$\text{Maximum}(*, \text{NaN}) \rightarrow \text{NaN}$$
$$\text{Maximum}(\text{NaN}, *) \rightarrow \text{NaN}$$
$$\text{Maximum}(x, y) \rightarrow \text{Encode}_f(\max(X, Y)) \quad \textbf{where} \quad X = \text{Decode}_f(x), Y = \text{Decode}_f(y)$$

Note

The operations **minimumNumber**, **maximumNumber**, defined by IEEE Std 754-2019 are not defined in this document.

4.10.2 Comparisons

Comparison operators take two P3109 operands and return a boolean. Any NaN operand yields the result False.

The below specifications are expressed with the following substitutions for compareOp and CompareExpr :

compareOp	CompareExpr
compareLess	$X < Y$
compareLessEqual	$X \leq Y$
compareEqual	$X = Y$
$\text{compareGreaterEqual}$	$X > Y$
compareGreater	$X \geq Y$

Signature

$\text{compareOp}_{f_x, f_y}(x, y) \rightarrow \text{Boolean}$

Parameters

f_x : format of x

f_y : format of y

Operands

x : P3109 value, format f_x

y : P3109 value, format f_y

Output

b : boolean value

Behavior

$\text{compareOp}(\text{NaN}, \text{NaN}) \rightarrow \text{False}$

$\text{compareOp}(\text{NaN}, y) \rightarrow \text{False}$

$\text{compareOp}(x, \text{NaN}) \rightarrow \text{False}$

$\text{compareOp}(x, y) \rightarrow \text{CompareExpr}$

where

$X = \text{Decode}_{f_x}(x)$

$Y = \text{Decode}_{f_y}(y)$

4.10.3 Predicates and classification

Conforming implementations shall provide the classification predicates and the classifier operation defined below.

The classification operations comprise: 1) a set of predicate functions with a boolean return value, taking a single P3109 value as input; 2) a classifier operation $\text{class}(x)$ that returns a single value of enumeration type, describing the input value's properties.

Signature

$$\text{isClass}_f(x) \rightarrow b$$

Parameters

f : format of x

Operands

x : P3109 value, format f

Output

b : boolean value

Behavior

$\text{isZero}(\text{NaN}) \rightarrow \text{False}$

$\text{isZero}(x) \rightarrow \text{Decode}_f(x) = 0$

$\text{isOne}(\text{NaN}) \rightarrow \text{False}$

$\text{isOne}(x) \rightarrow \text{Decode}_f(x) = 1$

$\text{isNaN}(\text{NaN}) \rightarrow \text{True}$

$\text{isNaN}(x) \rightarrow \text{False}$

$\text{isSignMinus}(\text{NaN}) \rightarrow \text{True}$

$\text{isSignMinus}(x) \rightarrow \text{Decode}_f(x) < 0$

$\text{isNormal}(x \in \{0, -\text{Inf}, \text{Inf}, \text{NaN}\}) \rightarrow \text{False}$

$\text{isNormal}(x) \rightarrow (x \bmod 2^{K_f-1}) \div 2^{P_f-1} > 0$

$\text{isSubnormal}(x \in \{0, -\text{Inf}, \text{Inf}, \text{NaN}\}) \rightarrow \text{False}$

$\text{isSubnormal}(x) \rightarrow (x \bmod 2^{K_f-1}) \div 2^{P_f-1} = 0$

$\text{isFinite}(x) \rightarrow x \notin \{-\text{Inf}, \text{Inf}, \text{NaN}\}$

$\text{isInfinite}(x) \rightarrow x \in \{-\text{Inf}, \text{Inf}\}$

$\text{isSignaling}(x) \rightarrow \text{False}$

$\text{isCanonical}(x) \rightarrow \text{True}$

—All P3109 formats have one NaN, which does not signal

—There are no non-canonical P3109 interchange formats

4.10.4 Classifier operation

The Classifier operation $\text{class}(x)$ tells which of the eight classes x falls into as defined by Table 4.

Signature
$$\text{class}_f(x) \rightarrow c$$
Parameters f : format of x **Operands** x : P3109 value, format f **Output** c : enumeration**Behavior** $\text{class}(x) \rightarrow \text{ClassEnum}$

Table 4: Classifier operation

ClassEnum	Condition
clsNaN	isNaN(x)
clsNegativeInfinity	isInfinite(x) \wedge isSignMinus(x)
clsNegativeNormal	isNormal(x) \wedge isSignMinus(x)
clsNegativeSubnormal	isSubnormal(x) \wedge isSignMinus(x)
clsZero	isZero(x)
clsPositiveSubnormal	isSubnormal(x) \wedge \neg isSignMinus(x)
clsPositiveNormal	isNormal(x) \wedge \neg isSignMinus(x)
clsPositiveInfinity	isInfinite(x) \wedge \neg isSignMinus(x)

4.10.5 Total order predicate

The $\text{totalOrder}(x, y)$ predicate provides a total ordering over each P3109 format's value set.

Signature

$$\text{totalOrder}_{f_x, f_y}(x, y) \rightarrow b$$

Operands

f_x : format of x

f_y : format of y

Operands

x : P3109 value, format f_x

y : P3109 value, format f_y

Output

b : boolean value

Behavior

$$\text{totalOrder}(\text{NaN}, x) \rightarrow \text{True}$$
$$\text{totalOrder}(x, \text{NaN}) \rightarrow \text{False}$$
$$\text{totalOrder}(x, y) \rightarrow \text{compareLessEqual}_{f_x, f_y}(x, y)$$

Note

The above definition is consistent with the IEEE Std 754-2019 definition of totalOrder . In particular, among P3109 formats, there is a single NaN and it always compares as the most-negative value.

4.10.6 Comparison predicates

Conforming implementations shall provide the comparison predicates defined by Table 5 and the `totalOrder(x, y)` predicate.

Comparison operations are two-argument predicates, and their negations, that return True or False. Comparisons may be ordered or unordered. A comparison is considered unordered iff either argument is NaN. All other comparisons are ordered.

For $\{=, >, \geq, <, \leq, \leqslant\}$, if any argument is NaN, the result is False.

For $\{\neq, \not>, \not\geq, \not<, \not\leq, \not\leqslant\}$, if any argument is NaN, the result is True.

Otherwise, the result of a comparison shall match the mathematical result.

Table 5: Comparison predicates and negations

Math symbol	Predicate <i>true relations</i>	Math symbol	Negation of predicate <i>true relations</i>
=	compareEqual <i>equal</i>	\neq , NOT =	compareNotEqual <i>less than, greater than, unordered</i>
>	compareGreater <i>greater than</i>	$\not>$, NOT >	compareNotGreater <i>less than, equal, unordered</i>
\geq	compareGreaterEqual <i>greater than, equal</i>	$\not\geq$, NOT \geq	compareLessUnordered <i>less than, unordered</i>
<	compareLess <i>less than</i>	$\not<$, NOT <	compareNotLess <i>greater than, equal, unordered</i>
\leq	compareLessEqual <i>less than, equal</i>	$\not\leq$, NOT \leq	compareGreaterUnordered <i>greater than, unordered</i>
\leqslant	compareOrdered <i>less than, equal, greater than</i>	$\not\leqslant$, NOT \leqslant	compareUnordered <i>unordered</i>

Appendices

A Rationales

A.1 Use of infinity in computation of attention masks

This section expands the rationale in §3.6. A common use for ∞ is to create masks, for example, in Transformer models in machine learning [5].

These values, assembled in mask matrix M with values $M_{ij} \in \{0, -\infty\}$ are typically added to computed values A , in a computation such as:

$$\log \left(\sum \exp(\tau \times (A + M)) \right)$$

where τ is a “temperature” or “base” parameter [6]. This calculation depends on the property $\exp(\tau \times (A_{ij} - \infty)) = 0$.

If a floating-point encoding does not provide infinity, then instead M_{ij} will be replaced by a large float (e.g. 224 is the largest finite binary8p4 value). This is not in itself a difficulty: if all the A values are bounded (e.g. the results of a softmax operation are bounded above by 1.0), then $\exp(1.0 - 224.0)$ is an extremely small number, which will certainly round to zero. Therefore, an explicit representation of infinity is *not* needed in order for this computation to yield its desired value.

However, careful implementations do not execute the calculation as written, and instead fuse the $\log(\sum_i \exp(v_i))$ operation into a single operation $\text{logsumexp}(v)$, whose implementation makes use of the identity transformation

$$\text{logsumexp}(v) \rightarrow \text{logsumexp}(v - \max(v)) + \max(v)$$

Without the “sticky” properties of Inf, this would produce incorrect answers.

For example, in a format where $\text{maxFinite}=240$ without Inf, and $\text{maxFinite}=224$ with Inf:

$$\text{logsumexp}(t * [-224, -\infty]) \rightarrow \text{logsumexp}(t * [0, -\infty])$$

while

$$\text{logsumexp}(t * [-224, -240]) \rightarrow \text{logsumexp}(t * [0, -16])$$

If $t = 1$ and all calculations are done in 8-bit floating-point, then the answer will be the same, because $\exp(-16) \approx 1.1 \times 10^{-7}$, which will round to zero in all precisions $P > 2$; but if t is small, or calculations are done in mixed precision, as is common with 8-bit floating-point, the loss of “stickiness” will silently yield unexpected answers. It is not expected that the full calculation shall be done in 8-bit floating-point, but the subtraction of the maximum value (and computation of the maximum) might reasonably be in 8-bit floating-point.

A.2 Eight Bit Formats

Eight bit floating-point representations have received much attention for their usefulness and efficacy in machine learning, especially deep learning. Various 8-bit floating-point formats have been proposed, investigated in research papers, and some have been modeled in software. Four precisions [3, 4, 5, and 6 bit precisions] have become generally accepted as providing greater operational benefits than the others. Overall [there are exceptions], current efforts focus

more on precisions of 3 and 4 bits (exponent fields of 5 and 4 bits, respectively). The specifics of third-party proposals for 8-bit floating-point representations vary and can differ from one precision to another. Regardless, the precisions emphasized in current research and other third-party work do share the same focus.

A.3 Essential Formats

A.3.1 binary8p3

This format has 3 subnormal and 123 normal magnitudes. The subnormal magnitudes cover $[1.0 * 2^{-17}, 3.0 * 2^{-17}]$. The normal magnitudes cover $[1.0 * 2^{-15}, 4195.0]$. There are 31 normal binades with 4 magnitudes per complete binade [binade 2^{15} has 3 magnitudes, the 4^{th} is used for Inf].

A.3.2 binary8p4

This format has 7 subnormal and 119 normal magnitudes. The subnormal magnitudes cover $[1.0 * 2^{-10}, 7.0 * 2^{-10}]$. The normal magnitudes cover $[1.0 * 2^{-7}, 224.0]$. There are 15 normal binades with 8 magnitudes per complete binade [binade 2^7 has 7 magnitudes, the 8^{th} is used for Inf].

A.3.3 binary8p5

This format has 15 subnormal and 111 normal magnitudes. The subnormal magnitudes cover $[1.0 * 2^{-7}, 15.0 * 2^{-7}]$. The normal magnitudes cover $[0.125, 15.0]$. There are 7 normal binades with 16 magnitudes per complete binade [binade 2^3 has 15 magnitudes, the 16^{th} is used for Inf].

A.3.4 binary8p6

This format has 31 subnormal and 95 normal magnitudes. The subnormal magnitudes cover $[1.0 * 2^{-6}, 31.0 * 2^{-6}]$. The normal magnitudes cover $[0.5, 3.875]$. There are 3 normal binades with 32 magnitudes per complete binade [binade 2^1 has 31 magnitudes, the 32^{nd} is used for Inf].

B External Formats

This table summarizes the points of agreement and of difference between the formats proposed in this document and a number of existing formats, some of which have hardware implementations.

OCP: Open Compute Platform [7], describing hardware implementations including nVidia, Intel, and ARM.

AGQ: AMD, Graphcore, Qualcomm[8], implemented in Graphcore's C600 product, and AMD's gfx940.

TSL: Tesla Dojo Technology [9], A Guide to Tesla's Configurable floating-point Formats & Arithmetic

Format	P3109			OCP		AGQ		TSL	
Subformat	P3	P4	P5	E5	E4	E5	E4	E4	E5
Special values shared by all subformats	Y			N		Y		N	
Exactly one NaN	Y			N		Y		Y	
Positive and negative infinity	Y			Y	N	N		N	
Include negative zero	N			Y		N		N	
Max exponent emax	15	7	3	15	8	15	7	N/A	N/A

C 8-bit Value Tables

Value tables mapping 8-bit strings to value sets are provided in this section.

A typical entry is of the form:

$\text{HEX} = \text{BINARY} = \text{BINARY_FLOAT} = \text{DECIMAL}$
 $0x0b = 0.0001.011 = +0b1.011 \times 2^{-7} = 0.0107421875$

Where the fields are interpreted as follows:

HEX Hexadecimal encoding of the code point
 BINARY Binary expansion of the code point, underscores separate
 < sign > . < exponent > . < significand >
 BINARY_FLOAT The precise float value as a binary fraction followed by 2^e with exponent e
 DECIMAL A decimal expansion of the value. If the expansion is not an exact
 representation of the precise float value, the equals sign is replaced by
 "approximately equals" (\approx).

In addition, entries for subnormal and special values are rendered in color as follows:

$0x05 = 0.0000.101 = +0b0.101 \times 2^{-7} = 0.0048828125$ Subnormal value
 $0x80 = 1.0000.000 = \text{NaN}$ Special value (NaN, +Inf, -Inf)

C.1 Value Table: P1, $\epsilon_{\min} = -62$, $\epsilon_{\max} = 63$ (nonsymmetric)

0x00 = 0.0000000_ = 0.0

0x01 = 0.0000001_ = +0b1.0×2⁻⁶² ≈ 2.1684043E-19
0x02 = 0.0000010_ = +0b1.0×2⁻⁶¹ ≈ 4.3368087E-19
0x03 = 0.0000011_ = +0b1.0×2⁻⁶⁰ ≈ 8.6736174E-19
0x04 = 0.0000100_ = +0b1.0×2⁻⁵⁹ ≈ 1.7347235E-18
0x05 = 0.0000101_ = +0b1.0×2⁻⁵⁸ ≈ 3.469447E-18
0x06 = 0.0000110_ = +0b1.0×2⁻⁵⁷ ≈ 6.9388939E-18
0x07 = 0.0000111_ = +0b1.0×2⁻⁵⁶ ≈ 1.3877788E-17
0x08 = 0.0001000_ = +0b1.0×2⁻⁵⁵ ≈ 2.775576E-17
0x09 = 0.0001001_ = +0b1.0×2⁻⁵⁴ ≈ 5.5511151E-17
0x0a = 0.0001010_ = +0b1.0×2⁻⁵³ ≈ 1.110223E-16
0x0b = 0.0001011_ = +0b1.0×2⁻⁵² ≈ 2.220446E-16
0x0c = 0.0001100_ = +0b1.0×2⁻⁵¹ ≈ 4.4408921E-16
0x0d = 0.0001101_ = +0b1.0×2⁻⁵⁰ ≈ 8.8817842E-16
0x0e = 0.0001110_ = +0b1.0×2⁻⁴⁹ ≈ 1.7763568E-15
0x0f = 0.0001111_ = +0b1.0×2⁻⁴⁸ ≈ 3.5527137E-15
0x10 = 0.0010000_ = +0b1.0×2⁻⁴⁷ ≈ 7.1054274E-15
0x11 = 0.0010001_ = +0b1.0×2⁻⁴⁶ ≈ 1.4210855E-14
0x12 = 0.0010010_ = +0b1.0×2⁻⁴⁵ ≈ 2.8421709E-14
0x13 = 0.0010011_ = +0b1.0×2⁻⁴⁴ ≈ 5.6843419E-14
0x14 = 0.0010100_ = +0b1.0×2⁻⁴³ ≈ 1.1368684E-13
0x15 = 0.0010101_ = +0b1.0×2⁻⁴² ≈ 2.2737368E-13
0x16 = 0.0010110_ = +0b1.0×2⁻⁴¹ ≈ 4.5474735E-13
0x17 = 0.0010111_ = +0b1.0×2⁻⁴⁰ ≈ 9.094947E-13
0x18 = 0.0011000_ = +0b1.0×2⁻³⁹ ≈ 1.8189894E-12
0x19 = 0.0011001_ = +0b1.0×2⁻³⁸ ≈ 3.6379788E-12
0x1a = 0.0011010_ = +0b1.0×2⁻³⁷ ≈ 7.2759576E-12
0x1b = 0.0011011_ = +0b1.0×2⁻³⁶ ≈ 1.4551915E-11
0x1c = 0.0011100_ = +0b1.0×2⁻³⁵ ≈ 2.910383E-11
0x1d = 0.0011101_ = +0b1.0×2⁻³⁴ ≈ 5.8207661E-11
0x1e = 0.0011110_ = +0b1.0×2⁻³³ ≈ 1.1641532E-10
0x1f = 0.0011111_ = +0b1.0×2⁻³² ≈ 2.3283064E-10
0x20 = 0.0100000_ = +0b1.0×2⁻³¹ ≈ 4.6566129E-10
0x21 = 0.0100001_ = +0b1.0×2⁻³⁰ ≈ 9.3132257E-10
0x22 = 0.0100010_ = +0b1.0×2⁻²⁹ ≈ 1.8626451E-09
0x23 = 0.0100011_ = +0b1.0×2⁻²⁸ ≈ 3.7252903E-09
0x24 = 0.0100100_ = +0b1.0×2⁻²⁷ ≈ 7.4505806E-09
0x25 = 0.0100101_ = +0b1.0×2⁻²⁶ ≈ 1.4901161E-08
0x26 = 0.0100110_ = +0b1.0×2⁻²⁵ ≈ 2.9802322E-08
0x27 = 0.0100111_ = +0b1.0×2⁻²⁴ ≈ 5.9604645E-08
0x28 = 0.0101000_ = +0b1.0×2⁻²³ ≈ 1.1920929E-07
0x29 = 0.0101001_ = +0b1.0×2⁻²² ≈ 2.3841858E-07
0x2a = 0.0101010_ = +0b1.0×2⁻²¹ ≈ 4.7683716E-07
0x2b = 0.0101011_ = +0b1.0×2⁻²⁰ ≈ 9.5367432E-07
0x2c = 0.0101100_ = +0b1.0×2⁻¹⁹ ≈ 1.9073486E-06
0x2d = 0.0101101_ = +0b1.0×2⁻¹⁸ ≈ 3.8146973E-06
0x2e = 0.0101110_ = +0b1.0×2⁻¹⁷ ≈ 7.6293945E-06
0x2f = 0.0101111_ = +0b1.0×2⁻¹⁶ ≈ 1.5258789E-05
0x30 = 0.0110000_ = +0b1.0×2⁻¹⁵ ≈ 3.0517578E-05
0x31 = 0.0110001_ = +0b1.0×2⁻¹⁴ ≈ 6.1035156E-05
0x32 = 0.0110010_ = +0b1.0×2⁻¹³ ≈ 0.00012207031
0x33 = 0.0110011_ = +0b1.0×2⁻¹² ≈ 0.000244140625
0x34 = 0.0110100_ = +0b1.0×2⁻¹¹ ≈ 0.00048828125
0x35 = 0.0110101_ = +0b1.0×2⁻¹⁰ ≈ 0.0009765625
0x36 = 0.0110110_ = +0b1.0×2⁻⁹ ≈ 0.001953125
0x37 = 0.0110111_ = +0b1.0×2⁻⁸ ≈ 0.00390625
0x38 = 0.0111000_ = +0b1.0×2⁻⁷ ≈ 0.0078125
0x39 = 0.0111001_ = +0b1.0×2⁻⁶ ≈ 0.015625
0x3a = 0.0111010_ = +0b1.0×2⁻⁵ ≈ 0.03125
0x3b = 0.0111011_ = +0b1.0×2⁻⁴ ≈ 0.0625
0x3c = 0.0111100_ = +0b1.0×2⁻³ ≈ 0.125
0x3d = 0.0111101_ = +0b1.0×2⁻² ≈ 0.25
0x3e = 0.0111110_ = +0b1.0×2⁻¹ ≈ 0.5
0x3f = 0.0111111_ = +0b1.0×2⁰ = 1.0

