

IEEE Working Group P3109 Interim Report on 8-bit Binary Floating-point Formats

Questions and comments via GitHub issues at
<https://github.com/P3109/Public>

Initial release: 18 September 2023

Version 0.6.1: 18 February 2024 (compiled 2024-02-19)

DRAFT DOCUMENT

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 | Introduction | 3 |
| 1.1 | Typographical conventions and notation | 3 |
| 2 | Values | 4 |
| 2.1 | Notes on emax | 5 |
| 2.2 | Subnormals | 6 |
| 2.3 | Not a number (NaN) | 6 |
| 2.4 | Zero | 7 |
| 2.5 | Infinities | 7 |
| 2.6 | Extremal values | 8 |
| 2.7 | Equivalents | 8 |
| 3 | Classification operators | 9 |
| 3.1 | Classification predicates | 9 |
| 3.2 | Classifier operator | 9 |
| 4 | Comparison Predicates | 11 |
| 4.1 | Details of comparison predicates | 11 |
| 4.2 | Details of totalOrder predicate | 11 |
| A | Rationales | 13 |
| A.1 | Exponent bias | 13 |
| A.2 | Infinity | 13 |
| A.2.1 | Mask Values | 13 |
| A.2.2 | Overflow to Infinity | 14 |
| A.3 | Eight Bit Formats | 14 |
| A.3.1 | binary8p3 | 14 |
| A.3.2 | binary8p4 | 14 |
| A.3.3 | binary8p5 | 14 |
| A.3.4 | binary8p6 | 14 |
| B | External Formats | 15 |
| C | Value Tables | 15 |
| C.1 | Value Table: P1, P = 1, emax = 63 | 16 |
| C.2 | Value Table: P2, P = 2, emax = 31 | 17 |
| C.3 | Value Table: P3, P = 3, emax = 15 | 18 |
| C.4 | Value Table: P4, P = 4, emax = 7 | 19 |
| C.5 | Value Table: P5, P = 5, emax = 3 | 20 |
| C.6 | Value Table: P6, P = 6, emax = 1 | 21 |
| C.7 | Value Table: P7 (Linear), P = 7, emax = 0 | 22 |

1 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.

The scope of this interim release is interchange