0x40 = 0.1000000_ = +0b1.0×2¹ = 2.0
0x41 = 0.1000001_ = +0b1.0×2² = 4.0
0x42 = 0.1000010_ = +0b1.0×2³ = 8.0
0x43 = 0.1000011_ = +0b1.0×2⁴ = 16.0
0x44 = 0.1000100_ = +0b1.0×2⁵ = 32.0
0x45 = 0.1000101_ = +0b1.0×2⁶ = 64.0
0x46 = 0.1000110_ = +0b1.0×2⁷ = 128.0
0x47 = 0.1000111_ = +0b1.0×2⁸ = 256.0
0x48 = 0.1001000_ = +0b1.0×2⁹ = 512.0
0x49 = 0.1001001_ = +0b1.0×2¹⁰ = 1024.0
0x4a = 0.1001010_ = +0b1.0×2¹¹ = 2048.0
0x4b = 0.1001011_ = +0b1.0×2¹² = 4096.0
0x4c = 0.1001100_ = +0b1.0×2¹³ = 8192.0
0x4d = 0.1001101_ = +0b1.0×2¹⁴ = 16384.0
0x4e = 0.1001110_ = +0b1.0×2¹⁵ = 32768.0
0x4f = 0.1001111_ = +0b1.0×2¹⁶ = 65536.0
0x50 = 0.1010000_ = +0b1.0×2¹⁷ = 131072.0
0x51 = 0.1010001_ = +0b1.0×2¹⁸ = 262144.0
0x52 = 0.1010010_ = +0b1.0×2¹⁹ = 524288.0
0x53 = 0.1010011_ = +0b1.0×2²⁰ = 1048576.0
0x54 = 0.1010100_ = +0b1.0×2²¹ = 2097152.0
0x55 = 0.1010101_ = +0b1.0×2²² = 4194304.0
0x56 = 0.1010110_ = +0b1.0×2²³ = 8388608.0
0x57 = 0.1010111_ = +0b1.0×2²⁴ = 16777216.0
0x58 = 0.1011000_ = +0b1.0×2²⁵ = 33554432.0
0x59 = 0.1011001_ = +0b1.0×2²⁶ = 67108864.0
0x5a = 0.1011010_ = +0b1.0×2²⁷ = 134217728.0
0x5b = 0.1011011_ = +0b1.0×2²⁸ = 268435456.0
0x5c = 0.1011100_ = +0b1.0×2²⁹ = 536870912.0
0x5d = 0.1011101_ = +0b1.0×2³⁰ = 1073741824.0
0x5e = 0.1011110_ = +0b1.0×2³¹ = 2147483648.0
0x5f = 0.1011111_ = +0b1.0×2³² = 4294967296.0
0x60 = 0.1100000_ = +0b1.0×2³³ = 8589934592.0
0x61 = 0.1100001_ = +0b1.0×2³⁴ = 17179869184.0
0x62 = 0.1100010_ = +0b1.0×2³⁵ = 34359738368.0
0x63 = 0.1100011_ = +0b1.0×2³⁶ = 68719476736.0
0x64 = 0.1100100_ = +0b1.0×2³⁷ = 137438953472.0
0x65 = 0.1100101_ = +0b1.0×2³⁸ = 274877906944.0
0x66 = 0.1100110_ = +0b1.0×2³⁹ = 549755813888.0
0x67 = 0.1100111_ = +0b1.0×2⁴⁰ ≈ 1.0995116e+12
0x68 = 0.1101000_ = +0b1.0×2⁴¹ ≈ 2.1990233e+12
0x69 = 0.1101001_ = +0b1.0×2⁴² ≈ 4.3980465e+12
0x6a = 0.1101010_ = +0b1.0×2⁴³ ≈ 8.796093e+12
0x6b = 0.1101011_ = +0b1.0×2⁴⁴ ≈ 1.7592186e+13
0x6c = 0.1101100_ = +0b1.0×2⁴⁵ ≈ 3.5184372e+13
0x6d = 0.1101101_ = +0b1.0×2⁴⁶ ≈ 7.0368744e+13
0x6e = 0.1101110_ = +0b1.0×2⁴⁷ ≈ 1.4073749e+14
0x6f = 0.1101111_ = +0b1.0×2⁴⁸ ≈ 2.8147498e+14
0x70 = 0.1110000_ = +0b1.0×2⁴⁹ ≈ 5.6294995e+14
0x71 = 0.1110001_ = +0b1.0×2⁵⁰ ≈ 1.1258999e+15
0x72 = 0.1110010_ = +0b1.0×2⁵¹ ≈ 2.2517998e+15
0x73 = 0.1110011_ = +0b1.0×2⁵² ≈ 4.5035996e+15
0x74 = 0.1110100_ = +0b1.0×2⁵³ ≈ 9.0071993e+15
0x75 = 0.1110101_ = +0b1.0×2⁵⁴ ≈ 1.8014399e+16
0x76 = 0.1110110_ = +0b1.0×2⁵⁵ ≈ 3.6028797e+16
0x77 = 0.1110111_ = +0b1.0×2⁵⁶ ≈ 7.2057594e+16
0x78 = 0.1111000_ = +0b1.0×2⁵⁷ ≈ 1.4411519e+17
0x79 = 0.1111001_ = +0b1.0×2⁵⁸ ≈ 2.8823038e+17
0x7a = 0.1111010_ = +0b1.0×2⁵⁹ ≈ 5.7646075e+17
0x7b = 0.1111011_ = +0b1.0×2⁶⁰ ≈ 1.1529215e+18
0x7c = 0.1111100_ = +0b1.0×2⁶¹ ≈ 2.305843e+18
0x7d = 0.1111101_ = +0b1.0×2⁶² ≈ 4.611686e+18
0x7e = 0.1111110_ = +0b1.0×2⁶³ ≈ 9.223372e+18
0x7f = 0.1111111_ = +Inf

0x80 = 1.0000000_ = NaN

0x81 = 1.0000001_ = -0b1.0×2⁻⁶² ≈ -2.1684043E-19
0x82 = 1.0000010_ = -0b1.0×2⁻⁶¹ ≈ -4.3368087E-19
0x83 = 1.0000011_ = -0b1.0×2⁻⁶⁰ ≈ -8.6736174E-19
0x84 = 1.0000100_ = -0b1.0×2⁻⁵⁹ ≈ -1.7347235E-18
0x85 = 1.0000101_ = -0b1.0×2⁻⁵⁸ ≈ -3.469447E-18
0x86 = 1.0000110_ = -0b1.0×2⁻⁵⁷ ≈ -6.9388939E-18
0x87 = 1.0000111_ = -0b1.0×2⁻⁵⁶ ≈ -1.3877788E-17
0x88 = 1.0001000_ = -0b1.0×2⁻⁵⁵ ≈ -2.775576E-17
0x89 = 1.0001001_ = -0b1.0×2⁻⁵⁴ ≈ -5.5511151E-17
0x8a = 1.0001010_ = -0b1.0×2⁻⁵³ ≈ -1.110223E-16
0x8b = 1.0001011_ = -0b1.0×2⁻⁵² ≈ -2.220446E-16
0x8c = 1.0001100_ = -0b1.0×2⁻⁵¹ ≈ -4.4408921E-16
0x8d = 1.0001101_ = -0b1.0×2⁻⁵⁰ ≈ -8.8817842E-16
0x8e = 1.0001110_ = -0b1.0×2⁻⁴⁹ ≈ -1.7763568E-15
0x8f = 1.0001111_ = -0b1.0×2⁻⁴⁸ ≈ -3.5527137E-15
0x90 = 1.0010000_ = -0b1.0×2⁻⁴⁷ ≈ -7.1054274E-15
0x91 = 1.0010001_ = -0b1.0×2⁻⁴⁶ ≈ -1.4210855E-14
0x92 = 1.0010010_ = -0b1.0×2⁻⁴⁵ ≈ -2.8421709E-14
0x93 = 1.0010011_ = -0b1.0×2⁻⁴⁴ ≈ -5.6843419E-14
0x94 = 1.0010100_ = -0b1.0×2⁻⁴³ ≈ -1.1368684E-13
0x95 = 1.0010101_ = -0b1.0×2⁻⁴² ≈ -2.2737368E-13
0x96 = 1.0010110_ = -0b1.0×2⁻⁴¹ ≈ -4.5474735E-13
0x97 = 1.0010111_ = -0b1.0×2⁻⁴⁰ ≈ -9.094947E-13
0x98 = 1.0011000_ = -0b1.0×2⁻³⁹ ≈ -1.8189894E-12
0x99 = 1.0011001_ = -0b1.0×2⁻³⁸ ≈ -3.6379788E-12
0x9a = 1.0011010_ = -0b1.0×2⁻³⁷ ≈ -7.2759576E-12
0x9b = 1.0011011_ = -0b1.0×2⁻³⁶ ≈ -1.4551915E-11
0x9c = 1.0011100_ = -0b1.0×2⁻³⁵ ≈ -2.910383E-11
0x9d = 1.0011101_ = -0b1.0×2⁻³⁴ ≈ -5.8207661E-11
0x9e = 1.0011110_ = -0b1.0×2⁻³³ ≈ -1.1641532E-10
0x9f = 1.0011111_ = -0b1.0×2⁻³² ≈ -2.3283064E-10
0xa0 = 1.0100000_ = -0b1.0×2⁻³¹ ≈ -4.6566129E-10
0xa1 = 1.0100001_ = -0b1.0×2⁻³⁰ ≈ -9.3132257E-10
0xa2 = 1.0100010_ = -0b1.0×2⁻²⁹ ≈ -1.8626451E-09
0xa3 = 1.0100011_ = -0b1.0×2⁻²⁸ ≈ -3.7252903E-09
0xa4 = 1.0100100_ = -0b1.0×2⁻²⁷ ≈ -7.4505806E-09
0xa5 = 1.0100101_ = -0b1.0×2⁻²⁶ ≈ -1.4901161E-08
0xa6 = 1.0100110_ = -0b1.0×2⁻²⁵ ≈ -2.9802322E-08
0xa7 = 1.0100111_ = -0b1.0×2⁻²⁴ ≈ -5.9604645E-08
0xa8 = 1.0101000_ = -0b1.0×2⁻²³ ≈ -1.1920929E-07
0xa9 = 1.0101001_ = -0b1.0×2⁻²² ≈ -2.3841858E-07
0xaa = 1.0101010_ = -0b1.0×2⁻²¹ ≈ -4.7683716E-07
0xab = 1.0101011_ = -0b1.0×2⁻²⁰ ≈ -9.5367432E-07
0xac = 1.0101100_ = -0b1.0×2⁻¹⁹ ≈ -1.9073486E-06
0xad = 1.0101101_ = -0b1.0×2⁻¹⁸ ≈ -3.8146973E-06
0xae = 1.0101110_ = -0b1.0×2⁻¹⁷ ≈ -7.6293945E-06
0xaf = 1.0101111_ = -0b1.0×2⁻¹⁶ ≈ -1.5258789E-05
0xb0 = 1.0110000_ = -0b1.0×2⁻¹⁵ ≈ -3.0517578E-05
0xb1 = 1.0110001_ = -0b1.0×2⁻¹⁴ ≈ -6.1035156E-05
0xb2 = 1.0110010_ = -0b1.0×2⁻¹³ ≈ -0.00012207031
0xb3 = 1.0110011_ = -0b1.0×2⁻¹² ≈ -0.000244140625
0xb4 = 1.0110100_ = -0b1.0×2⁻¹¹ ≈ -0.00048828125
0xb5 = 1.0110101_ = -0b1.0×2⁻¹⁰ ≈ -0.0009765625
0xb6 = 1.0110110_ = -0b1.0×2⁻⁹ ≈ -0.001953125
0xb7 = 1.0110111_ = -0b1.0×2⁻⁸ ≈ -0.00390625
0xb8 = 1.0111000_ = -0b1.0×2⁻⁷ ≈ -0.0078125
0xb9 = 1.0111001_ = -0b1.0×2⁻⁶ ≈ -0.015625
0xba = 1.0111010_ = -0b1.0×2⁻⁵ ≈ -0.03125
0xbb = 1.0111011_ = -0b1.0×2⁻⁴ ≈ -0.0625
0xbc = 1.0111100_ = -0b1.0×2⁻³ ≈ -0.125
0xbd = 1.0111101_ = -0b1.0×2⁻² ≈ -0.25
0xbe = 1.0111110_ = -0b1.0×2⁻¹ ≈ -0.5
0xbf = 1.0111111_ = -0b1.0×2⁰ = -1.0

0xc0 = 1.1000000_ = -0b1.0×2¹ = -2.0
0xc1 = 1.1000001_ = -0b1.0×2² = -4.0
0xc2 = 1.1000010_ = -0b1.0×2³ = -8.0
0xc3 = 1.1000011_ = -0b1.0×2⁴ = -16.0
0xc4 = 1.1000100_ = -0b1.0×2⁵ = -32.0
0xc5 = 1.1000101_ = -0b1.0×2⁶ = -64.0
0xc6 = 1.1000110_ = -0b1.0×2⁷ = -128.0
0xc7 = 1.1000111_ = -0b1.0×2⁸ = -256.0
0xc8 = 1.1001000_ = -0b1.0×2⁹ = -512.0
0xc9 = 1.1001001_ = -0b1.0×2¹⁰ = -1024.0
0xca = 1.1001010_ = -0b1.0×2¹¹ = -2048.0
0xcb = 1.1001011_ = -0b1.0×2¹² = -4096.0
0xcc = 1.1001100_ = -0b1.0×2¹³ = -8192.0
0xcd = 1.1001101_ = -0b1.0×2¹⁴ = -16384.0
0xce = 1.1001110_ = -0b1.0×2¹⁵ = -32768.0
0xcf = 1.1001111_ = -0b1.0×2¹⁶ = -65536.0
0xd0 = 1.1010000_ = -0b1.0×2¹⁷ = -131072.0
0xd1 = 1.1010001_ = -0b1.0×2¹⁸ = -262144.0
0xd2 = 1.1010010_ = -0b1.0×2¹⁹ = -524288.0
0xd3 = 1.1010011_ = -0b1.0×2²⁰ = -1048576.0
0xd4 = 1.1010100_ = -0b1.0×2²¹ = -2097152.0
0xd5 = 1.1010101_ = -0b1.0×2²² = -4194304.0
0xd6 = 1.1010110_ = -0b1.0×2²³ = -8388608.0
0xd7 = 1.1010111_ = -0b1.0×2²⁴ = -16777216.0
0xd8 = 1.1011000_ = -0b1.0×2²⁵ = -33554432.0
0xd9 = 1.1011001_ = -0b1.0×2²⁶ = -67108864.0
0xda = 1.1011010_ = -0b1.0×2²⁷ = -134217728.0
0xdb = 1.1011011_ = -0b1.0×2²⁸ = -268435456.0
0xdc = 1.1011100_ = -0b1.0×2²⁹ = -536870912.0
0xdd = 1.1011101_ = -0b1.0×2³⁰ = -1073741824.0
0xde = 1.1011110_ = -0b1.0×2³¹ = -2147483648.0
0xdf = 1.1011111_ = -0b1.0×2³² = -4294967296.0
0xe0 = 1.1100000_ = -0b1.0×2³³ = -8589934592.0
0xe1 = 1.1100001_ = -0b1.0×2³⁴ = -17179869184.0
0xe2 = 1.1100010_ = -0b1.0×2³⁵ = -34359738368.0
0xe3 = 1.1100011_ = -0b1.0×2³⁶ = -68719476736.0
0xe4 = 1.1100100_ = -0b1.0×2³⁷ ≈ -1.3743895e+11
0xe5 = 1.1100101_ = -0b1.0×2³⁸ ≈ -2.7487791e+11
0xe6 = 1.1100110_ = -0b1.0×2³⁹ ≈ -5.4975581e+11
0xe7 = 1.1100111_ = -0b1.0×2⁴⁰ ≈ -1.0995116e+12
0xe8 = 1.1101000_ = -0b1.0×2⁴¹ ≈ -2.1990233e+12
0xe9 = 1.1101001_ = -0b1.0×2⁴² ≈ -4.3980465e+12
0xea = 1.1101010_ = -0b1.0×2⁴³ ≈ -8.796093e+12
0xeb = 1.1101011_ = -0b1.0×2⁴⁴ ≈ -1.7592186e+13
0xec = 1.1101100_ = -0b1.0×2⁴⁵ ≈ -3.5184372e+13
0xed = 1.1101101_ = -0b1.0×2⁴⁶ ≈ -7.0368744e+13
0xee = 1.1101110_ = -0b1.0×2⁴⁷ ≈ -1.4073749e+14
0xef = 1.1101111_ = -0b1.0×2⁴⁸ ≈ -2.8147498e+14
0xf0 = 1.1110000_ = -0b1.0×2⁴⁹ ≈ -5.6294995e+14
0xf1 = 1.1110001_ = -0b1.0×2⁵⁰ ≈ -1.1258999e+15
0xf2 = 1.1110010_ = -0b1.0×2⁵¹ ≈ -2.2517998e+15
0xf3 = 1.1110011_ = -0b1.0×2⁵² ≈ -4.5035996e+15
0xf4 = 1.1110100_ = -0b1.0×2⁵³ ≈ -9.0071993e+15
0xf5 = 1.1110101_ = -0b1.0×2⁵⁴ ≈ -1.8014399e+16
0xf6 = 1.1110110_ = -0b1.0×2⁵⁵ ≈ -3.6028797e+16
0xf7 = 1.1110111_ = -0b1.0×2⁵⁶ ≈ -7.2057594e+16
0xf8 = 1.1111000_ = -0b1.0×2⁵⁷ ≈ -1.4411519e+17
0xf9 = 1.1111001_ = -0b1.0×2⁵⁸ ≈ -2.8823038e+17
0xfa = 1.1111010_ = -0b1.0×2⁵⁹ ≈ -5.7646075e+17
0xfb = 1.1111011_ = -0b1.0×2⁶⁰ ≈ -1.1529215e+18
0xfc = 1.1111100_ = -0b1.0×2⁶¹ ≈ -2.305843e+18
0xfd = 1.1111101_ = -0b1.0×2⁶² ≈ -4.611686e+18
0xfe = 1.1111110_ = -0b1.0×2⁶³ ≈ -9.223372e+18
0xff = 1.1111111_ = -Inf

C.2 Value Table: P2, $e_{\min} = -31$, $e_{\max} = 31$

0x00 = 0.000000.0 = 0.0	0x40 = 0.100000.0 = +0b1.0×2 ⁰ = 1.0	0x80 = 1.000000.0 = NaN	0xc0 = 1.100000.0 = -0b1.0×2 ⁰ = -1.0
0x01 = 0.000000.1 = +0b0.1×2 ⁻³¹ ≈ 2.3283064E-10	0x41 = 0.100000.1 = +0b1.1×2 ⁰ = 1.5	0x81 = 1.000000.1 = -0b0.1×2 ⁻³¹ ≈ -2.3283064E-10	0xc1 = 1.100000.1 = -0b1.1×2 ⁰ = -1.5
0x02 = 0.000001.0 = +0b1.0×2 ⁻³¹ ≈ 4.6566129E-10	0x42 = 0.100001.0 = +0b1.0×2 ¹ = 2.0	0x82 = 1.000001.0 = -0b1.0×2 ⁻³¹ ≈ -4.6566129E-10	0xc2 = 1.100001.0 = -0b1.0×2 ¹ = -2.0
0x03 = 0.000001.1 = +0b1.1×2 ⁻³¹ ≈ 6.9849193E-10	0x43 = 0.100001.1 = +0b1.1×2 ¹ = 3.0	0x83 = 1.000001.1 = -0b1.1×2 ⁻³¹ ≈ -6.9849193E-10	0xc3 = 1.100001.1 = -0b1.1×2 ¹ = -3.0
0x04 = 0.000010.0 = +0b1.0×2 ⁻³⁰ ≈ 9.3132257E-10	0x44 = 0.100010.0 = +0b1.0×2 ² = 4.0	0x84 = 1.000010.0 = -0b1.0×2 ⁻³⁰ ≈ -9.3132257E-10	0xc4 = 1.100010.0 = -0b1.0×2 ² = -4.0
0x05 = 0.000010.1 = +0b1.1×2 ⁻³⁰ ≈ 1.3969839E-09	0x45 = 0.100010.1 = +0b1.1×2 ² = 6.0	0x85 = 1.000010.1 = -0b1.1×2 ⁻³⁰ ≈ -1.3969839E-09	0xc5 = 1.100010.1 = -0b1.1×2 ² = -6.0
0x06 = 0.000011.0 = +0b1.0×2 ⁻²⁹ ≈ 1.8626451E-09	0x46 = 0.100011.0 = +0b1.0×2 ³ = 8.0	0x86 = 1.000011.0 = -0b1.0×2 ⁻²⁹ ≈ -1.8626451E-09	0xc6 = 1.100011.0 = -0b1.0×2 ³ = -8.0
0x07 = 0.000011.1 = +0b1.1×2 ⁻²⁹ ≈ 2.7939677E-09	0x47 = 0.100011.1 = +0b1.1×2 ³ = 12.0	0x87 = 1.000011.1 = -0b1.1×2 ⁻²⁹ ≈ -2.7939677E-09	0xc7 = 1.100011.1 = -0b1.1×2 ³ = -12.0
0x08 = 0.000010.1 = +0b1.0×2 ⁻²⁸ ≈ 3.7252903E-09	0x48 = 0.100100.0 = +0b1.0×2 ⁴ = 16.0	0x88 = 1.000100.0 = -0b1.0×2 ⁻²⁸ ≈ -3.7252903E-09	0xc8 = 1.100100.0 = -0b1.0×2 ⁴ = -16.0
0x09 = 0.000100.1 = +0b1.1×2 ⁻²⁸ ≈ 5.5879354E-09	0x49 = 0.100100.1 = +0b1.1×2 ⁴ = 24.0	0x89 = 1.000100.1 = -0b1.1×2 ⁻²⁸ ≈ -5.5879354E-09	0xc9 = 1.100100.1 = -0b1.1×2 ⁴ = -24.0
0x0a = 0.000101.0 = +0b1.0×2 ⁻²⁷ ≈ 7.4505806E-09	0x4a = 0.100101.0 = +0b1.0×2 ⁵ = 32.0	0x8a = 1.000101.0 = -0b1.0×2 ⁻²⁷ ≈ -7.4505806E-09	0xca = 1.100101.0 = -0b1.0×2 ⁵ = -32.0
0x0b = 0.000101.1 = +0b1.1×2 ⁻²⁷ ≈ 1.1175871E-08	0x4b = 0.100101.1 = +0b1.1×2 ⁵ = 48.0	0x8b = 1.000101.1 = -0b1.1×2 ⁻²⁷ ≈ -1.1175871E-08	0xcb = 1.100101.1 = -0b1.1×2 ⁵ = -48.0
0x0c = 0.000110.0 = +0b1.0×2 ⁻²⁶ ≈ 1.4901161E-08	0x4c = 0.100110.0 = +0b1.0×2 ⁶ = 64.0	0x8c = 1.000110.0 = -0b1.0×2 ⁻²⁶ ≈ -1.4901161E-08	0xcc = 1.100110.0 = -0b1.0×2 ⁶ = -64.0
0x0d = 0.000110.1 = +0b1.1×2 ⁻²⁶ ≈ 2.2351742E-08	0x4d = 0.100110.1 = +0b1.1×2 ⁶ = 96.0	0x8d = 1.000110.1 = -0b1.1×2 ⁻²⁶ ≈ -2.2351742E-08	0xcd = 1.100110.1 = -0b1.1×2 ⁶ = -96.0
0x0e = 0.000111.0 = +0b1.0×2 ⁻²⁵ ≈ 2.9802322E-08	0x4e = 0.100111.0 = +0b1.0×2 ⁷ = 128.0	0x8e = 1.000111.0 = -0b1.0×2 ⁻²⁵ ≈ -2.9802322E-08	0xce = 1.100111.0 = -0b1.0×2 ⁷ = -128.0
0x0f = 0.000111.1 = +0b1.1×2 ⁻²⁵ ≈ 4.4703484E-08	0x4f = 0.100111.1 = +0b1.1×2 ⁷ = 192.0	0x8f = 1.000111.1 = -0b1.1×2 ⁻²⁵ ≈ -4.4703484E-08	0xcf = 1.100111.1 = -0b1.1×2 ⁷ = -192.0
0x10 = 0.001000.0 = +0b1.0×2 ⁻²⁴ ≈ 5.9604645E-08	0x50 = 0.101000.0 = +0b1.0×2 ⁸ = 256.0	0x90 = 1.001000.0 = -0b1.0×2 ⁻²⁴ ≈ -5.9604645E-08	0xd0 = 1.101000.0 = -0b1.0×2 ⁸ = -256.0
0x11 = 0.001000.1 = +0b1.1×2 ⁻²⁴ ≈ 8.9406967E-08	0x51 = 0.101000.1 = +0b1.1×2 ⁸ = 384.0	0x91 = 1.001000.1 = -0b1.1×2 ⁻²⁴ ≈ -8.9406967E-08	0xd1 = 1.101000.1 = -0b1.1×2 ⁸ = -384.0
0x12 = 0.001001.0 = +0b1.0×2 ⁻²³ ≈ 1.1920929E-07	0x52 = 0.101001.0 = +0b1.0×2 ⁹ = 512.0	0x92 = 1.001001.0 = -0b1.0×2 ⁻²³ ≈ -1.1920929E-07	0xd2 = 1.101001.0 = -0b1.0×2 ⁹ = -512.0
0x13 = 0.001001.1 = +0b1.1×2 ⁻²³ ≈ 1.7881393E-07	0x53 = 0.101001.1 = +0b1.1×2 ⁹ = 768.0	0x93 = 1.001001.1 = -0b1.1×2 ⁻²³ ≈ -1.7881393E-07	0xd3 = 1.101001.1 = -0b1.1×2 ⁹ = -768.0
0x14 = 0.001010.0 = +0b1.0×2 ⁻²² ≈ 2.3841858E-07	0x54 = 0.101010.0 = +0b1.0×2 ¹⁰ = 1024.0	0x94 = 1.001010.0 = -0b1.0×2 ⁻²² ≈ -2.3841858E-07	0xd4 = 1.101010.0 = -0b1.0×2 ¹⁰ = -1024.0
0x15 = 0.001010.1 = +0b1.1×2 ⁻²² ≈ 3.5762787E-07	0x55 = 0.101010.1 = +0b1.1×2 ¹⁰ = 1536.0	0x95 = 1.001010.1 = -0b1.1×2 ⁻²² ≈ -3.5762787E-07	0xd5 = 1.101010.1 = -0b1.1×2 ¹⁰ = -1536.0
0x16 = 0.001011.0 = +0b1.0×2 ⁻²¹ ≈ 4.7683716E-07	0x56 = 0.101011.0 = +0b1.0×2 ¹¹ = 2048.0	0x96 = 1.001011.0 = -0b1.0×2 ⁻²¹ ≈ -4.7683716E-07	0xd6 = 1.101011.0 = -0b1.0×2 ¹¹ = -2048.0
0x17 = 0.001011.1 = +0b1.1×2 ⁻²¹ ≈ 7.1525574E-07	0x57 = 0.101011.1 = +0b1.1×2 ¹¹ = 3072.0	0x97 = 1.001011.1 = -0b1.1×2 ⁻²¹ ≈ -7.1525574E-07	0xd7 = 1.101011.1 = -0b1.1×2 ¹¹ = -3072.0
0x18 = 0.001100.0 = +0b1.0×2 ⁻²⁰ ≈ 9.5367432E-07	0x58 = 0.101100.0 = +0b1.0×2 ¹² = 4096.0	0x98 = 1.001100.0 = -0b1.0×2 ⁻²⁰ ≈ -9.5367432E-07	0xd8 = 1.101100.0 = -0b1.0×2 ¹² = -4096.0
0x19 = 0.001100.1 = +0b1.1×2 ⁻²⁰ ≈ 1.4305115E-06	0x59 = 0.101100.1 = +0b1.1×2 ¹² = 6144.0	0x99 = 1.001100.1 = -0b1.1×2 ⁻²⁰ ≈ -1.4305115E-06	0xd9 = 1.101100.1 = -0b1.1×2 ¹² = -6144.0
0x1a = 0.001101.0 = +0b1.0×2 ⁻¹⁹ ≈ 1.9073486E-06	0x5a = 0.101101.0 = +0b1.0×2 ¹³ = 8192.0	0x9a = 1.001101.0 = -0b1.0×2 ⁻¹⁹ ≈ -1.9073486E-06	0xda = 1.101101.0 = -0b1.0×2 ¹³ = -8192.0
0x1b = 0.001101.1 = +0b1.1×2 ⁻¹⁹ ≈ 2.8610229E-06	0x5b = 0.101101.1 = +0b1.1×2 ¹³ = 12288.0	0x9b = 1.001101.1 = -0b1.1×2 ⁻¹⁹ ≈ -2.8610229E-06	0xdb = 1.101101.1 = -0b1.1×2 ¹³ = -12288.0
0x1c = 0.001110.0 = +0b1.0×2 ⁻¹⁸ ≈ 3.8146973E-06	0x5c = 0.101110.0 = +0b1.0×2 ¹⁴ = 16384.0	0x9c = 1.001110.0 = -0b1.0×2 ⁻¹⁸ ≈ -3.8146973E-06	0xdc = 1.101110.0 = -0b1.0×2 ¹⁴ = -16384.0
0x1d = 0.001110.1 = +0b1.1×2 ⁻¹⁸ ≈ 5.7220459E-06	0x5d = 0.101110.1 = +0b1.1×2 ¹⁴ = 24576.0	0x9d = 1.001110.1 = -0b1.1×2 ⁻¹⁸ ≈ -5.7220459E-06	0xdd = 1.101110.1 = -0b1.1×2 ¹⁴ = -24576.0
0x1e = 0.001111.0 = +0b1.0×2 ⁻¹⁷ ≈ 7.6293945E-06	0x5e = 0.101111.0 = +0b1.0×2 ¹⁵ = 32768.0	0x9e = 1.001111.0 = -0b1.0×2 ⁻¹⁷ ≈ -7.6293945E-06	0xde = 1.101111.0 = -0b1.0×2 ¹⁵ = -32768.0
0x1f = 0.001111.1 = +0b1.1×2 ⁻¹⁷ ≈ 1.1444092E-05	0x5f = 0.101111.1 = +0b1.1×2 ¹⁵ = 49152.0	0x9f = 1.001111.1 = -0b1.1×2 ⁻¹⁷ ≈ -1.1444092E-05	0xdf = 1.101111.1 = -0b1.1×2 ¹⁵ = -49152.0
0x20 = 0.010000.0 = +0b1.0×2 ⁻¹⁶ ≈ 1.5258789E-05	0x60 = 0.110000.0 = +0b1.0×2 ¹⁶ = 65536.0	0xa0 = 1.010000.0 = -0b1.0×2 ⁻¹⁶ ≈ -1.5258789E-05	0xe0 = 1.110000.0 = -0b1.0×2 ¹⁶ = -65536.0
0x21 = 0.010000.1 = +0b1.1×2 ⁻¹⁶ ≈ 2.2888184E-05	0x61 = 0.110000.1 = +0b1.1×2 ¹⁶ = 98304.0	0xa1 = 1.010000.1 = -0b1.1×2 ⁻¹⁶ ≈ -2.2888184E-05	0xe1 = 1.110000.1 = -0b1.1×2 ¹⁶ = -98304.0
0x22 = 0.010001.0 = +0b1.0×2 ⁻¹⁵ ≈ 3.0517578E-05	0x62 = 0.110001.0 = +0b1.0×2 ¹⁷ = 131072.0	0xa2 = 1.010001.0 = -0b1.0×2 ⁻¹⁵ ≈ -3.0517578E-05	0xe2 = 1.110001.0 = -0b1.0×2 ¹⁷ = -131072.0
0x23 = 0.010001.1 = +0b1.1×2 ⁻¹⁵ ≈ 4.5776367E-05	0x63 = 0.110001.1 = +0b1.1×2 ¹⁷ = 196608.0	0xa3 = 1.010001.1 = -0b1.1×2 ⁻¹⁵ ≈ -4.5776367E-05	0xe3 = 1.110001.1 = -0b1.1×2 ¹⁷ = -196608.0
0x24 = 0.010010.0 = +0b1.0×2 ⁻¹⁴ ≈ 6.1035156E-05	0x64 = 0.110010.0 = +0b1.0×2 ¹⁸ = 262144.0	0xa4 = 1.010010.0 = -0b1.0×2 ⁻¹⁴ ≈ -6.1035156E-05	0xe4 = 1.110010.0 = -0b1.0×2 ¹⁸ = -262144.0
0x25 = 0.010010.1 = +0b1.1×2 ⁻¹⁴ ≈ 9.1552734E-05	0x65 = 0.110010.1 = +0b1.1×2 ¹⁸ = 393216.0	0xa5 = 1.010010.1 = -0b1.1×2 ⁻¹⁴ ≈ -9.1552734E-05	0xe5 = 1.110010.1 = -0b1.1×2 ¹⁸ = -393216.0
0x26 = 0.010011.0 = +0b1.0×2 ⁻¹³ ≈ 0.00012207031	0x66 = 0.110011.0 = +0b1.0×2 ¹⁹ = 524288.0	0xa6 = 1.010011.0 = -0b1.0×2 ⁻¹³ ≈ -0.00012207031	0xe6 = 1.110011.0 = -0b1.0×2 ¹⁹ = -524288.0
0x27 = 0.010011.1 = +0b1.1×2 ⁻¹³ ≈ 0.00018310547	0x67 = 0.110011.1 = +0b1.1×2 ¹⁹ = 786432.0	0xa7 = 1.010011.1 = -0b1.1×2 ⁻¹³ ≈ -0.00018310547	0xe7 = 1.110011.1 = -0b1.1×2 ¹⁹ = -786432.0
0x28 = 0.010100.0 = +0b1.0×2 ⁻¹² ≈ 0.000244140625	0x68 = 0.110100.0 = +0b1.0×2 ²⁰ = 1048576.0	0xa8 = 1.010100.0 = -0b1.0×2 ⁻¹² ≈ -0.000244140625	0xe8 = 1.110100.0 = -0b1.0×2 ²⁰ = -1048576.0
0x29 = 0.010100.1 = +0b1.1×2 ⁻¹² ≈ 0.00036621094	0x69 = 0.110100.1 = +0b1.1×2 ²⁰ = 1572864.0	0xa9 = 1.010100.1 = -0b1.1×2 ⁻¹² ≈ -0.00036621094	0xe9 = 1.110100.1 = -0b1.1×2 ²⁰ = -1572864.0
0x2a = 0.010101.0 = +0b1.0×2 ⁻¹¹ ≈ 0.00048828125	0x6a = 0.110101.0 = +0b1.0×2 ²¹ = 2097152.0	0xaa = 1.010101.0 = -0b1.0×2 ⁻¹¹ ≈ -0.00048828125	0xea = 1.110101.0 = -0b1.0×2 ²¹ = -2097152.0
0x2b = 0.010101.1 = +0b1.1×2 ⁻¹¹ ≈ 0.000732421875	0x6b = 0.110101.1 = +0b1.1×2 ²¹ = 3145728.0	0xab = 1.010101.1 = -0b1.1×2 ⁻¹¹ ≈ -0.000732421875	0xeb = 1.110101.1 = -0b1.1×2 ²¹ = -3145728.0
0x2c = 0.010110.0 = +0b1.0×2 ⁻¹⁰ ≈ 0.0009765625	0x6c = 0.110110.0 = +0b1.0×2 ²² = 4194304.0	0xac = 1.010110.0 = -0b1.0×2 ⁻¹⁰ ≈ -0.0009765625	0xec = 1.110110.0 = -0b1.0×2 ²² = -4194304.0
0x2d = 0.010110.1 = +0b1.1×2 ⁻¹⁰ ≈ 0.00146484375	0x6d = 0.110110.1 = +0b1.1×2 ²² = 6291456.0	0xad = 1.010110.1 = -0b1.1×2 ⁻¹⁰ ≈ -0.00146484375	0xed = 1.110110.1 = -0b1.1×2 ²² = -6291456.0
0x2e = 0.010111.0 = +0b1.0×2 ⁻⁹ ≈ 0.001953125	0x6e = 0.110111.0 = +0b1.0×2 ²³ = 8388608.0	0xae = 1.010111.0 = -0b1.0×2 ⁻⁹ ≈ -0.001953125	0xee = 1.110111.0 = -0b1.0×2 ²³ = -8388608.0
0x2f = 0.010111.1 = +0b1.1×2 ⁻⁹ ≈ 0.0029296875	0x6f = 0.110111.1 = +0b1.1×2 ²³ = 12582912.0	0xaf = 1.010111.1 = -0b1.1×2 ⁻⁹ ≈ -0.0029296875	0xef = 1.110111.1 = -0b1.1×2 ²³ = -12582912.0
0x30 = 0.011000.0 = +0b1.0×2 ⁻⁸ ≈ 0.00390625	0x70 = 0.111000.0 = +0b1.0×2 ²⁴ = 16777216.0	0xb0 = 1.011000.0 = -0b1.0×2 ⁻⁸ ≈ -0.00390625	0xf0 = 1.111000.0 = -0b1.0×2 ²⁴ = -16777216.0
0x31 = 0.011000.1 = +0b1.1×2 ⁻⁸ ≈ 0.005859375	0x71 = 0.111000.1 = +0b1.1×2 ²⁴ = 25165824.0	0xb1 = 1.011000.1 = -0b1.1×2 ⁻⁸ ≈ -0.005859375	0xf1 = 1.111000.1 = -0b1.1×2 ²⁴ = -25165824.0
0x32 = 0.011001.0 = +0b1.0×2 ⁻⁷ ≈ 0.0078125	0x72 = 0.111001.0 = +0b1.0×2 ²⁵ = 33554432.0	0xb2 = 1.011001.0 = -0b1.0×2 ⁻⁷ ≈ -0.0078125	0xf2 = 1.111001.0 = -0b1.0×2 ²⁵ = -33554432.0
0x33 = 0.011001.1 = +0b1.1×2 ⁻⁷ ≈ 0.01171875	0x73 = 0.111001.1 = +0b1.1×2 ²⁵ = 50331648.0	0xb3 = 1.011001.1 = -0b1.1×2 ⁻⁷ ≈ -0.01171875	0xf3 = 1.111001.1 = -0b1.1×2 ²⁵ = -50331648.0
0x34 = 0.011010.0 = +0b1.0×2 ⁻⁶ ≈ 0.015625	0x74 = 0.111010.0 = +0b1.0×2 ²⁶ = 67108864.0	0xb4 = 1.011010.0 = -0b1.0×2 ⁻⁶ ≈ -0.015625	0xf4 = 1.111010.0 = -0b1.0×2 ²⁶ = -67108864.0
0x35 = 0.011010.1 = +0b1.1×2 ⁻⁶ ≈ 0.0234375	0x75 = 0.111010.1 = +0b1.1×2 ²⁶ = 100663296.0	0xb5 = 1.011010.1 = -0b1.1×2 ⁻⁶ ≈ -0.0234375	0xf5 = 1.111010.1 = -0b1.1×2 ²⁶ = -100663296.0
0x36 = 0.011011.0 = +0b1.0×2 ⁻⁵ ≈ 0.03125	0x76 = 0.111011.0 = +0b1.0×2 ²⁷ = 134217728.0	0xb6 = 1.011011.0 = -0b1.0×2 ⁻⁵ ≈ -0.03125	0xf6 = 1.111011.0 = -0b1.0×2 ²⁷ = -134217728.0
0x37 = 0.011011.1 = +0b1.1×2 ⁻⁵ ≈ 0.046875	0x77 = 0.111011.1 = +0b1.1×2 ²⁷ = 201326592.0	0xb7 = 1.011011.1 = -0b1.1×2 ⁻⁵ ≈ -0.046875	0xf7 = 1.111011.1 = -0b1.1×2 ²⁷ = -201326592.0
0x38 = 0.011100.0 = +0b1.0×2 ⁻⁴ ≈ 0.0625	0x78 = 0.111100.0 = +0b1.0×2 ²⁸ = 268435456.0	0xb8 = 1.011100.0 = -0b1.0×2 ⁻⁴ ≈ -0.0625	0xf8 = 1.111100.0 = -0b1.0×2 ²⁸ = -268435456.0
0x39 = 0.011100.1 = +0b1.1×2 ⁻⁴ ≈ 0.09375	0x79 = 0.111100.1 = +0b1.1×2 ²⁸ = 402653184.0	0xb9 = 1.011100.1 = -0b1.1×2 ⁻⁴ ≈ -0.09375	0xf9 = 1.111100.1 = -0b1.1×2 ²⁸ = -402653184.0
0x3a = 0.011101.0 = +0b1.0×2 ⁻³ ≈ 0.125	0x7a = 0.111101.0 = +0b1.0×2 ²⁹ = 536870912.0	0xba = 1.011101.0 = -0b1.0×2 ⁻³ ≈ -0.125	0xfa = 1.111101.0 = -0b1.0×2 ²⁹ = -536870912.0
0x3b = 0.011101.1 = +0b1.1×2 ⁻³ ≈ 0.1875	0x7b = 0.111101.1 = +0b1.1×2 ²⁹ = 805306368.0	0xbb = 1.011101.1 = -0b1.1×2 ⁻³ ≈ -0.1875	0xfb = 1.111101.1 = -0b1.1×2 ²⁹ = -805306368.0
0x3c = 0.011110.0 = +0b1.0×2 ⁻² ≈ 0.25	0x7c = 0.111110.0 = +0b1.0×2 ³⁰ = 1073741824.0	0xbc = 1.011110.0 = -0b1.0×2 ⁻² ≈ -0.25	0xfc = 1.111110.0 = -0b1.0×2 ³⁰ = -1073741824.0
0x3d = 0.011110.1 = +0b1.1×2 ⁻² ≈ 0.375	0x7d = 0.111110.1 = +0b1.1×2 ³⁰ = 1610612736.0	0xbd = 1.011110.1 = -0b1.1×2 ⁻² ≈ -0.375	0xfd = 1.111110.1 = -0b1.1×2 ³⁰ = -1610612736.0
0x3e = 0.011111.0 = +0b1.0×2 ^{-1</}			

C.3 Value Table: P3, $\epsilon_{\min} = -15$, $\epsilon_{\max} = 15$

0x00 = 0.00000.00 = 0.0	0x40 = 0.10000.00 = +0b1.00×2 ⁰ = 1.0	0x80 = 1.00000.00 = NaN	0xc0 = 1.10000.00 = -0b1.00×2 ⁰ = -1.0
0x01 = 0.00000.01 = +0b0.01×2 ⁻¹⁵ ≈ 7.6293945E-06	0x41 = 0.10000.01 = +0b1.01×2 ⁰ = 1.25	0x81 = 1.00000.01 = -0b0.01×2 ⁻¹⁵ ≈ -7.6293945E-06	0xc1 = 1.10000.01 = -0b1.01×2 ⁰ = -1.25
0x02 = 0.00000.10 = +0b0.10×2 ⁻¹⁵ ≈ 1.5258789E-05	0x42 = 0.10000.10 = +0b1.10×2 ⁰ = 1.5	0x82 = 1.00000.10 = -0b0.10×2 ⁻¹⁵ ≈ -1.5258789E-05	0xc2 = 1.10000.10 = -0b1.10×2 ⁰ = -1.5
0x03 = 0.00000.11 = +0b0.11×2 ⁻¹⁵ ≈ 2.2888184E-05	0x43 = 0.10000.11 = +0b1.11×2 ⁰ = 1.75	0x83 = 1.00000.11 = -0b0.11×2 ⁻¹⁵ ≈ -2.2888184E-05	0xc3 = 1.10000.11 = -0b1.11×2 ⁰ = -1.75
0x04 = 0.00001.00 = +0b1.00×2 ⁻¹⁵ ≈ 3.0517578E-05	0x44 = 0.10001.00 = +0b1.00×2 ¹ = 2.0	0x84 = 1.00001.00 = -0b1.00×2 ⁻¹⁵ ≈ -3.0517578E-05	0xc4 = 1.10001.00 = -0b1.00×2 ¹ = -2.0
0x05 = 0.00001.01 = +0b1.01×2 ⁻¹⁵ ≈ 3.8146973E-05	0x45 = 0.10001.01 = +0b1.01×2 ¹ = 2.5	0x85 = 1.00001.01 = -0b1.01×2 ⁻¹⁵ ≈ -3.8146973E-05	0xc5 = 1.10001.01 = -0b1.01×2 ¹ = -2.5
0x06 = 0.00001.10 = +0b1.10×2 ⁻¹⁵ ≈ 4.5776367E-05	0x46 = 0.10001.10 = +0b1.10×2 ¹ = 3.0	0x86 = 1.00001.10 = -0b1.10×2 ⁻¹⁵ ≈ -4.5776367E-05	0xc6 = 1.10001.10 = -0b1.10×2 ¹ = -3.0
0x07 = 0.00001.11 = +0b1.11×2 ⁻¹⁵ ≈ 5.3405762E-05	0x47 = 0.10001.11 = +0b1.11×2 ¹ = 3.5	0x87 = 1.00001.11 = -0b1.11×2 ⁻¹⁵ ≈ -5.3405762E-05	0xc7 = 1.10001.11 = -0b1.11×2 ¹ = -3.5
0x08 = 0.00010.00 = +0b1.00×2 ⁻¹⁴ ≈ 6.1035156E-05	0x48 = 0.10010.00 = +0b1.00×2 ² = 4.0	0x88 = 1.00010.00 = -0b1.00×2 ⁻¹⁴ ≈ -6.1035156E-05	0xc8 = 1.10010.00 = -0b1.00×2 ² = -4.0
0x09 = 0.00010.01 = +0b1.01×2 ⁻¹⁴ ≈ 7.6293945E-05	0x49 = 0.10010.01 = +0b1.01×2 ² = 5.0	0x89 = 1.00010.01 = -0b1.01×2 ⁻¹⁴ ≈ -7.6293945E-05	0xc9 = 1.10010.01 = -0b1.01×2 ² = -5.0
0x0a = 0.00010.10 = +0b1.10×2 ⁻¹⁴ ≈ 9.1552734E-05	0x4a = 0.10010.10 = +0b1.10×2 ² = 6.0	0x8a = 1.00010.10 = -0b1.10×2 ⁻¹⁴ ≈ -9.1552734E-05	0xca = 1.10010.10 = -0b1.10×2 ² = -6.0
0x0b = 0.00010.11 = +0b1.11×2 ⁻¹⁴ ≈ 0.00010681152	0x4b = 0.10010.11 = +0b1.11×2 ² = 7.0	0x8b = 1.00010.11 = -0b1.11×2 ⁻¹⁴ ≈ -0.00010681152	0xcb = 1.10010.11 = -0b1.11×2 ² = -7.0
0x0c = 0.00011.00 = +0b1.00×2 ⁻¹³ ≈ 0.00012207031	0x4c = 0.10011.00 = +0b1.00×2 ³ = 8.0	0x8c = 1.00011.00 = -0b1.00×2 ⁻¹³ ≈ -0.00012207031	0xcc = 1.10011.00 = -0b1.00×2 ³ = -8.0
0x0d = 0.00011.01 = +0b1.01×2 ⁻¹³ ≈ 0.00015258789	0x4d = 0.10011.01 = +0b1.01×2 ³ = 10.0	0x8d = 1.00011.01 = -0b1.01×2 ⁻¹³ ≈ -0.00015258789	0xcd = 1.10011.01 = -0b1.01×2 ³ = -10.0
0x0e = 0.00011.10 = +0b1.10×2 ⁻¹³ ≈ 0.00018310547	0x4e = 0.10011.10 = +0b1.10×2 ³ = 12.0	0x8e = 1.00011.10 = -0b1.10×2 ⁻¹³ ≈ -0.00018310547	0xce = 1.10011.10 = -0b1.10×2 ³ = -12.0
0x0f = 0.00011.11 = +0b1.11×2 ⁻¹³ ≈ 0.00021362305	0x4f = 0.10011.11 = +0b1.11×2 ³ = 14.0	0x8f = 1.00011.11 = -0b1.11×2 ⁻¹³ ≈ -0.00021362305	0xcf = 1.10011.11 = -0b1.11×2 ³ = -14.0
0x10 = 0.00100.00 = +0b1.00×2 ⁻¹² = 0.000244140625	0x50 = 0.10100.00 = +0b1.00×2 ⁴ = 16.0	0x90 = 1.00100.00 = -0b1.00×2 ⁻¹² ≈ -0.000244140625	0xd0 = 1.10100.00 = -0b1.00×2 ⁴ = -16.0
0x11 = 0.00100.01 = +0b1.01×2 ⁻¹² ≈ 0.00030517578	0x51 = 0.10100.01 = +0b1.01×2 ⁴ = 20.0	0x91 = 1.00100.01 = -0b1.01×2 ⁻¹² ≈ -0.00030517578	0xd1 = 1.10100.01 = -0b1.01×2 ⁴ = -20.0
0x12 = 0.00100.10 = +0b1.10×2 ⁻¹² ≈ 0.00036621094	0x52 = 0.10100.10 = +0b1.10×2 ⁴ = 24.0	0x92 = 1.00100.10 = -0b1.10×2 ⁻¹² ≈ -0.00036621094	0xd2 = 1.10100.10 = -0b1.10×2 ⁴ = -24.0
0x13 = 0.00100.11 = +0b1.11×2 ⁻¹² ≈ 0.00042724609	0x53 = 0.10100.11 = +0b1.11×2 ⁴ = 28.0	0x93 = 1.00100.11 = -0b1.11×2 ⁻¹² ≈ -0.00042724609	0xd3 = 1.10100.11 = -0b1.11×2 ⁴ = -28.0
0x14 = 0.00101.00 = +0b1.00×2 ⁻¹¹ ≈ 0.00048828125	0x54 = 0.10101.00 = +0b1.00×2 ⁵ = 32.0	0x94 = 1.00101.00 = -0b1.00×2 ⁻¹¹ ≈ -0.00048828125	0xd4 = 1.10101.00 = -0b1.00×2 ⁵ = -32.0
0x15 = 0.00101.01 = +0b1.01×2 ⁻¹¹ ≈ 0.00061035156	0x55 = 0.10101.01 = +0b1.01×2 ⁵ = 40.0	0x95 = 1.00101.01 = -0b1.01×2 ⁻¹¹ ≈ -0.00061035156	0xd5 = 1.10101.01 = -0b1.01×2 ⁵ = -40.0
0x16 = 0.00101.10 = +0b1.10×2 ⁻¹¹ ≈ 0.000732421875	0x56 = 0.10101.10 = +0b1.10×2 ⁵ = 48.0	0x96 = 1.00101.10 = -0b1.10×2 ⁻¹¹ ≈ -0.000732421875	0xd6 = 1.10101.10 = -0b1.10×2 ⁵ = -48.0
0x17 = 0.00101.11 = +0b1.11×2 ⁻¹¹ ≈ 0.00085449219	0x57 = 0.10101.11 = +0b1.11×2 ⁵ = 56.0	0x97 = 1.00101.11 = -0b1.11×2 ⁻¹¹ ≈ -0.00085449219	0xd7 = 1.10101.11 = -0b1.11×2 ⁵ = -56.0
0x18 = 0.00110.00 = +0b1.00×2 ⁻¹⁰ = 0.0009765625	0x58 = 0.10110.00 = +0b1.00×2 ⁶ = 64.0	0x98 = 1.00110.00 = -0b1.00×2 ⁻¹⁰ = -0.0009765625	0xd8 = 1.10110.00 = -0b1.00×2 ⁶ = -64.0
0x19 = 0.00110.01 = +0b1.01×2 ⁻¹⁰ = 0.001220703125	0x59 = 0.10110.01 = +0b1.01×2 ⁶ = 80.0	0x99 = 1.00110.01 = -0b1.01×2 ⁻¹⁰ ≈ -0.0012207031	0xd9 = 1.10110.01 = -0b1.01×2 ⁶ = -80.0
0x1a = 0.00110.10 = +0b1.10×2 ⁻¹⁰ = 0.00146484375	0x5a = 0.10110.10 = +0b1.10×2 ⁶ = 96.0	0x9a = 1.00110.10 = -0b1.10×2 ⁻¹⁰ = -0.00146484375	0xda = 1.10110.10 = -0b1.10×2 ⁶ = -96.0
0x1b = 0.00110.11 = +0b1.11×2 ⁻¹⁰ = 0.001708984375	0x5b = 0.10110.11 = +0b1.11×2 ⁶ = 112.0	0x9b = 1.00110.11 = -0b1.11×2 ⁻¹⁰ ≈ -0.0017089844	0xdb = 1.10110.11 = -0b1.11×2 ⁶ = -112.0
0x1c = 0.00111.00 = +0b1.00×2 ⁻⁹ = 0.001953125	0x5c = 0.10111.00 = +0b1.00×2 ⁷ = 128.0	0x9c = 1.00111.00 = -0b1.00×2 ⁻⁹ = -0.001953125	0xdc = 1.10111.00 = -0b1.00×2 ⁷ = -128.0
0x1d = 0.00111.01 = +0b1.01×2 ⁻⁹ = 0.00244140625	0x5d = 0.10111.01 = +0b1.01×2 ⁷ = 160.0	0x9d = 1.00111.01 = -0b1.01×2 ⁻⁹ = -0.00244140625	0xdd = 1.10111.01 = -0b1.01×2 ⁷ = -160.0
0x1e = 0.00111.10 = +0b1.10×2 ⁻⁹ = 0.0029296875	0x5e = 0.10111.10 = +0b1.10×2 ⁷ = 192.0	0x9e = 1.00111.10 = -0b1.10×2 ⁻⁹ = -0.0029296875	0xde = 1.10111.10 = -0b1.10×2 ⁷ = -192.0
0x1f = 0.00111.11 = +0b1.11×2 ⁻⁹ = 0.00341796875	0x5f = 0.10111.11 = +0b1.11×2 ⁷ = 224.0	0x9f = 1.00111.11 = -0b1.11×2 ⁻⁹ = -0.00341796875	0xdf = 1.10111.11 = -0b1.11×2 ⁷ = -224.0
0x20 = 0.01000.00 = +0b1.00×2 ⁻⁸ = 0.00390625	0x60 = 0.11000.00 = +0b1.00×2 ⁸ = 256.0	0xa0 = 1.01000.00 = -0b1.00×2 ⁻⁸ = -0.00390625	0xe0 = 1.11000.00 = -0b1.00×2 ⁸ = -256.0
0x21 = 0.01000.01 = +0b1.01×2 ⁻⁸ = 0.0048828125	0x61 = 0.11000.01 = +0b1.01×2 ⁸ = 320.0	0xa1 = 1.01000.01 = -0b1.01×2 ⁻⁸ = -0.0048828125	0xe1 = 1.11000.01 = -0b1.01×2 ⁸ = -320.0
0x22 = 0.01000.10 = +0b1.10×2 ⁻⁸ = 0.005859375	0x62 = 0.11000.10 = +0b1.10×2 ⁸ = 384.0	0xa2 = 1.01000.10 = -0b1.10×2 ⁻⁸ = -0.005859375	0xe2 = 1.11000.10 = -0b1.10×2 ⁸ = -384.0
0x23 = 0.01000.11 = +0b1.11×2 ⁻⁸ = 0.0068359375	0x63 = 0.11000.11 = +0b1.11×2 ⁸ = 448.0	0xa3 = 1.01000.11 = -0b1.11×2 ⁻⁸ = -0.0068359375	0xe3 = 1.11000.11 = -0b1.11×2 ⁸ = -448.0
0x24 = 0.01001.00 = +0b1.00×2 ⁻⁷ = 0.0078125	0x64 = 0.11001.00 = +0b1.00×2 ⁹ = 512.0	0xa4 = 1.01001.00 = -0b1.00×2 ⁻⁷ = -0.0078125	0xe4 = 1.11001.00 = -0b1.00×2 ⁹ = -512.0
0x25 = 0.01001.01 = +0b1.01×2 ⁻⁷ = 0.009765625	0x65 = 0.11001.01 = +0b1.01×2 ⁹ = 640.0	0xa5 = 1.01001.01 = -0b1.01×2 ⁻⁷ = -0.009765625	0xe5 = 1.11001.01 = -0b1.01×2 ⁹ = -640.0
0x26 = 0.01001.10 = +0b1.10×2 ⁻⁷ = 0.01171875	0x66 = 0.11001.10 = +0b1.10×2 ⁹ = 768.0	0xa6 = 1.01001.10 = -0b1.10×2 ⁻⁷ = -0.01171875	0xe6 = 1.11001.10 = -0b1.10×2 ⁹ = -768.0
0x27 = 0.01001.11 = +0b1.11×2 ⁻⁷ = 0.013671875	0x67 = 0.11001.11 = +0b1.11×2 ⁹ = 896.0	0xa7 = 1.01001.11 = -0b1.11×2 ⁻⁷ = -0.013671875	0xe7 = 1.11001.11 = -0b1.11×2 ⁹ = -896.0
0x28 = 0.01010.00 = +0b1.00×2 ⁻⁶ = 0.015625	0x68 = 0.11010.00 = +0b1.00×2 ¹⁰ = 1024.0	0xa8 = 1.01010.00 = -0b1.00×2 ⁻⁶ = -0.015625	0xe8 = 1.11010.00 = -0b1.00×2 ¹⁰ = -1024.0
0x29 = 0.01010.01 = +0b1.01×2 ⁻⁶ = 0.01953125	0x69 = 0.11010.01 = +0b1.01×2 ¹⁰ = 1280.0	0xa9 = 1.01010.01 = -0b1.01×2 ⁻⁶ = -0.01953125	0xe9 = 1.11010.01 = -0b1.01×2 ¹⁰ = -1280.0
0x2a = 0.01010.10 = +0b1.10×2 ⁻⁶ = 0.0234375	0x6a = 0.11010.10 = +0b1.10×2 ¹⁰ = 1536.0	0xaa = 1.01010.10 = -0b1.10×2 ⁻⁶ = -0.0234375	0xea = 1.11010.10 = -0b1.10×2 ¹⁰ = -1536.0
0x2b = 0.01010.11 = +0b1.11×2 ⁻⁶ = 0.02734375	0x6b = 0.11010.11 = +0b1.11×2 ¹⁰ = 1792.0	0xab = 1.01010.11 = -0b1.11×2 ⁻⁶ = -0.02734375	0xeb = 1.11010.11 = -0b1.11×2 ¹⁰ = -1792.0
0x2c = 0.01011.00 = +0b1.00×2 ⁻⁵ = 0.03125	0x6c = 0.11011.00 = +0b1.00×2 ¹¹ = 2048.0	0xac = 1.01011.00 = -0b1.00×2 ⁻⁵ = -0.03125	0xec = 1.11011.00 = -0b1.00×2 ¹¹ = -2048.0
0x2d = 0.01011.01 = +0b1.01×2 ⁻⁵ = 0.0390625	0x6d = 0.11011.01 = +0b1.01×2 ¹¹ = 2560.0	0xad = 1.01011.01 = -0b1.01×2 ⁻⁵ = -0.0390625	0xed = 1.11011.01 = -0b1.01×2 ¹¹ = -2560.0
0x2e = 0.01011.10 = +0b1.10×2 ⁻⁵ = 0.046875	0x6e = 0.11011.10 = +0b1.10×2 ¹¹ = 3072.0	0xae = 1.01011.10 = -0b1.10×2 ⁻⁵ = -0.046875	0xee = 1.11011.10 = -0b1.10×2 ¹¹ = -3072.0
0x2f = 0.01011.11 = +0b1.11×2 ⁻⁵ = 0.0546875	0x6f = 0.11011.11 = +0b1.11×2 ¹¹ = 3584.0	0xaf = 1.01011.11 = -0b1.11×2 ⁻⁵ = -0.0546875	0xef = 1.11011.11 = -0b1.11×2 ¹¹ = -3584.0
0x30 = 0.01100.00 = +0b1.00×2 ⁻⁴ = 0.0625	0x70 = 0.11100.00 = +0b1.00×2 ¹² = 4096.0	0xb0 = 1.01100.00 = -0b1.00×2 ⁻⁴ = -0.0625	0xf0 = 1.11100.00 = -0b1.00×2 ¹² = -4096.0
0x31 = 0.01100.01 = +0b1.01×2 ⁻⁴ = 0.078125	0x71 = 0.11100.01 = +0b1.01×2 ¹² = 5120.0	0xb1 = 1.01100.01 = -0b1.01×2 ⁻⁴ = -0.078125	0xf1 = 1.11100.01 = -0b1.01×2 ¹² = -5120.0
0x32 = 0.01100.10 = +0b1.10×2 ⁻⁴ = 0.09375	0x72 = 0.11100.10 = +0b1.10×2 ¹² = 6144.0	0xb2 = 1.01100.10 = -0b1.10×2 ⁻⁴ = -0.09375	0xf2 = 1.11100.10 = -0b1.10×2 ¹² = -6144.0
0x33 = 0.01100.11 = +0b1.11×2 ⁻⁴ = 0.109375	0x73 = 0.11100.11 = +0b1.11×2 ¹² = 7168.0	0xb3 = 1.01100.11 = -0b1.11×2 ⁻⁴ = -0.109375	0xf3 = 1.11100.11 = -0b1.11×2 ¹² = -7168.0
0x34 = 0.01101.00 = +0b1.00×2 ⁻³ = 0.125	0x74 = 0.11101.00 = +0b1.00×2 ¹³ = 8192.0	0xb4 = 1.01101.00 = -0b1.00×2 ⁻³ = -0.125	0xf4 = 1.11101.00 = -0b1.00×2 ¹³ = -8192.0
0x35 = 0.01101.01 = +0b1.01×2 ⁻³ = 0.15625	0x75 = 0.11101.01 = +0b1.01×2 ¹³ = 10240.0	0xb5 = 1.01101.01 = -0b1.01×2 ⁻³ = -0.15625	0xf5 = 1.11101.01 = -0b1.01×2 ¹³ = -10240.0
0x36 = 0.01101.10 = +0b1.10×2 ⁻³ = 0.1875	0x76 = 0.11101.10 = +0b1.10×2 ¹³ = 12288.0	0xb6 = 1.01101.10 = -0b1.10×2 ⁻³ = -0.1875	0xf6 = 1.11101.10 = -0b1.10×2 ¹³ = -12288.0
0x37 = 0.01101.11 = +0b1.11×2 ⁻³ = 0.21875	0x77 = 0.11101.11 = +0b1.11×2 ¹³ = 14336.0	0xb7 = 1.01101.11 = -0b1.11×2 ⁻³ = -0.21875	0xf7 = 1.11101.11 = -0b1.11×2 ¹³ = -14336.0
0x38 = 0.01110.00 = +0b1.00×2 ⁻² = 0.25	0x78 = 0.11110.00 = +0b1.00×2 ¹⁴ = 16384.0	0xb8 = 1.01110.00 = -0b1.00×2 ⁻² = -0.25	0xf8 = 1.11110.00 = -0b1.00×2 ¹⁴ = -16384.0
0x39 = 0.01110.01 = +0b1.01×2 ⁻² = 0.3125	0x79 = 0.11110.01 = +0b1.01×2 ¹⁴ = 20480.0	0xb9 = 1.01110.01 = -0b1.01×2 ⁻² = -0.3125	0xf9 = 1.11110.01 = -0b1.01×2 ¹⁴ = -20480.0
0x3a = 0.01110.10 = +0b1.10×2 ⁻² = 0.375	0x7a = 0.11110.10 = +0b1.10×2 ¹⁴ = 24576.0	0xba = 1.01110.10 = -0b1.10×2 ⁻² = -0.375	0xfa = 1.11110.10 = -0b1.10×2 ¹⁴ = -24576.0
0x3b = 0.01110.11 = +0b1.11×2 ⁻² = 0.4375	0x7b = 0.11110.11 = +0b1.11×2 ¹⁴ = 28672.0	0xbb = 1.01110.11 = -0b1.11×2 ⁻² = -0.4375	0xfb = 1.11110.11 = -0b1.11×2 ¹⁴ = -28672.0
0x3c = 0.01111.00 = +0b1.00×2 ⁻¹ = 0.5	0x7c = 0.11111.00 = +0b1.00×2 ¹⁵ = 32768.0	0xbc = 1.01111.00 = -0b1.00×2 ⁻¹ = -0.5	0xfc = 1.11111.00 = -0b1.00×2 ¹⁵ = -32768.0
0x3d = 0.01111.01 = +0b1.01×2 ⁻¹ = 0.625	0x7d = 0.11111.01 = +0b1.01×2 ¹⁵ = 40960.0	0xbd = 1.01111.01 = -0b1.01×2 ⁻¹ = -0.625	0xfd = 1.11111.01 = -0b1.01×2 ¹⁵ = -40960.0
0x3e = 0.01111.10 = +0b1.10×2 ⁻¹ = 0.75	0x7e = 0.11111.10 = +0b1.10×2 ¹⁵ = 49152.0	0xbe = 1.01111.10 = -0b1.10×2 ⁻¹ = -0.75	0xfe = 1.11111.10 = -0b1.10×2 ¹⁵ = -49152.0
0x3f = 0.01111.11 = +0b1.11×2 ⁻¹ = 0.875	0x7f = 0.11111.11 = +Inf	0xbf = 1.01111.11 = -0b1.11×2 ⁻¹ = -	

C.4 Value Table: P4, emin = -7, emax = 7

0x00 = 0.0000.000 = 0.0	0x40 = 0.1000.000 = +0b1.000×2 ⁷⁰ = 1.0	0x80 = 1.0000.000 = NaN	0xc0 = 1.1000.000 = -0b1.000×2 ⁷⁰ = -1.0
0x01 = 0.0000.001 = +0b0.001×2 ⁻⁷ = 0.0009765625	0x41 = 0.1000.001 = +0b1.001×2 ⁷⁰ = 1.125	0x81 = 1.0000.001 = -0b0.001×2 ⁻⁷ = -0.0009765625	0xc1 = 1.1000.001 = -0b1.001×2 ⁷⁰ = -1.125
0x02 = 0.0000.010 = +0b0.010×2 ⁻⁷ = 0.001953125	0x42 = 0.1000.010 = +0b1.010×2 ⁷⁰ = 1.25	0x82 = 1.0000.010 = -0b0.010×2 ⁻⁷ = -0.001953125	0xc2 = 1.1000.010 = -0b1.010×2 ⁷⁰ = -1.25
0x03 = 0.0000.011 = +0b0.011×2 ⁻⁷ = 0.0029296875	0x43 = 0.1000.011 = +0b1.011×2 ⁷⁰ = 1.375	0x83 = 1.0000.011 = -0b0.011×2 ⁻⁷ = -0.0029296875	0xc3 = 1.1000.011 = -0b1.011×2 ⁷⁰ = -1.375
0x04 = 0.0000.100 = +0b0.100×2 ⁻⁷ = 0.00390625	0x44 = 0.1000.100 = +0b1.100×2 ⁷⁰ = 1.5	0x84 = 1.0000.100 = -0b0.100×2 ⁻⁷ = -0.00390625	0xc4 = 1.1000.100 = -0b1.100×2 ⁷⁰ = -1.5
0x05 = 0.0000.101 = +0b0.101×2 ⁻⁷ = 0.0048828125	0x45 = 0.1000.101 = +0b1.101×2 ⁷⁰ = 1.625	0x85 = 1.0000.101 = -0b0.101×2 ⁻⁷ = -0.0048828125	0xc5 = 1.1000.101 = -0b1.101×2 ⁷⁰ = -1.625
0x06 = 0.0000.110 = +0b0.110×2 ⁻⁷ = 0.005859375	0x46 = 0.1000.110 = +0b1.110×2 ⁷⁰ = 1.75	0x86 = 1.0000.110 = -0b0.110×2 ⁻⁷ = -0.005859375	0xc6 = 1.1000.110 = -0b1.110×2 ⁷⁰ = -1.75
0x07 = 0.0000.111 = +0b0.111×2 ⁻⁷ = 0.0068359375	0x47 = 0.1000.111 = +0b1.111×2 ⁷⁰ = 1.875	0x87 = 1.0000.111 = -0b0.111×2 ⁻⁷ = -0.0068359375	0xc7 = 1.1000.111 = -0b1.111×2 ⁷⁰ = -1.875
0x08 = 0.0001.000 = +0b1.000×2 ⁻⁶ = 0.0078125	0x48 = 0.1001.000 = +0b1.000×2 ⁷¹ = 2.0	0x88 = 1.0001.000 = -0b1.000×2 ⁻⁶ = -0.0078125	0xc8 = 1.1001.000 = -0b1.000×2 ⁷¹ = -2.0
0x09 = 0.0001.001 = +0b1.001×2 ⁻⁶ = 0.0087890625	0x49 = 0.1001.001 = +0b1.001×2 ⁷¹ = 2.25	0x89 = 1.0001.001 = -0b1.001×2 ⁻⁶ = -0.0087890625	0xc9 = 1.1001.001 = -0b1.001×2 ⁷¹ = -2.25
0x0a = 0.0001.010 = +0b1.010×2 ⁻⁶ = 0.009765625	0x4a = 0.1001.010 = +0b1.010×2 ⁷¹ = 2.5	0x8a = 1.0001.010 = -0b1.010×2 ⁻⁶ = -0.009765625	0xca = 1.1001.010 = -0b1.010×2 ⁷¹ = -2.5
0x0b = 0.0001.011 = +0b1.011×2 ⁻⁶ = 0.0107421875	0x4b = 0.1001.011 = +0b1.011×2 ⁷¹ = 2.75	0x8b = 1.0001.011 = -0b1.011×2 ⁻⁶ = -0.0107421875	0xcb = 1.1001.011 = -0b1.011×2 ⁷¹ = -2.75
0x0c = 0.0001.100 = +0b1.100×2 ⁻⁶ = 0.01171875	0x4c = 0.1001.100 = +0b1.100×2 ⁷¹ = 3.0	0x8c = 1.0001.100 = -0b1.100×2 ⁻⁶ = -0.01171875	0xcc = 1.1001.100 = -0b1.100×2 ⁷¹ = -3.0
0x0d = 0.0001.101 = +0b1.101×2 ⁻⁶ = 0.0126953125	0x4d = 0.1001.101 = +0b1.101×2 ⁷¹ = 3.25	0x8d = 1.0001.101 = -0b1.101×2 ⁻⁶ = -0.0126953125	0xcd = 1.1001.101 = -0b1.101×2 ⁷¹ = -3.25
0x0e = 0.0001.110 = +0b1.110×2 ⁻⁶ = 0.013671875	0x4e = 0.1001.110 = +0b1.110×2 ⁷¹ = 3.5	0x8e = 1.0001.110 = -0b1.110×2 ⁻⁶ = -0.013671875	0xce = 1.1001.110 = -0b1.110×2 ⁷¹ = -3.5
0x0f = 0.0001.111 = +0b1.111×2 ⁻⁶ = 0.0146484375	0x4f = 0.1001.111 = +0b1.111×2 ⁷¹ = 3.75	0x8f = 1.0001.111 = -0b1.111×2 ⁻⁶ = -0.0146484375	0xcf = 1.1001.111 = -0b1.111×2 ⁷¹ = -3.75
0x10 = 0.0010.000 = +0b1.000×2 ⁻⁵ = 0.015625	0x50 = 0.1010.000 = +0b1.000×2 ⁷² = 4.0	0x90 = 1.0010.000 = -0b1.000×2 ⁻⁵ = -0.015625	0xd0 = 1.1010.000 = -0b1.000×2 ⁷² = -4.0
0x11 = 0.0010.001 = +0b1.001×2 ⁻⁵ = 0.017578125	0x51 = 0.1010.001 = +0b1.001×2 ⁷² = 4.5	0x91 = 1.0010.001 = -0b1.001×2 ⁻⁵ = -0.017578125	0xd1 = 1.1010.001 = -0b1.001×2 ⁷² = -4.5
0x12 = 0.0010.010 = +0b1.010×2 ⁻⁵ = 0.01953125	0x52 = 0.1010.010 = +0b1.010×2 ⁷² = 5.0	0x92 = 1.0010.010 = -0b1.010×2 ⁻⁵ = -0.01953125	0xd2 = 1.1010.010 = -0b1.010×2 ⁷² = -5.0
0x13 = 0.0010.011 = +0b1.011×2 ⁻⁵ = 0.021484375	0x53 = 0.1010.011 = +0b1.011×2 ⁷² = 5.5	0x93 = 1.0010.011 = -0b1.011×2 ⁻⁵ = -0.021484375	0xd3 = 1.1010.011 = -0b1.011×2 ⁷² = -5.5
0x14 = 0.0010.100 = +0b1.100×2 ⁻⁵ = 0.0234375	0x54 = 0.1010.100 = +0b1.100×2 ⁷² = 6.0	0x94 = 1.0010.100 = -0b1.100×2 ⁻⁵ = -0.0234375	0xd4 = 1.1010.100 = -0b1.100×2 ⁷² = -6.0
0x15 = 0.0010.101 = +0b1.101×2 ⁻⁵ = 0.025390625	0x55 = 0.1010.101 = +0b1.101×2 ⁷² = 6.5	0x95 = 1.0010.101 = -0b1.101×2 ⁻⁵ = -0.025390625	0xd5 = 1.1010.101 = -0b1.101×2 ⁷² = -6.5
0x16 = 0.0010.110 = +0b1.110×2 ⁻⁵ = 0.02734375	0x56 = 0.1010.110 = +0b1.110×2 ⁷² = 7.0	0x96 = 1.0010.110 = -0b1.110×2 ⁻⁵ = -0.02734375	0xd6 = 1.1010.110 = -0b1.110×2 ⁷² = -7.0
0x17 = 0.0010.111 = +0b1.111×2 ⁻⁵ = 0.029296875	0x57 = 0.1010.111 = +0b1.111×2 ⁷² = 7.5	0x97 = 1.0010.111 = -0b1.111×2 ⁻⁵ = -0.029296875	0xd7 = 1.1010.111 = -0b1.111×2 ⁷² = -7.5
0x18 = 0.0011.000 = +0b1.000×2 ⁻⁴ = 0.03125	0x58 = 0.1011.000 = +0b1.000×2 ⁷³ = 8.0	0x98 = 1.0011.000 = -0b1.000×2 ⁻⁴ = -0.03125	0xd8 = 1.1011.000 = -0b1.000×2 ⁷³ = -8.0
0x19 = 0.0011.001 = +0b1.001×2 ⁻⁴ = 0.03515625	0x59 = 0.1011.001 = +0b1.001×2 ⁷³ = 9.0	0x99 = 1.0011.001 = -0b1.001×2 ⁻⁴ = -0.03515625	0xd9 = 1.1011.001 = -0b1.001×2 ⁷³ = -9.0
0x1a = 0.0011.010 = +0b1.010×2 ⁻⁴ = 0.0390625	0x5a = 0.1011.010 = +0b1.010×2 ⁷³ = 10.0	0x9a = 1.0011.010 = -0b1.010×2 ⁻⁴ = -0.0390625	0xda = 1.1011.010 = -0b1.010×2 ⁷³ = -10.0
0x1b = 0.0011.011 = +0b1.011×2 ⁻⁴ = 0.04296875	0x5b = 0.1011.011 = +0b1.011×2 ⁷³ = 11.0	0x9b = 1.0011.011 = -0b1.011×2 ⁻⁴ = -0.04296875	0xdb = 1.1011.011 = -0b1.011×2 ⁷³ = -11.0
0x1c = 0.0011.100 = +0b1.100×2 ⁻⁴ = 0.046875	0x5c = 0.1011.100 = +0b1.100×2 ⁷³ = 12.0	0x9c = 1.0011.100 = -0b1.100×2 ⁻⁴ = -0.046875	0xdc = 1.1011.100 = -0b1.100×2 ⁷³ = -12.0
0x1d = 0.0011.101 = +0b1.101×2 ⁻⁴ = 0.05078125	0x5d = 0.1011.101 = +0b1.101×2 ⁷³ = 13.0	0x9d = 1.0011.101 = -0b1.101×2 ⁻⁴ = -0.05078125	0xdd = 1.1011.101 = -0b1.101×2 ⁷³ = -13.0
0x1e = 0.0011.110 = +0b1.110×2 ⁻⁴ = 0.0546875	0x5e = 0.1011.110 = +0b1.110×2 ⁷³ = 14.0	0x9e = 1.0011.110 = -0b1.110×2 ⁻⁴ = -0.0546875	0xde = 1.1011.110 = -0b1.110×2 ⁷³ = -14.0
0x1f = 0.0011.111 = +0b1.111×2 ⁻⁴ = 0.05859375	0x5f = 0.1011.111 = +0b1.111×2 ⁷³ = 15.0	0x9f = 1.0011.111 = -0b1.111×2 ⁻⁴ = -0.05859375	0xdf = 1.1011.111 = -0b1.111×2 ⁷³ = -15.0
0x20 = 0.0100.000 = +0b1.000×2 ⁻³ = 0.0625	0x60 = 0.1100.000 = +0b1.000×2 ⁷⁴ = 16.0	0xa0 = 1.0100.000 = -0b1.000×2 ⁻³ = -0.0625	0xe0 = 1.1100.000 = -0b1.000×2 ⁷⁴ = -16.0
0x21 = 0.0100.001 = +0b1.001×2 ⁻³ = 0.0703125	0x61 = 0.1100.001 = +0b1.001×2 ⁷⁴ = 18.0	0xa1 = 1.0100.001 = -0b1.001×2 ⁻³ = -0.0703125	0xe1 = 1.1100.001 = -0b1.001×2 ⁷⁴ = -18.0
0x22 = 0.0100.010 = +0b1.010×2 ⁻³ = 0.078125	0x62 = 0.1100.010 = +0b1.010×2 ⁷⁴ = 20.0	0xa2 = 1.0100.010 = -0b1.010×2 ⁻³ = -0.078125	0xe2 = 1.1100.010 = -0b1.010×2 ⁷⁴ = -20.0
0x23 = 0.0100.011 = +0b1.011×2 ⁻³ = 0.0859375	0x63 = 0.1100.011 = +0b1.011×2 ⁷⁴ = 22.0	0xa3 = 1.0100.011 = -0b1.011×2 ⁻³ = -0.0859375	0xe3 = 1.1100.011 = -0b1.011×2 ⁷⁴ = -22.0
0x24 = 0.0100.100 = +0b1.100×2 ⁻³ = 0.09375	0x64 = 0.1100.100 = +0b1.100×2 ⁷⁴ = 24.0	0xa4 = 1.0100.100 = -0b1.100×2 ⁻³ = -0.09375	0xe4 = 1.1100.100 = -0b1.100×2 ⁷⁴ = -24.0
0x25 = 0.0100.101 = +0b1.101×2 ⁻³ = 0.1015625	0x65 = 0.1100.101 = +0b1.101×2 ⁷⁴ = 26.0	0xa5 = 1.0100.101 = -0b1.101×2 ⁻³ = -0.1015625	0xe5 = 1.1100.101 = -0b1.101×2 ⁷⁴ = -26.0
0x26 = 0.0100.110 = +0b1.110×2 ⁻³ = 0.109375	0x66 = 0.1100.110 = +0b1.110×2 ⁷⁴ = 28.0	0xa6 = 1.0100.110 = -0b1.110×2 ⁻³ = -0.109375	0xe6 = 1.1100.110 = -0b1.110×2 ⁷⁴ = -28.0
0x27 = 0.0100.111 = +0b1.111×2 ⁻³ = 0.1171875	0x67 = 0.1100.111 = +0b1.111×2 ⁷⁴ = 30.0	0xa7 = 1.0100.111 = -0b1.111×2 ⁻³ = -0.1171875	0xe7 = 1.1100.111 = -0b1.111×2 ⁷⁴ = -30.0
0x28 = 0.0101.000 = +0b1.000×2 ⁻² = 0.125	0x68 = 0.1101.000 = +0b1.000×2 ⁷⁵ = 32.0	0xa8 = 1.0101.000 = -0b1.000×2 ⁻² = -0.125	0xe8 = 1.1101.000 = -0b1.000×2 ⁷⁵ = -32.0
0x29 = 0.0101.001 = +0b1.001×2 ⁻² = 0.140625	0x69 = 0.1101.001 = +0b1.001×2 ⁷⁵ = 36.0	0xa9 = 1.0101.001 = -0b1.001×2 ⁻² = -0.140625	0xe9 = 1.1101.001 = -0b1.001×2 ⁷⁵ = -36.0
0x2a = 0.0101.010 = +0b1.010×2 ⁻² = 0.15625	0x6a = 0.1101.010 = +0b1.010×2 ⁷⁵ = 40.0	0xaa = 1.0101.010 = -0b1.010×2 ⁻² = -0.15625	0xea = 1.1101.010 = -0b1.010×2 ⁷⁵ = -40.0
0x2b = 0.0101.011 = +0b1.011×2 ⁻² = 0.171875	0x6b = 0.1101.011 = +0b1.011×2 ⁷⁵ = 44.0	0xab = 1.0101.011 = -0b1.011×2 ⁻² = -0.171875	0xeb = 1.1101.011 = -0b1.011×2 ⁷⁵ = -44.0
0x2c = 0.0101.100 = +0b1.100×2 ⁻² = 0.1875	0x6c = 0.1101.100 = +0b1.100×2 ⁷⁵ = 48.0	0xac = 1.0101.100 = -0b1.100×2 ⁻² = -0.1875	0xec = 1.1101.100 = -0b1.100×2 ⁷⁵ = -48.0
0x2d = 0.0101.101 = +0b1.101×2 ⁻² = 0.203125	0x6d = 0.1101.101 = +0b1.101×2 ⁷⁵ = 52.0	0xad = 1.0101.101 = -0b1.101×2 ⁻² = -0.203125	0xed = 1.1101.101 = -0b1.101×2 ⁷⁵ = -52.0
0x2e = 0.0101.110 = +0b1.110×2 ⁻² = 0.21875	0x6e = 0.1101.110 = +0b1.110×2 ⁷⁵ = 56.0	0xae = 1.0101.110 = -0b1.110×2 ⁻² = -0.21875	0xee = 1.1101.110 = -0b1.110×2 ⁷⁵ = -56.0
0x2f = 0.0101.111 = +0b1.111×2 ⁻² = 0.234375	0x6f = 0.1101.111 = +0b1.111×2 ⁷⁵ = 60.0	0xaf = 1.0101.111 = -0b1.111×2 ⁻² = -0.234375	0xef = 1.1101.111 = -0b1.111×2 ⁷⁵ = -60.0
0x30 = 0.0110.000 = +0b1.000×2 ⁻¹ = 0.25	0x70 = 0.1110.000 = +0b1.000×2 ⁷⁶ = 64.0	0xb0 = 1.0110.000 = -0b1.000×2 ⁻¹ = -0.25	0xf0 = 1.1110.000 = -0b1.000×2 ⁷⁶ = -64.0
0x31 = 0.0110.001 = +0b1.001×2 ⁻¹ = 0.28125	0x71 = 0.1110.001 = +0b1.001×2 ⁷⁶ = 72.0	0xb1 = 1.0110.001 = -0b1.001×2 ⁻¹ = -0.28125	0xf1 = 1.1110.001 = -0b1.001×2 ⁷⁶ = -72.0
0x32 = 0.0110.010 = +0b1.010×2 ⁻¹ = 0.3125	0x72 = 0.1110.010 = +0b1.010×2 ⁷⁶ = 80.0	0xb2 = 1.0110.010 = -0b1.010×2 ⁻¹ = -0.3125	0xf2 = 1.1110.010 = -0b1.010×2 ⁷⁶ = -80.0
0x33 = 0.0110.011 = +0b1.011×2 ⁻¹ = 0.34375	0x73 = 0.1110.011 = +0b1.011×2 ⁷⁶ = 88.0	0xb3 = 1.0110.011 = -0b1.011×2 ⁻¹ = -0.34375	0xf3 = 1.1110.011 = -0b1.011×2 ⁷⁶ = -88.0
0x34 = 0.0110.100 = +0b1.100×2 ⁻¹ = 0.375	0x74 = 0.1110.100 = +0b1.100×2 ⁷⁶ = 96.0	0xb4 = 1.0110.100 = -0b1.100×2 ⁻¹ = -0.375	0xf4 = 1.1110.100 = -0b1.100×2 ⁷⁶ = -96.0
0x35 = 0.0110.101 = +0b1.101×2 ⁻¹ = 0.40625	0x75 = 0.1110.101 = +0b1.101×2 ⁷⁶ = 104.0	0xb5 = 1.0110.101 = -0b1.101×2 ⁻¹ = -0.40625	0xf5 = 1.1110.101 = -0b1.101×2 ⁷⁶ = -104.0
0x36 = 0.0110.110 = +0b1.110×2 ⁻¹ = 0.4375	0x76 = 0.1110.110 = +0b1.110×2 ⁷⁶ = 112.0	0xb6 = 1.0110.110 = -0b1.110×2 ⁻¹ = -0.4375	0xf6 = 1.1110.110 = -0b1.110×2 ⁷⁶ = -112.0
0x37 = 0.0110.111 = +0b1.111×2 ⁻¹ = 0.46875	0x77 = 0.1110.111 = +0b1.111×2 ⁷⁶ = 120.0	0xb7 = 1.0110.111 = -0b1.111×2 ⁻¹ = -0.46875	0xf7 = 1.1110.111 = -0b1.111×2 ⁷⁶ = -120.0
0x38 = 0.0111.000 = +0b1.000×2 ⁰ = 0.5	0x78 = 0.1111.000 = +0b1.000×2 ⁷⁷ = 128.0	0xb8 = 1.0111.000 = -0b1.000×2 ⁰ = -0.5	0xf8 = 1.1111.000 = -0b1.000×2 ⁷⁷ = -128.0
0x39 = 0.0111.001 = +0b1.001×2 ⁰ = 0.5625	0x79 = 0.1111.001 = +0b1.001×2 ⁷⁷ = 144.0	0xb9 = 1.0111.001 = -0b1.001×2 ⁰ = -0.5625	0xf9 = 1.1111.001 = -0b1.001×2 ⁷⁷ = -144.0
0x3a = 0.0111.010 = +0b1.010×2 ⁰ = 0.625	0x7a = 0.1111.010 = +0b1.010×2 ⁷⁷ = 160.0	0xba = 1.0111.010 = -0b1.010×2 ⁰ = -0.625	0xfa = 1.1111.010 = -0b1.010×2 ⁷⁷ = -160.0
0x3b = 0.0111.011 = +0b1.011×2 ⁰ = 0.6875	0x7b = 0.1111.011 = +0b1.011×2 ⁷⁷ = 176.0	0xbb = 1.0111.011 = -0b1.011×2 ⁰ = -0.6875	0xfb = 1.1111.011 = -0b1.011×2 ⁷⁷ = -176.0
0x3c = 0.0111.100 = +0b1.100×2 ⁰ = 0.75	0x7c = 0.1111.100 = +0b1.100×2 ⁷⁷ = 192.0	0xbc = 1.0111.100 = -0b1.100×2 ⁰ = -0.75	0xfc = 1.1111.100 = -0b1.100×2 ⁷⁷ = -192.0
0x3d = 0.0111.101 = +0b1.101×2 ⁰ = 0.8125	0x7d = 0.1111.101 = +0b1.101×2 ⁷⁷ = 208.0	0xbd = 1.0111.101 = -0b1.101×2 ⁰ = -0.8125	0xfd = 1.1111.101 = -0b1.101×2 ⁷⁷ = -208.0
0x3e = 0.0111.110 = +0b1.110×2 ⁰ = 0.875	0x7e = 0.1111.110 = +0b1.110×2 ⁷⁷ = 224.0	0xbe = 1.0111.110 = -0b1.110×2 ⁰ = -0.875	0xfe = 1.1111.110 = -0b1.110×2 ⁷⁷ = -224.0
0x3f = 0.0111.111 = +0b1.111×2 ⁰ = 0.9375	0x7f = 0.1111.111 = +Inf	0xbf = 1.0111.111 = -0b1.111×2 ⁰ = -0.9375	0xff = 1.1111.111 = -Inf

C.5 Value Table: P5, emin = -3, emax = 3

0x00 = 0.000.0000 = 0.0	0x40 = 0.100.0000 = +0b1.0000×2 ⁰ = 1.0	0x80 = 1.000.0000 = NaN	0xc0 = 1.100.0000 = -0b1.0000×2 ⁰ = -1.0
0x01 = 0.000.0001 = +0b0.0001×2 ⁻³ = 0.0078125	0x41 = 0.100.0001 = +0b1.0001×2 ⁰ = 1.0625	0x81 = 1.000.0001 = -0b0.0001×2 ⁻³ = -0.0078125	0xc1 = 1.100.0001 = -0b1.0001×2 ⁰ = -1.0625
0x02 = 0.000.0010 = +0b0.0010×2 ⁻³ = 0.015625	0x42 = 0.100.0010 = +0b1.0010×2 ⁰ = 1.125	0x82 = 1.000.0010 = -0b0.0010×2 ⁻³ = -0.015625	0xc2 = 1.100.0010 = -0b1.0010×2 ⁰ = -1.125
0x03 = 0.000.0011 = +0b0.0011×2 ⁻³ = 0.0234375	0x43 = 0.100.0011 = +0b1.0011×2 ⁰ = 1.1875	0x83 = 1.000.0011 = -0b0.0011×2 ⁻³ = -0.0234375	0xc3 = 1.100.0011 = -0b1.0011×2 ⁰ = -1.1875
0x04 = 0.000.0100 = +0b0.0100×2 ⁻³ = 0.03125	0x44 = 0.100.0100 = +0b1.0100×2 ⁰ = 1.25	0x84 = 1.000.0100 = -0b0.0100×2 ⁻³ = -0.03125	0xc4 = 1.100.0100 = -0b1.0100×2 ⁰ = -1.25
0x05 = 0.000.0101 = +0b0.0101×2 ⁻³ = 0.0390625	0x45 = 0.100.0101 = +0b1.0101×2 ⁰ = 1.3125	0x85 = 1.000.0101 = -0b0.0101×2 ⁻³ = -0.0390625	0xc5 = 1.100.0101 = -0b1.0101×2 ⁰ = -1.3125
0x06 = 0.000.0110 = +0b0.0110×2 ⁻³ = 0.046875	0x46 = 0.100.0110 = +0b1.0110×2 ⁰ = 1.375	0x86 = 1.000.0110 = -0b0.0110×2 ⁻³ = -0.046875	0xc6 = 1.100.0110 = -0b1.0110×2 ⁰ = -1.375
0x07 = 0.000.0111 = +0b0.0111×2 ⁻³ = 0.0546875	0x47 = 0.100.0111 = +0b1.0111×2 ⁰ = 1.4375	0x87 = 1.000.0111 = -0b0.0111×2 ⁻³ = -0.0546875	0xc7 = 1.100.0111 = -0b1.0111×2 ⁰ = -1.4375
0x08 = 0.000.1000 = +0b0.1011×2 ⁻³ = 0.0625	0x48 = 0.100.1000 = +0b1.1000×2 ⁰ = 1.5	0x88 = 1.000.1000 = -0b0.1000×2 ⁻³ = -0.0625	0xc8 = 1.100.1000 = -0b1.1000×2 ⁰ = -1.5
0x09 = 0.000.1001 = +0b0.1001×2 ⁻³ = 0.0703125	0x49 = 0.100.1001 = +0b1.1001×2 ⁰ = 1.5625	0x89 = 1.000.1001 = -0b0.1001×2 ⁻³ = -0.0703125	0xc9 = 1.100.1001 = -0b1.1001×2 ⁰ = -1.5625
0x0a = 0.000.1010 = +0b0.1010×2 ⁻³ = 0.078125	0x4a = 0.100.1010 = +0b1.1010×2 ⁰ = 1.625	0x8a = 1.000.1010 = -0b0.1010×2 ⁻³ = -0.078125	0xca = 1.100.1010 = -0b1.1010×2 ⁰ = -1.625
0x0b = 0.000.1011 = +0b0.1011×2 ⁻³ = 0.0859375	0x4b = 0.100.1011 = +0b1.1011×2 ⁰ = 1.6875	0x8b = 1.000.1011 = -0b0.1011×2 ⁻³ = -0.0859375	0xcb = 1.100.1011 = -0b1.1011×2 ⁰ = -1.6875
0x0c = 0.000.1100 = +0b0.1100×2 ⁻³ = 0.09375	0x4c = 0.100.1100 = +0b1.1100×2 ⁰ = 1.75	0x8c = 1.000.1100 = -0b0.1100×2 ⁻³ = -0.09375	0xcc = 1.100.1100 = -0b1.1100×2 ⁰ = -1.75
0x0d = 0.000.1101 = +0b0.1101×2 ⁻³ = 0.1015625	0x4d = 0.100.1101 = +0b1.1101×2 ⁰ = 1.8125	0x8d = 1.000.1101 = -0b0.1101×2 ⁻³ = -0.1015625	0xcd = 1.100.1101 = -0b1.1101×2 ⁰ = -1.8125
0x0e = 0.000.1110 = +0b0.1110×2 ⁻³ = 0.109375	0x4e = 0.100.1110 = +0b1.1110×2 ⁰ = 1.875	0x8e = 1.000.1110 = -0b0.1110×2 ⁻³ = -0.109375	0xce = 1.100.1110 = -0b1.1110×2 ⁰ = -1.875
0x0f = 0.000.1111 = +0b0.1111×2 ⁻³ = 0.1171875	0x4f = 0.100.1111 = +0b1.1111×2 ⁰ = 1.9375	0x8f = 1.000.1111 = -0b0.1111×2 ⁻³ = -0.1171875	0xcf = 1.100.1111 = -0b1.1111×2 ⁰ = -1.9375
0x10 = 0.001.0000 = +0b1.0000×2 ⁻³ = 0.125	0x50 = 0.101.0000 = +0b1.0000×2 ¹ = 2.0	0x90 = 1.001.0000 = -0b1.0000×2 ⁻³ = -0.125	0xd0 = 1.101.0000 = -0b1.0000×2 ¹ = -2.0
0x11 = 0.001.0001 = +0b1.0001×2 ⁻³ = 0.1328125	0x51 = 0.101.0001 = +0b1.0001×2 ¹ = 2.125	0x91 = 1.001.0001 = -0b1.0001×2 ⁻³ = -0.1328125	0xd1 = 1.101.0001 = -0b1.0001×2 ¹ = -2.125
0x12 = 0.001.0010 = +0b1.0010×2 ⁻³ = 0.140625	0x52 = 0.101.0010 = +0b1.0010×2 ¹ = 2.25	0x92 = 1.001.0010 = -0b1.0010×2 ⁻³ = -0.140625	0xd2 = 1.101.0010 = -0b1.0010×2 ¹ = -2.25
0x13 = 0.001.0011 = +0b1.0011×2 ⁻³ = 0.1484375	0x53 = 0.101.0011 = +0b1.0011×2 ¹ = 2.375	0x93 = 1.001.0011 = -0b1.0011×2 ⁻³ = -0.1484375	0xd3 = 1.101.0011 = -0b1.0011×2 ¹ = -2.375
0x14 = 0.001.0100 = +0b1.0100×2 ⁻³ = 0.15625	0x54 = 0.101.0100 = +0b1.0100×2 ¹ = 2.5	0x94 = 1.001.0100 = -0b1.0100×2 ⁻³ = -0.15625	0xd4 = 1.101.0100 = -0b1.0100×2 ¹ = -2.5
0x15 = 0.001.0101 = +0b1.0101×2 ⁻³ = 0.1640625	0x55 = 0.101.0101 = +0b1.0101×2 ¹ = 2.625	0x95 = 1.001.0101 = -0b1.0101×2 ⁻³ = -0.1640625	0xd5 = 1.101.0101 = -0b1.0101×2 ¹ = -2.625
0x16 = 0.001.0110 = +0b1.0110×2 ⁻³ = 0.171875	0x56 = 0.101.0110 = +0b1.0110×2 ¹ = 2.75	0x96 = 1.001.0110 = -0b1.0110×2 ⁻³ = -0.171875	0xd6 = 1.101.0110 = -0b1.0110×2 ¹ = -2.75
0x17 = 0.001.0111 = +0b1.0111×2 ⁻³ = 0.1796875	0x57 = 0.101.0111 = +0b1.0111×2 ¹ = 2.875	0x97 = 1.001.0111 = -0b1.0111×2 ⁻³ = -0.1796875	0xd7 = 1.101.0111 = -0b1.0111×2 ¹ = -2.875
0x18 = 0.001.1000 = +0b1.1000×2 ⁻³ = 0.1875	0x58 = 0.101.1000 = +0b1.1000×2 ¹ = 3.0	0x98 = 1.001.1000 = -0b1.1000×2 ⁻³ = -0.1875	0xd8 = 1.101.1000 = -0b1.1000×2 ¹ = -3.0
0x19 = 0.001.1001 = +0b1.1001×2 ⁻³ = 0.1953125	0x59 = 0.101.1001 = +0b1.1001×2 ¹ = 3.125	0x99 = 1.001.1001 = -0b1.1001×2 ⁻³ = -0.1953125	0xd9 = 1.101.1001 = -0b1.1001×2 ¹ = -3.125
0x1a = 0.001.1010 = +0b1.1010×2 ⁻³ = 0.203125	0x5a = 0.101.1010 = +0b1.1010×2 ¹ = 3.25	0x9a = 1.001.1010 = -0b1.1010×2 ⁻³ = -0.203125	0xda = 1.101.1010 = -0b1.1010×2 ¹ = -3.25
0x1b = 0.001.1011 = +0b1.1011×2 ⁻³ = 0.2109375	0x5b = 0.101.1011 = +0b1.1011×2 ¹ = 3.375	0x9b = 1.001.1011 = -0b1.1011×2 ⁻³ = -0.2109375	0xdb = 1.101.1011 = -0b1.1011×2 ¹ = -3.375
0x1c = 0.001.1100 = +0b1.1100×2 ⁻³ = 0.21875	0x5c = 0.101.1100 = +0b1.1100×2 ¹ = 3.5	0x9c = 1.001.1100 = -0b1.1100×2 ⁻³ = -0.21875	0xdc = 1.101.1100 = -0b1.1100×2 ¹ = -3.5
0x1d = 0.001.1101 = +0b1.1101×2 ⁻³ = 0.2265625	0x5d = 0.101.1101 = +0b1.1101×2 ¹ = 3.625	0x9d = 1.001.1101 = -0b1.1101×2 ⁻³ = -0.2265625	0xdd = 1.101.1101 = -0b1.1101×2 ¹ = -3.625
0x1e = 0.001.1110 = +0b1.1110×2 ⁻³ = 0.234375	0x5e = 0.101.1110 = +0b1.1110×2 ¹ = 3.75	0x9e = 1.001.1110 = -0b1.1110×2 ⁻³ = -0.234375	0xde = 1.101.1110 = -0b1.1110×2 ¹ = -3.75
0x1f = 0.001.1111 = +0b1.1111×2 ⁻³ = 0.2421875	0x5f = 0.101.1111 = +0b1.1111×2 ¹ = 3.875	0x9f = 1.001.1111 = -0b1.1111×2 ⁻³ = -0.2421875	0xdf = 1.101.1111 = -0b1.1111×2 ¹ = -3.875
0x20 = 0.010.0000 = +0b1.0000×2 ⁻² = 0.25	0x60 = 0.110.0000 = +0b1.0000×2 ² = 4.0	0xa0 = 1.010.0000 = -0b1.0000×2 ⁻² = -0.25	0xe0 = 1.110.0000 = -0b1.0000×2 ² = -4.0
0x21 = 0.010.0001 = +0b1.0001×2 ⁻² = 0.265625	0x61 = 0.110.0001 = +0b1.0001×2 ² = 4.25	0xa1 = 1.010.0001 = -0b1.0001×2 ⁻² = -0.265625	0xe1 = 1.110.0001 = -0b1.0001×2 ² = -4.25
0x22 = 0.010.0010 = +0b1.0010×2 ⁻² = 0.28125	0x62 = 0.110.0010 = +0b1.0010×2 ² = 4.5	0xa2 = 1.010.0010 = -0b1.0010×2 ⁻² = -0.28125	0xe2 = 1.110.0010 = -0b1.0010×2 ² = -4.5
0x23 = 0.010.0011 = +0b1.0011×2 ⁻² = 0.296875	0x63 = 0.110.0011 = +0b1.0011×2 ² = 4.75	0xa3 = 1.010.0011 = -0b1.0011×2 ⁻² = -0.296875	0xe3 = 1.110.0011 = -0b1.0011×2 ² = -4.75
0x24 = 0.010.0100 = +0b1.0100×2 ⁻² = 0.3125	0x64 = 0.110.0100 = +0b1.0100×2 ² = 5.0	0xa4 = 1.010.0100 = -0b1.0100×2 ⁻² = -0.3125	0xe4 = 1.110.0100 = -0b1.0100×2 ² = -5.0
0x25 = 0.010.0101 = +0b1.0101×2 ⁻² = 0.328125	0x65 = 0.110.0101 = +0b1.0101×2 ² = 5.25	0xa5 = 1.010.0101 = -0b1.0101×2 ⁻² = -0.328125	0xe5 = 1.110.0101 = -0b1.0101×2 ² = -5.25
0x26 = 0.010.0110 = +0b1.0110×2 ⁻² = 0.34375	0x66 = 0.110.0110 = +0b1.0110×2 ² = 5.5	0xa6 = 1.010.0110 = -0b1.0110×2 ⁻² = -0.34375	0xe6 = 1.110.0110 = -0b1.0110×2 ² = -5.5
0x27 = 0.010.0111 = +0b1.0111×2 ⁻² = 0.359375	0x67 = 0.110.0111 = +0b1.0111×2 ² = 5.75	0xa7 = 1.010.0111 = -0b1.0111×2 ⁻² = -0.359375	0xe7 = 1.110.0111 = -0b1.0111×2 ² = -5.75
0x28 = 0.010.1000 = +0b1.1000×2 ⁻² = 0.375	0x68 = 0.110.1000 = +0b1.1000×2 ² = 6.0	0xa8 = 1.010.1000 = -0b1.1000×2 ⁻² = -0.375	0xe8 = 1.110.1000 = -0b1.1000×2 ² = -6.0
0x29 = 0.010.1001 = +0b1.1001×2 ⁻² = 0.390625	0x69 = 0.110.1001 = +0b1.1001×2 ² = 6.25	0xa9 = 1.010.1001 = -0b1.1001×2 ⁻² = -0.390625	0xe9 = 1.110.1001 = -0b1.1001×2 ² = -6.25
0x2a = 0.010.1010 = +0b1.1010×2 ⁻² = 0.40625	0x6a = 0.110.1010 = +0b1.1010×2 ² = 6.5	0xaa = 1.010.1010 = -0b1.1010×2 ⁻² = -0.40625	0xea = 1.110.1010 = -0b1.1010×2 ² = -6.5
0x2b = 0.010.1011 = +0b1.1011×2 ⁻² = 0.421875	0x6b = 0.110.1011 = +0b1.1011×2 ² = 6.75	0xab = 1.010.1011 = -0b1.1011×2 ⁻² = -0.421875	0xeb = 1.110.1011 = -0b1.1011×2 ² = -6.75
0x2c = 0.010.1100 = +0b1.1100×2 ⁻² = 0.4375	0x6c = 0.110.1100 = +0b1.1100×2 ² = 7.0	0xac = 1.010.1100 = -0b1.1100×2 ⁻² = -0.4375	0xec = 1.110.1100 = -0b1.1100×2 ² = -7.0
0x2d = 0.010.1101 = +0b1.1101×2 ⁻² = 0.453125	0x6d = 0.110.1101 = +0b1.1101×2 ² = 7.25	0xad = 1.010.1101 = -0b1.1101×2 ⁻² = -0.453125	0xed = 1.110.1101 = -0b1.1101×2 ² = -7.25
0x2e = 0.010.1110 = +0b1.1110×2 ⁻² = 0.46875	0x6e = 0.110.1110 = +0b1.1110×2 ² = 7.5	0xae = 1.010.1110 = -0b1.1110×2 ⁻² = -0.46875	0xee = 1.110.1110 = -0b1.1110×2 ² = -7.5
0x2f = 0.010.1111 = +0b1.1111×2 ⁻² = 0.484375	0x6f = 0.110.1111 = +0b1.1111×2 ² = 7.75	0xaf = 1.010.1111 = -0b1.1111×2 ⁻² = -0.484375	0xef = 1.110.1111 = -0b1.1111×2 ² = -7.75
0x30 = 0.011.0000 = +0b1.0000×2 ⁻¹ = 0.5	0x70 = 0.111.0000 = +0b1.0000×2 ³ = 8.0	0xb0 = 1.011.0000 = -0b1.0000×2 ⁻¹ = -0.5	0xf0 = 1.111.0000 = -0b1.0000×2 ³ = -8.0
0x31 = 0.011.0001 = +0b1.0001×2 ⁻¹ = 0.53125	0x71 = 0.111.0001 = +0b1.0001×2 ³ = 8.5	0xb1 = 1.011.0001 = -0b1.0001×2 ⁻¹ = -0.53125	0xf1 = 1.111.0001 = -0b1.0001×2 ³ = -8.5
0x32 = 0.011.0010 = +0b1.0010×2 ⁻¹ = 0.5625	0x72 = 0.111.0010 = +0b1.0010×2 ³ = 9.0	0xb2 = 1.011.0010 = -0b1.0010×2 ⁻¹ = -0.5625	0xf2 = 1.111.0010 = -0b1.0010×2 ³ = -9.0
0x33 = 0.011.0011 = +0b1.0011×2 ⁻¹ = 0.59375	0x73 = 0.111.0011 = +0b1.0011×2 ³ = 9.5	0xb3 = 1.011.0011 = -0b1.0011×2 ⁻¹ = -0.59375	0xf3 = 1.111.0011 = -0b1.0011×2 ³ = -9.5
0x34 = 0.011.0100 = +0b1.0100×2 ⁻¹ = 0.625	0x74 = 0.111.0100 = +0b1.0100×2 ³ = 10.0	0xb4 = 1.011.0100 = -0b1.0100×2 ⁻¹ = -0.625	0xf4 = 1.111.0100 = -0b1.0100×2 ³ = -10.0
0x35 = 0.011.0101 = +0b1.0101×2 ⁻¹ = 0.65625	0x75 = 0.111.0101 = +0b1.0101×2 ³ = 10.5	0xb5 = 1.011.0101 = -0b1.0101×2 ⁻¹ = -0.65625	0xf5 = 1.111.0101 = -0b1.0101×2 ³ = -10.5
0x36 = 0.011.0110 = +0b1.0110×2 ⁻¹ = 0.6875	0x76 = 0.111.0110 = +0b1.0110×2 ³ = 11.0	0xb6 = 1.011.0110 = -0b1.0110×2 ⁻¹ = -0.6875	0xf6 = 1.111.0110 = -0b1.0110×2 ³ = -11.0
0x37 = 0.011.0111 = +0b1.0111×2 ⁻¹ = 0.71875	0x77 = 0.111.0111 = +0b1.0111×2 ³ = 11.5	0xb7 = 1.011.0111 = -0b1.0111×2 ⁻¹ = -0.71875	0xf7 = 1.111.0111 = -0b1.0111×2 ³ = -11.5
0x38 = 0.011.1000 = +0b1.1000×2 ⁻¹ = 0.75	0x78 = 0.111.1000 = +0b1.1000×2 ³ = 12.0	0xb8 = 1.011.1000 = -0b1.1000×2 ⁻¹ = -0.75	0xf8 = 1.111.1000 = -0b1.1000×2 ³ = -12.0
0x39 = 0.011.1001 = +0b1.1001×2 ⁻¹ = 0.78125	0x79 = 0.111.1001 = +0b1.1001×2 ³ = 12.5	0xb9 = 1.011.1001 = -0b1.1001×2 ⁻¹ = -0.78125	0xf9 = 1.111.1001 = -0b1.1001×2 ³ = -12.5
0x3a = 0.011.1010 = +0b1.1010×2 ⁻¹ = 0.8125	0x7a = 0.111.1010 = +0b1.1010×2 ³ = 13.0	0xba = 1.011.1010 = -0b1.1010×2 ⁻¹ = -0.8125	0xfa = 1.111.1010 = -0b1.1010×2 ³ = -13.0
0x3b = 0.011.1011 = +0b1.1011×2 ⁻¹ = 0.84375	0x7b = 0.111.1011 = +0b1.1011×2 ³ = 13.5	0xbb = 1.011.1011 = -0b1.1011×2 ⁻¹ = -0.84375	0xfb = 1.111.1011 = -0b1.1011×2 ³ = -13.5
0x3c = 0.011.1100 = +0b1.1100×2 ⁻¹ = 0.875	0x7c = 0.111.1100 = +0b1.1100×2 ³ = 14.0	0xbc = 1.011.1100 = -0b1.1100×2 ⁻¹ = -0.875	0xfc = 1.111.1100 = -0b1.1100×2 ³ = -14.0
0x3d = 0.011.1101 = +0b1.1101×2 ⁻¹ = 0.90625	0x7d = 0.111.1101 = +0b1.1101×2 ³ = 14.5	0xbd = 1.011.1101 = -0b1.1101×2 ⁻¹ = -0.90625	0xfd = 1.111.1101 = -0b1.1101×2 ³ = -14.5
0x3e = 0.011.1110 = +0b1.1110×2 ⁻¹ = 0.9375	0x7e = 0.111.1110 = +0b1.1110×2 ³ = 15.0	0xbe = 1.011.1110 = -0b1.1110×2 ⁻¹ = -0.9375	0xfe = 1.111.1110 = -0b1.1110×2 ³ = -15.0
0x3f = 0.011.1111 = +0b1.1111×2 ⁻¹ = 0.96875	0x7f = 0.111.1111 = +Inf	0xbf = 1.011.1111 = -0b1.1111×2 ⁻¹ = -0.96875	0xff = 1.111.1111 = -Inf

C.6 Value Table: P6, emin = -1, emax = 1

```
0x00 = 0.00.00000 = 0.0
0x01 = 0.00.00001 = +0b0.00001×2-1 = 0.015625
0x02 = 0.00.00010 = +0b0.00010×2-1 = 0.03125
0x03 = 0.00.00011 = +0b0.00011×2-1 = 0.046875
0x04 = 0.00.00100 = +0b0.00100×2-1 = 0.0625
0x05 = 0.00.00101 = +0b0.00101×2-1 = 0.078125
0x06 = 0.00.00110 = +0b0.00110×2-1 = 0.09375
0x07 = 0.00.00111 = +0b0.00111×2-1 = 0.109375
0x08 = 0.00.01000 = +0b0.01000×2-1 = 0.125
0x09 = 0.00.01001 = +0b0.01001×2-1 = 0.140625
0x0a = 0.00.01010 = +0b0.01010×2-1 = 0.15625
0x0b = 0.00.01011 = +0b0.01011×2-1 = 0.171875
0x0c = 0.00.01100 = +0b0.01100×2-1 = 0.1875
0x0d = 0.00.01101 = +0b0.01101×2-1 = 0.203125
0x0e = 0.00.01110 = +0b0.01110×2-1 = 0.21875
0x0f = 0.00.01111 = +0b0.01111×2-1 = 0.234375
0x10 = 0.00.10000 = +0b0.10000×2-1 = 0.25
0x11 = 0.00.10001 = +0b0.10001×2-1 = 0.265625
0x12 = 0.00.10010 = +0b0.10010×2-1 = 0.28125
0x13 = 0.00.10011 = +0b0.10011×2-1 = 0.296875
0x14 = 0.00.10100 = +0b0.10100×2-1 = 0.3125
0x15 = 0.00.10101 = +0b0.10101×2-1 = 0.328125
0x16 = 0.00.10110 = +0b0.10110×2-1 = 0.34375
0x17 = 0.00.10111 = +0b0.10111×2-1 = 0.359375
0x18 = 0.00.11000 = +0b0.11000×2-1 = 0.375
0x19 = 0.00.11001 = +0b0.11001×2-1 = 0.390625
0x1a = 0.00.11010 = +0b0.11010×2-1 = 0.40625
0x1b = 0.00.11011 = +0b0.11011×2-1 = 0.421875
0x1c = 0.00.11100 = +0b0.11100×2-1 = 0.4375
0x1d = 0.00.11101 = +0b0.11101×2-1 = 0.453125
0x1e = 0.00.11110 = +0b0.11110×2-1 = 0.46875
0x1f = 0.00.11111 = +0b0.11111×2-1 = 0.484375
0x20 = 0.01.00000 = +0b1.00000×2-1 = 0.5
0x21 = 0.01.00001 = +0b1.00001×2-1 = 0.515625
0x22 = 0.01.00010 = +0b1.00010×2-1 = 0.53125
0x23 = 0.01.00011 = +0b1.00011×2-1 = 0.546875
0x24 = 0.01.00100 = +0b1.00100×2-1 = 0.5625
0x25 = 0.01.00101 = +0b1.00101×2-1 = 0.578125
0x26 = 0.01.00110 = +0b1.00110×2-1 = 0.59375
0x27 = 0.01.00111 = +0b1.00111×2-1 = 0.609375
0x28 = 0.01.01000 = +0b1.01000×2-1 = 0.625
0x29 = 0.01.01001 = +0b1.01001×2-1 = 0.640625
0x2a = 0.01.01010 = +0b1.01010×2-1 = 0.65625
0x2b = 0.01.01011 = +0b1.01011×2-1 = 0.671875
0x2c = 0.01.01100 = +0b1.01100×2-1 = 0.6875
0x2d = 0.01.01101 = +0b1.01101×2-1 = 0.703125
0x2e = 0.01.01110 = +0b1.01110×2-1 = 0.71875
0x2f = 0.01.01111 = +0b1.01111×2-1 = 0.734375
0x30 = 0.01.10000 = +0b1.10000×2-1 = 0.75
0x31 = 0.01.10001 = +0b1.10001×2-1 = 0.765625
0x32 = 0.01.10010 = +0b1.10010×2-1 = 0.78125
0x33 = 0.01.10011 = +0b1.10011×2-1 = 0.796875
0x34 = 0.01.10100 = +0b1.10100×2-1 = 0.8125
0x35 = 0.01.10101 = +0b1.10101×2-1 = 0.828125
0x36 = 0.01.10110 = +0b1.10110×2-1 = 0.84375
0x37 = 0.01.10111 = +0b1.10111×2-1 = 0.859375
0x38 = 0.01.11000 = +0b1.11000×2-1 = 0.875
0x39 = 0.01.11001 = +0b1.11001×2-1 = 0.890625
0x3a = 0.01.11010 = +0b1.11010×2-1 = 0.90625
0x3b = 0.01.11011 = +0b1.11011×2-1 = 0.921875
0x3c = 0.01.11100 = +0b1.11100×2-1 = 0.9375
0x3d = 0.01.11101 = +0b1.11101×2-1 = 0.953125
0x3e = 0.01.11110 = +0b1.11110×2-1 = 0.96875
0x3f = 0.01.11111 = +0b1.11111×2-1 = 0.984375

0x40 = 0.10.00000 = +0b1.00000×20 = 1.0
0x41 = 0.10.00001 = +0b1.00001×20 = 1.03125
0x42 = 0.10.00010 = +0b1.00010×20 = 1.0625
0x43 = 0.10.00011 = +0b1.00011×20 = 1.09375
0x44 = 0.10.00100 = +0b1.00100×20 = 1.125
0x45 = 0.10.00101 = +0b1.00101×20 = 1.15625
0x46 = 0.10.00110 = +0b1.00110×20 = 1.1875
0x47 = 0.10.00111 = +0b1.00111×20 = 1.21875
0x48 = 0.10.01000 = +0b1.01000×20 = 1.25
0x49 = 0.10.01001 = +0b1.01001×20 = 1.28125
0x4a = 0.10.01010 = +0b1.01010×20 = 1.3125
0x4b = 0.10.01011 = +0b1.01011×20 = 1.34375
0x4c = 0.10.01100 = +0b1.01100×20 = 1.375
0x4d = 0.10.01101 = +0b1.01101×20 = 1.40625
0x4e = 0.10.01110 = +0b1.01110×20 = 1.4375
0x4f = 0.10.01111 = +0b1.01111×20 = 1.46875
0x50 = 0.10.10000 = +0b1.10000×20 = 1.5
0x51 = 0.10.10001 = +0b1.10001×20 = 1.53125
0x52 = 0.10.10010 = +0b1.10010×20 = 1.5625
0x53 = 0.10.10011 = +0b1.10011×20 = 1.59375
0x54 = 0.10.10100 = +0b1.10100×20 = 1.625
0x55 = 0.10.10101 = +0b1.10101×20 = 1.65625
0x56 = 0.10.10110 = +0b1.10110×20 = 1.6875
0x57 = 0.10.10111 = +0b1.10111×20 = 1.71875
0x58 = 0.10.11000 = +0b1.11000×20 = 1.75
0x59 = 0.10.11001 = +0b1.11001×20 = 1.78125
0x5a = 0.10.11010 = +0b1.11010×20 = 1.8125
0x5b = 0.10.11011 = +0b1.11011×20 = 1.84375
0x5c = 0.10.11100 = +0b1.11100×20 = 1.875
0x5d = 0.10.11101 = +0b1.11101×20 = 1.90625
0x5e = 0.10.11110 = +0b1.11110×20 = 1.9375
0x5f = 0.10.11111 = +0b1.11111×20 = 1.96875
0x60 = 0.11.00000 = +0b1.00000×21 = 2.0
0x61 = 0.11.00001 = +0b1.00001×21 = 2.0625
0x62 = 0.11.00010 = +0b1.00010×21 = 2.125
0x63 = 0.11.00011 = +0b1.00011×21 = 2.1875
0x64 = 0.11.00100 = +0b1.00100×21 = 2.25
0x65 = 0.11.00101 = +0b1.00101×21 = 2.3125
0x66 = 0.11.00110 = +0b1.00110×21 = 2.375
0x67 = 0.11.00111 = +0b1.00111×21 = 2.4375
0x68 = 0.11.01000 = +0b1.01000×21 = 2.5
0x69 = 0.11.01001 = +0b1.01001×21 = 2.5625
0x6a = 0.11.01010 = +0b1.01010×21 = 2.625
0x6b = 0.11.01011 = +0b1.01011×21 = 2.6875
0x6c = 0.11.01100 = +0b1.01100×21 = 2.75
0x6d = 0.11.01101 = +0b1.01101×21 = 2.8125
0x6e = 0.11.01110 = +0b1.01110×21 = 2.875
0x6f = 0.11.01111 = +0b1.01111×21 = 2.9375
0x70 = 0.11.10000 = +0b1.10000×21 = 3.0
0x71 = 0.11.10001 = +0b1.10001×21 = 3.0625
0x72 = 0.11.10010 = +0b1.10010×21 = 3.125
0x73 = 0.11.10011 = +0b1.10011×21 = 3.1875
0x74 = 0.11.10100 = +0b1.10100×21 = 3.25
0x75 = 0.11.10101 = +0b1.10101×21 = 3.3125
0x76 = 0.11.10110 = +0b1.10110×21 = 3.375
0x77 = 0.11.10111 = +0b1.10111×21 = 3.4375
0x78 = 0.11.11000 = +0b1.11000×21 = 3.5
0x79 = 0.11.11001 = +0b1.11001×21 = 3.5625
0x7a = 0.11.11010 = +0b1.11010×21 = 3.625
0x7b = 0.11.11011 = +0b1.11011×21 = 3.6875
0x7c = 0.11.11100 = +0b1.11100×21 = 3.75
0x7d = 0.11.11101 = +0b1.11101×21 = 3.8125
0x7e = 0.11.11110 = +0b1.11110×21 = 3.875
0x7f = 0.11.11111 = +Inf

0x80 = 1.00.00000 = NaN
0x81 = 1.00.00001 = -0b0.00001×2-1 = -0.015625
0x82 = 1.00.00010 = -0b0.00010×2-1 = -0.03125
0x83 = 1.00.00011 = -0b0.00011×2-1 = -0.046875
0x84 = 1.00.00100 = -0b0.00100×2-1 = -0.0625
0x85 = 1.00.00101 = -0b0.00101×2-1 = -0.078125
0x86 = 1.00.00110 = -0b0.00110×2-1 = -0.09375
0x87 = 1.00.00111 = -0b0.00111×2-1 = -0.109375
0x88 = 1.00.01000 = -0b0.01000×2-1 = -0.125
0x89 = 1.00.01001 = -0b0.01001×2-1 = -0.140625
0x8a = 1.00.01010 = -0b0.01010×2-1 = -0.15625
0x8b = 1.00.01011 = -0b0.01011×2-1 = -0.171875
0x8c = 1.00.01100 = -0b0.01100×2-1 = -0.1875
0x8d = 1.00.01101 = -0b0.01101×2-1 = -0.203125
0x8e = 1.00.01110 = -0b0.01110×2-1 = -0.21875
0x8f = 1.00.01111 = -0b0.01111×2-1 = -0.234375
0x90 = 1.00.10000 = -0b0.10000×2-1 = -0.25
0x91 = 1.00.10001 = -0b0.10001×2-1 = -0.265625
0x92 = 1.00.10010 = -0b0.10010×2-1 = -0.28125
0x93 = 1.00.10011 = -0b0.10011×2-1 = -0.296875
0x94 = 1.00.10100 = -0b0.10100×2-1 = -0.3125
0x95 = 1.00.10101 = -0b0.10101×2-1 = -0.328125
0x96 = 1.00.10110 = -0b0.10110×2-1 = -0.34375
0x97 = 1.00.10111 = -0b0.10111×2-1 = -0.359375
0x98 = 1.00.11000 = -0b0.11000×2-1 = -0.375
0x99 = 1.00.11001 = -0b0.11001×2-1 = -0.390625
0x9a = 1.00.11010 = -0b0.11010×2-1 = -0.40625
0x9b = 1.00.11011 = -0b0.11011×2-1 = -0.421875
0x9c = 1.00.11100 = -0b0.11100×2-1 = -0.4375
0x9d = 1.00.11101 = -0b0.11101×2-1 = -0.453125
0x9e = 1.00.11110 = -0b0.11110×2-1 = -0.46875
0x9f = 1.00.11111 = -0b0.11111×2-1 = -0.484375
0xa0 = 1.01.00000 = -0b1.00000×2-1 = -0.5
0xa1 = 1.01.00001 = -0b1.00001×2-1 = -0.515625
0xa2 = 1.01.00010 = -0b1.00010×2-1 = -0.53125
0xa3 = 1.01.00011 = -0b1.00011×2-1 = -0.546875
0xa4 = 1.01.00100 = -0b1.00100×2-1 = -0.5625
0xa5 = 1.01.00101 = -0b1.00101×2-1 = -0.578125
0xa6 = 1.01.00110 = -0b1.00110×2-1 = -0.59375
0xa7 = 1.01.00111 = -0b1.00111×2-1 = -0.609375
0xa8 = 1.01.01000 = -0b1.01000×2-1 = -0.625
0xa9 = 1.01.01001 = -0b1.01001×2-1 = -0.640625
0xaa = 1.01.01010 = -0b1.01010×2-1 = -0.65625
0xab = 1.01.01011 = -0b1.01011×2-1 = -0.671875
0xac = 1.01.01100 = -0b1.01100×2-1 = -0.6875
0xad = 1.01.01101 = -0b1.01101×2-1 = -0.703125
0xae = 1.01.01110 = -0b1.01110×2-1 = -0.71875
0xaf = 1.01.01111 = -0b1.01111×2-1 = -0.734375
0xb0 = 1.01.10000 = -0b1.10000×2-1 = -0.75
0xb1 = 1.01.10001 = -0b1.10001×2-1 = -0.765625
0xb2 = 1.01.10010 = -0b1.10010×2-1 = -0.78125
0xb3 = 1.01.10011 = -0b1.10011×2-1 = -0.796875
0xb4 = 1.01.10100 = -0b1.10100×2-1 = -0.8125
0xb5 = 1.01.10101 = -0b1.10101×2-1 = -0.828125
0xb6 = 1.01.10110 = -0b1.10110×2-1 = -0.84375
0xb7 = 1.01.10111 = -0b1.10111×2-1 = -0.859375
0xb8 = 1.01.11000 = -0b1.11000×2-1 = -0.875
0xb9 = 1.01.11001 = -0b1.11001×2-1 = -0.890625
0xba = 1.01.11010 = -0b1.11010×2-1 = -0.90625
0xbb = 1.01.11011 = -0b1.11011×2-1 = -0.921875
0xbc = 1.01.11100 = -0b1.11100×2-1 = -0.9375
0xbd = 1.01.11101 = -0b1.11101×2-1 = -0.953125
0xbe = 1.01.11110 = -0b1.11110×2-1 = -0.96875
0xbf = 1.01.11111 = -0b1.11111×2-1 = -0.984375

0xc0 = 1.10.00000 = -0b1.00000×20 = -1.0
0xc1 = 1.10.00001 = -0b1.00001×20 = -1.03125
0xc2 = 1.10.00010 = -0b1.00010×20 = -1.0625
0xc3 = 1.10.00011 = -0b1.00011×20 = -1.09375
0xc4 = 1.10.00100 = -0b1.00100×20 = -1.125
0xc5 = 1.10.00101 = -0b1.00101×20 = -1.15625
0xc6 = 1.10.00110 = -0b1.00110×20 = -1.1875
0xc7 = 1.10.00111 = -0b1.00111×20 = -1.21875
0xc8 = 1.10.01000 = -0b1.01000×20 = -1.25
0xc9 = 1.10.01001 = -0b1.01001×20 = -1.28125
0xca = 1.10.01010 = -0b1.01010×20 = -1.3125
0xcb = 1.10.01011 = -0b1.01011×20 = -1.34375
0xcc = 1.10.01100 = -0b1.01100×20 = -1.375
0xcd = 1.10.01101 = -0b1.01101×20 = -1.40625
0xce = 1.10.01110 = -0b1.01110×20 = -1.4375
0xcf = 1.10.01111 = -0b1.01111×20 = -1.46875
0xd0 = 1.10.10000 = -0b1.10000×20 = -1.5
0xd1 = 1.10.10001 = -0b1.10001×20 = -1.53125
0xd2 = 1.10.10010 = -0b1.10010×20 = -1.5625
0xd3 = 1.10.10011 = -0b1.10011×20 = -1.59375
0xd4 = 1.10.10100 = -0b1.10100×20 = -1.625
0xd5 = 1.10.10101 = -0b1.10101×20 = -1.65625
0xd6 = 1.10.10110 = -0b1.10110×20 = -1.6875
0xd7 = 1.10.10111 = -0b1.10111×20 = -1.71875
0xd8 = 1.10.11000 = -0b1.11000×20 = -1.75
0xd9 = 1.10.11001 = -0b1.11001×20 = -1.78125
0xda = 1.10.11010 = -0b1.11010×20 = -1.8125
0xdb = 1.10.11011 = -0b1.11011×20 = -1.84375
0xdc = 1.10.11100 = -0b1.11100×20 = -1.875
0xdd = 1.10.11101 = -0b1.11101×20 = -1.90625
0xde = 1.10.11110 = -0b1.11110×20 = -1.9375
0xdf = 1.10.11111 = -0b1.11111×20 = -1.96875
0xe0 = 1.11.00000 = -0b1.00000×21 = -2.0
0xe1 = 1.11.00001 = -0b1.00001×21 = -2.0625
0xe2 = 1.11.00010 = -0b1.00010×21 = -2.125
0xe3 = 1.11.00011 = -0b1.00011×21 = -2.1875
0xe4 = 1.11.00100 = -0b1.00100×21 = -2.25
0xe5 = 1.11.00101 = -0b1.00101×21 = -2.3125
0xe6 = 1.11.00110 = -0b1.00110×21 = -2.375
0xe7 = 1.11.00111 = -0b1.00111×21 = -2.4375
0xe8 = 1.11.01000 = -0b1.01000×21 = -2.5
0xe9 = 1.11.01001 = -0b1.01001×21 = -2.5625
0xea = 1.11.01010 = -0b1.01010×21 = -2.625
0xeb = 1.11.01011 = -0b1.01011×21 = -2.6875
0xec = 1.11.01100 = -0b1.01100×21 = -2.75
0xed = 1.11.01101 = -0b1.01101×21 = -2.8125
0xee = 1.11.01110 = -0b1.01110×21 = -2.875
0xef = 1.11.01111 = -0b1.01111×21 = -2.9375
0xf0 = 1.11.10000 = -0b1.10000×21 = -3.0
0xf1 = 1.11.10001 = -0b1.10001×21 = -3.0625
0xf2 = 1.11.10010 = -0b1.10010×21 = -3.125
0xf3 = 1.11.10011 = -0b1.10011×21 = -3.1875
0xf4 = 1.11.10100 = -0b1.10100×21 = -3.25
0xf5 = 1.11.10101 = -0b1.10101×21 = -3.3125
0xf6 = 1.11.10110 = -0b1.10110×21 = -3.375
0xf7 = 1.11.10111 = -0b1.10111×21 = -3.4375
0xf8 = 1.11.11000 = -0b1.11000×21 = -3.5
0xf9 = 1.11.11001 = -0b1.11001×21 = -3.5625
0xfa = 1.11.11010 = -0b1.11010×21 = -3.625
0xfb = 1.11.11011 = -0b1.11011×21 = -3.6875
0xfc = 1.11.11100 = -0b1.11100×21 = -3.75
0xfd = 1.11.11101 = -0b1.11101×21 = -3.8125
0xfe = 1.11.11110 = -0b1.11110×21 = -3.875
0xff = 1.11.11111 = -Inf
```

C.7 Value Table: P7, $\epsilon_{\min} = 0$, $\epsilon_{\max} = 0$ (linear)

0x00 = 0.0.000000 = 0.0	0x40 = 0.1.000000 = +0b1.000000×2 ⁰ = 1.0	0x80 = 1.0.000000 = NaN	0xc0 = 1.1.000000 = -0b1.000000×2 ⁰ = -1.0
0x01 = 0.0.000001 = +0b0.000001×2 ⁰ = 0.015625	0x41 = 0.1.000001 = +0b1.000001×2 ⁰ = 1.015625	0x81 = 1.0.000001 = -0b0.000001×2 ⁰ = -0.015625	0xc1 = 1.1.000001 = -0b1.000001×2 ⁰ = -1.015625
0x02 = 0.0.000010 = +0b0.000010×2 ⁰ = 0.03125	0x42 = 0.1.000010 = +0b1.000010×2 ⁰ = 1.03125	0x82 = 1.0.000010 = -0b0.000010×2 ⁰ = -0.03125	0xc2 = 1.1.000010 = -0b1.000010×2 ⁰ = -1.03125
0x03 = 0.0.000011 = +0b0.000011×2 ⁰ = 0.046875	0x43 = 0.1.000011 = +0b1.000011×2 ⁰ = 1.046875	0x83 = 1.0.000011 = -0b0.000011×2 ⁰ = -0.046875	0xc3 = 1.1.000011 = -0b1.000011×2 ⁰ = -1.046875
0x04 = 0.0.000100 = +0b0.000100×2 ⁰ = 0.0625	0x44 = 0.1.000100 = +0b1.000100×2 ⁰ = 1.0625	0x84 = 1.0.000100 = -0b0.000100×2 ⁰ = -0.0625	0xc4 = 1.1.000100 = -0b1.000100×2 ⁰ = -1.0625
0x05 = 0.0.000101 = +0b0.000101×2 ⁰ = 0.078125	0x45 = 0.1.000101 = +0b1.000101×2 ⁰ = 1.078125	0x85 = 1.0.000101 = -0b0.000101×2 ⁰ = -0.078125	0xc5 = 1.1.000101 = -0b1.000101×2 ⁰ = -1.078125
0x06 = 0.0.000110 = +0b0.000110×2 ⁰ = 0.09375	0x46 = 0.1.000110 = +0b1.000110×2 ⁰ = 1.09375	0x86 = 1.0.000110 = -0b0.000110×2 ⁰ = -0.09375	0xc6 = 1.1.000110 = -0b1.000110×2 ⁰ = -1.09375
0x07 = 0.0.000111 = +0b0.000111×2 ⁰ = 0.109375	0x47 = 0.1.000111 = +0b1.000111×2 ⁰ = 1.109375	0x87 = 1.0.000111 = -0b0.000111×2 ⁰ = -0.109375	0xc7 = 1.1.000111 = -0b1.000111×2 ⁰ = -1.109375
0x08 = 0.0.001000 = +0b0.001000×2 ⁰ = 0.125	0x48 = 0.1.001000 = +0b1.001000×2 ⁰ = 1.125	0x88 = 1.0.001000 = -0b0.001000×2 ⁰ = -0.125	0xc8 = 1.1.001000 = -0b1.001000×2 ⁰ = -1.125
0x09 = 0.0.001001 = +0b0.001001×2 ⁰ = 0.140625	0x49 = 0.1.001001 = +0b1.001001×2 ⁰ = 1.140625	0x89 = 1.0.001001 = -0b0.001001×2 ⁰ = -0.140625	0xc9 = 1.1.001001 = -0b1.001001×2 ⁰ = -1.140625
0x0a = 0.0.001010 = +0b0.001010×2 ⁰ = 0.15625	0x4a = 0.1.001010 = +0b1.001010×2 ⁰ = 1.15625	0x8a = 1.0.001010 = -0b0.001010×2 ⁰ = -0.15625	0xca = 1.1.001010 = -0b1.001010×2 ⁰ = -1.15625
0x0b = 0.0.001011 = +0b0.001011×2 ⁰ = 0.171875	0x4b = 0.1.001011 = +0b1.001011×2 ⁰ = 1.171875	0x8b = 1.0.001011 = -0b0.001011×2 ⁰ = -0.171875	0xcb = 1.1.001011 = -0b1.001011×2 ⁰ = -1.171875
0x0c = 0.0.001100 = +0b0.001100×2 ⁰ = 0.1875	0x4c = 0.1.001100 = +0b1.001100×2 ⁰ = 1.1875	0x8c = 1.0.001100 = -0b0.001100×2 ⁰ = -0.1875	0xcc = 1.1.001100 = -0b1.001100×2 ⁰ = -1.1875
0x0d = 0.0.001101 = +0b0.001101×2 ⁰ = 0.203125	0x4d = 0.1.001101 = +0b1.001101×2 ⁰ = 1.203125	0x8d = 1.0.001101 = -0b0.001101×2 ⁰ = -0.203125	0xcd = 1.1.001101 = -0b1.001101×2 ⁰ = -1.203125
0x0e = 0.0.001110 = +0b0.001110×2 ⁰ = 0.21875	0x4e = 0.1.001110 = +0b1.001110×2 ⁰ = 1.21875	0x8e = 1.0.001110 = -0b0.001110×2 ⁰ = -0.21875	0xce = 1.1.001110 = -0b1.001110×2 ⁰ = -1.21875
0x0f = 0.0.001111 = +0b0.001111×2 ⁰ = 0.234375	0x4f = 0.1.001111 = +0b1.001111×2 ⁰ = 1.234375	0x8f = 1.0.001111 = -0b0.001111×2 ⁰ = -0.234375	0xcf = 1.1.001111 = -0b1.001111×2 ⁰ = -1.234375
0x10 = 0.0.010000 = +0b0.010000×2 ⁰ = 0.25	0x50 = 0.1.010000 = +0b1.010000×2 ⁰ = 1.25	0x90 = 1.0.010000 = -0b0.010000×2 ⁰ = -0.25	0xd0 = 1.1.010000 = -0b1.010000×2 ⁰ = -1.25
0x11 = 0.0.010001 = +0b0.010001×2 ⁰ = 0.265625	0x51 = 0.1.010001 = +0b1.010001×2 ⁰ = 1.265625	0x91 = 1.0.010001 = -0b0.010001×2 ⁰ = -0.265625	0xd1 = 1.1.010001 = -0b1.010001×2 ⁰ = -1.265625
0x12 = 0.0.010010 = +0b0.010010×2 ⁰ = 0.28125	0x52 = 0.1.010010 = +0b1.010010×2 ⁰ = 1.28125	0x92 = 1.0.010010 = -0b0.010010×2 ⁰ = -0.28125	0xd2 = 1.1.010010 = -0b1.010010×2 ⁰ = -1.28125
0x13 = 0.0.010011 = +0b0.010011×2 ⁰ = 0.296875	0x53 = 0.1.010011 = +0b1.010011×2 ⁰ = 1.296875	0x93 = 1.0.010011 = -0b0.010011×2 ⁰ = -0.296875	0xd3 = 1.1.010011 = -0b1.010011×2 ⁰ = -1.296875
0x14 = 0.0.010100 = +0b0.010100×2 ⁰ = 0.3125	0x54 = 0.1.010100 = +0b1.010100×2 ⁰ = 1.3125	0x94 = 1.0.010100 = -0b0.010100×2 ⁰ = -0.3125	0xd4 = 1.1.010100 = -0b1.010100×2 ⁰ = -1.3125
0x15 = 0.0.010101 = +0b0.010101×2 ⁰ = 0.328125	0x55 = 0.1.010101 = +0b1.010101×2 ⁰ = 1.328125	0x95 = 1.0.010101 = -0b0.010101×2 ⁰ = -0.328125	0xd5 = 1.1.010101 = -0b1.010101×2 ⁰ = -1.328125
0x16 = 0.0.010110 = +0b0.010110×2 ⁰ = 0.34375	0x56 = 0.1.010110 = +0b1.010110×2 ⁰ = 1.34375	0x96 = 1.0.010110 = -0b0.010110×2 ⁰ = -0.34375	0xd6 = 1.1.010110 = -0b1.010110×2 ⁰ = -1.34375
0x17 = 0.0.010111 = +0b0.010111×2 ⁰ = 0.359375	0x57 = 0.1.010111 = +0b1.010111×2 ⁰ = 1.359375	0x97 = 1.0.010111 = -0b0.010111×2 ⁰ = -0.359375	0xd7 = 1.1.010111 = -0b1.010111×2 ⁰ = -1.359375
0x18 = 0.0.011000 = +0b0.011000×2 ⁰ = 0.375	0x58 = 0.1.011000 = +0b1.011000×2 ⁰ = 1.375	0x98 = 1.0.011000 = -0b0.011000×2 ⁰ = -0.375	0xd8 = 1.1.011000 = -0b1.011000×2 ⁰ = -1.375
0x19 = 0.0.011001 = +0b0.011001×2 ⁰ = 0.390625	0x59 = 0.1.011001 = +0b1.011001×2 ⁰ = 1.390625	0x99 = 1.0.011001 = -0b0.011001×2 ⁰ = -0.390625	0xd9 = 1.1.011001 = -0b1.011001×2 ⁰ = -1.390625
0x1a = 0.0.011010 = +0b0.011010×2 ⁰ = 0.40625	0x5a = 0.1.011010 = +0b1.011010×2 ⁰ = 1.40625	0x9a = 1.0.011010 = -0b0.011010×2 ⁰ = -0.40625	0xda = 1.1.011010 = -0b1.011010×2 ⁰ = -1.40625
0x1b = 0.0.011011 = +0b0.011011×2 ⁰ = 0.421875	0x5b = 0.1.011011 = +0b1.011011×2 ⁰ = 1.421875	0x9b = 1.0.011011 = -0b0.011011×2 ⁰ = -0.421875	0xdb = 1.1.011011 = -0b1.011011×2 ⁰ = -1.421875
0x1c = 0.0.011100 = +0b0.011100×2 ⁰ = 0.4375	0x5c = 0.1.011100 = +0b1.011100×2 ⁰ = 1.4375	0x9c = 1.0.011100 = -0b0.011100×2 ⁰ = -0.4375	0xdc = 1.1.011100 = -0b1.011100×2 ⁰ = -1.4375
0x1d = 0.0.011101 = +0b0.011101×2 ⁰ = 0.453125	0x5d = 0.1.011101 = +0b1.011101×2 ⁰ = 1.453125	0x9d = 1.0.011101 = -0b0.011101×2 ⁰ = -0.453125	0xdd = 1.1.011101 = -0b1.011101×2 ⁰ = -1.453125
0x1e = 0.0.011110 = +0b0.011110×2 ⁰ = 0.46875	0x5e = 0.1.011110 = +0b1.011110×2 ⁰ = 1.46875	0x9e = 1.0.011110 = -0b0.011110×2 ⁰ = -0.46875	0xde = 1.1.011110 = -0b1.011110×2 ⁰ = -1.46875
0x1f = 0.0.011111 = +0b0.011111×2 ⁰ = 0.484375	0x5f = 0.1.011111 = +0b1.011111×2 ⁰ = 1.484375	0x9f = 1.0.011111 = -0b0.011111×2 ⁰ = -0.484375	0xdf = 1.1.011111 = -0b1.011111×2 ⁰ = -1.484375
0x20 = 0.0.100000 = +0b0.100000×2 ⁰ = 0.5	0x60 = 0.1.100000 = +0b1.100000×2 ⁰ = 1.5	0xa0 = 1.0.100000 = -0b0.100000×2 ⁰ = -0.5	0xe0 = 1.1.100000 = -0b1.100000×2 ⁰ = -1.5
0x21 = 0.0.100001 = +0b0.100001×2 ⁰ = 0.515625	0x61 = 0.1.100001 = +0b1.100001×2 ⁰ = 1.515625	0xa1 = 1.0.100001 = -0b0.100001×2 ⁰ = -0.515625	0xe1 = 1.1.100001 = -0b1.100001×2 ⁰ = -1.515625
0x22 = 0.0.100010 = +0b0.100010×2 ⁰ = 0.53125	0x62 = 0.1.100010 = +0b1.100010×2 ⁰ = 1.53125	0xa2 = 1.0.100010 = -0b0.100010×2 ⁰ = -0.53125	0xe2 = 1.1.100010 = -0b1.100010×2 ⁰ = -1.53125
0x23 = 0.0.100011 = +0b0.100011×2 ⁰ = 0.546875	0x63 = 0.1.100011 = +0b1.100011×2 ⁰ = 1.546875	0xa3 = 1.0.100011 = -0b0.100011×2 ⁰ = -0.546875	0xe3 = 1.1.100011 = -0b1.100011×2 ⁰ = -1.546875
0x24 = 0.0.100100 = +0b0.100100×2 ⁰ = 0.5625	0x64 = 0.1.100100 = +0b1.100100×2 ⁰ = 1.5625	0xa4 = 1.0.100100 = -0b0.100100×2 ⁰ = -0.5625	0xe4 = 1.1.100100 = -0b1.100100×2 ⁰ = -1.5625
0x25 = 0.0.100101 = +0b0.100101×2 ⁰ = 0.578125	0x65 = 0.1.100101 = +0b1.100101×2 ⁰ = 1.578125	0xa5 = 1.0.100101 = -0b0.100101×2 ⁰ = -0.578125	0xe5 = 1.1.100101 = -0b1.100101×2 ⁰ = -1.578125
0x26 = 0.0.100110 = +0b0.100110×2 ⁰ = 0.59375	0x66 = 0.1.100110 = +0b1.100110×2 ⁰ = 1.59375	0xa6 = 1.0.100110 = -0b0.100110×2 ⁰ = -0.59375	0xe6 = 1.1.100110 = -0b1.100110×2 ⁰ = -1.59375
0x27 = 0.0.100111 = +0b0.100111×2 ⁰ = 0.609375	0x67 = 0.1.100111 = +0b1.100111×2 ⁰ = 1.609375	0xa7 = 1.0.100111 = -0b0.100111×2 ⁰ = -0.609375	0xe7 = 1.1.100111 = -0b1.100111×2 ⁰ = -1.609375
0x28 = 0.0.101000 = +0b0.101000×2 ⁰ = 0.625	0x68 = 0.1.101000 = +0b1.101000×2 ⁰ = 1.625	0xa8 = 1.0.101000 = -0b0.101000×2 ⁰ = -0.625	0xe8 = 1.1.101000 = -0b1.101000×2 ⁰ = -1.625
0x29 = 0.0.101001 = +0b0.101001×2 ⁰ = 0.640625	0x69 = 0.1.101001 = +0b1.101001×2 ⁰ = 1.640625	0xa9 = 1.0.101001 = -0b0.101001×2 ⁰ = -0.640625	0xe9 = 1.1.101001 = -0b1.101001×2 ⁰ = -1.640625
0x2a = 0.0.101010 = +0b0.101010×2 ⁰ = 0.65625	0x6a = 0.1.101010 = +0b1.101010×2 ⁰ = 1.65625	0xaa = 1.0.101010 = -0b0.101010×2 ⁰ = -0.65625	0xea = 1.1.101010 = -0b1.101010×2 ⁰ = -1.65625
0x2b = 0.0.101011 = +0b0.101011×2 ⁰ = 0.671875	0x6b = 0.1.101011 = +0b1.101011×2 ⁰ = 1.671875	0xab = 1.0.101011 = -0b0.101011×2 ⁰ = -0.671875	0xeb = 1.1.101011 = -0b1.101011×2 ⁰ = -1.671875
0x2c = 0.0.101100 = +0b0.101100×2 ⁰ = 0.6875	0x6c = 0.1.101100 = +0b1.101100×2 ⁰ = 1.6875	0xac = 1.0.101100 = -0b0.101100×2 ⁰ = -0.6875	0xec = 1.1.101100 = -0b1.101100×2 ⁰ = -1.6875
0x2d = 0.0.101101 = +0b0.101101×2 ⁰ = 0.703125	0x6d = 0.1.101101 = +0b1.101101×2 ⁰ = 1.703125	0xad = 1.0.101101 = -0b0.101101×2 ⁰ = -0.703125	0xed = 1.1.101101 = -0b1.101101×2 ⁰ = -1.703125
0x2e = 0.0.101110 = +0b0.101110×2 ⁰ = 0.71875	0x6e = 0.1.101110 = +0b1.101110×2 ⁰ = 1.71875	0xae = 1.0.101110 = -0b0.101110×2 ⁰ = -0.71875	0xee = 1.1.101110 = -0b1.101110×2 ⁰ = -1.71875
0x2f = 0.0.101111 = +0b0.101111×2 ⁰ = 0.734375	0x6f = 0.1.101111 = +0b1.101111×2 ⁰ = 1.734375	0xaf = 1.0.101111 = -0b0.101111×2 ⁰ = -0.734375	0xef = 1.1.101111 = -0b1.101111×2 ⁰ = -1.734375
0x30 = 0.0.110000 = +0b0.110000×2 ⁰ = 0.75	0x70 = 0.1.110000 = +0b1.110000×2 ⁰ = 1.75	0xb0 = 1.0.110000 = -0b0.110000×2 ⁰ = -0.75	0xf0 = 1.1.110000 = -0b1.110000×2 ⁰ = -1.75
0x31 = 0.0.110001 = +0b0.110001×2 ⁰ = 0.765625	0x71 = 0.1.110001 = +0b1.110001×2 ⁰ = 1.765625	0xb1 = 1.0.110001 = -0b0.110001×2 ⁰ = -0.765625	0xf1 = 1.1.110001 = -0b1.110001×2 ⁰ = -1.765625
0x32 = 0.0.110010 = +0b0.110010×2 ⁰ = 0.78125	0x72 = 0.1.110010 = +0b1.110010×2 ⁰ = 1.78125	0xb2 = 1.0.110010 = -0b0.110010×2 ⁰ = -0.78125	0xf2 = 1.1.110010 = -0b1.110010×2 ⁰ = -1.78125
0x33 = 0.0.110011 = +0b0.110011×2 ⁰ = 0.796875	0x73 = 0.1.110011 = +0b1.110011×2 ⁰ = 1.796875	0xb3 = 1.0.110011 = -0b0.110011×2 ⁰ = -0.796875	0xf3 = 1.1.110011 = -0b1.110011×2 ⁰ = -1.796875
0x34 = 0.0.110100 = +0b0.110100×2 ⁰ = 0.8125	0x74 = 0.1.110100 = +0b1.110100×2 ⁰ = 1.8125	0xb4 = 1.0.110100 = -0b0.110100×2 ⁰ = -0.8125	0xf4 = 1.1.110100 = -0b1.110100×2 ⁰ = -1.8125
0x35 = 0.0.110101 = +0b0.110101×2 ⁰ = 0.828125	0x75 = 0.1.110101 = +0b1.110101×2 ⁰ = 1.828125	0xb5 = 1.0.110101 = -0b0.110101×2 ⁰ = -0.828125	0xf5 = 1.1.110101 = -0b1.110101×2 ⁰ = -1.828125
0x36 = 0.0.110110 = +0b0.110110×2 ⁰ = 0.84375	0x76 = 0.1.110110 = +0b1.110110×2 ⁰ = 1.84375	0xb6 = 1.0.110110 = -0b0.110110×2 ⁰ = -0.84375	0xf6 = 1.1.110110 = -0b1.110110×2 ⁰ = -1.84375
0x37 = 0.0.110111 = +0b0.110111×2 ⁰ = 0.859375	0x77 = 0.1.110111 = +0b1.110111×2 ⁰ = 1.859375	0xb7 = 1.0.110111 = -0b0.110111×2 ⁰ = -0.859375	0xf7 = 1.1.110111 = -0b1.110111×2 ⁰ = -1.859375
0x38 = 0.0.111000 = +0b0.111000×2 ⁰ = 0.875	0x78 = 0.1.111000 = +0b1.111000×2 ⁰ = 1.875	0xb8 = 1.0.111000 = -0b0.111000×2 ⁰ = -0.875	0xf8 = 1.1.111000 = -0b1.111000×2 ⁰ = -1.875
0x39 = 0.0.111001 = +0b0.111001×2 ⁰ = 0.890625	0x79 = 0.1.111001 = +0b1.111001×2 ⁰ = 1.890625	0xb9 = 1.0.111001 = -0b0.111001×2 ⁰ = -0.890625	0xf9 = 1.1.111001 = -0b1.111001×2 ⁰ = -1.890625
0x3a = 0.0.111010 = +0b0.111010×2 ⁰ = 0.90625	0x7a = 0.1.111010 = +0b1.111010×2 ⁰ = 1.90625	0xba = 1.0.111010 = -0b0.111010×2 ⁰ = -0.90625	0xfa = 1.1.111010 = -0b1.111010×2 ⁰ = -1.90625
0x3b = 0.0.111011 = +0b0.111011×2 ⁰ = 0.921875	0x7b = 0.1.111011 = +0b1.111011×2 ⁰ = 1.921875	0xbb = 1.0.111011 = -0b0.111011×2 ⁰ = -0.921875	0xfb = 1.1.11101

References

- [1] W. Kahan, “Branch cuts for complex elementary functions or much ado about nothing’s sign bit,” *Institute of Mathematics and its Applications Conference*, 1987.
<https://people.freebsd.org/~das/kahan86branch.pdf>.
- [2] W. Kahan and J. W. Thomas, “Augmenting a programming language with complex arithmetic,” tech. rep., EECS Department, University of California, Berkeley, 1991.
- [3] Google, “Jax lax package: `_float_to_int_for_sort` .”
https://github.com/google/jax/blob/fc5960f2b8b7a0ef74dbae4e27c5c08ff1564cff/jax/_src/lax/lax.py#L3934.
- [4] B. Nouné, P. Jones, D. Justus, D. Masters, and C. Luschi, “Adaptive loss scaling for mixed precision training,” tech. rep., arXiv cs.LG, 2019.
<https://arxiv.org/abs/1910.12385>.
- [5] PyTorch authors, “Pytorch torchtext package: `_t5_multi_head_attention_forward` .”
<https://github.com/pytorch/text/blob/a933cbe5a008bc2cb61d985cf5864069194157eb/torchtext/prototype/models/t5/modules.py#L236>.
- [6] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*, ch. 6.2.2.3 Softmax Units for Multinoulli Output Distributions, pp. 180–184. MIT Press, 2016.
- [7] P. Micikevicius, S. Oberman, P. Dubey, M. Cornea, A. Rodriguez, I. Bratt, R. Grisenthwaite, N. Jouppi, C. Chou, A. Huffman, M. Schulte, R. Wittig, D. Jani, and S. Deng, “OCP 8-bit floating point specification (OFP8),” tech. rep., opencompute.org, 2023.
<https://www.opencompute.org/documents/ocp-8-bit-floating-point-specification-ofp8-revision-1-0-2>
- [8] B. Nouné, P. Jones, D. Justus, D. Masters, and C. Luschi, “8-bit numerical formats for deep neural networks,” tech. rep., arXiv cs.LG, 2022.
<https://arxiv.org/abs/2206.02915>.
- [9] Tesla, Inc., “Tesla Dojo Technology: A guide to Tesla’s configurable floating point formats and arithmetic,” 2023.
https://web.archive.org/web/20230503235751/https://tesla-cdn.thron.com/static/MXMU3S_tesla-dojo-technology_1WDVZN.pdf.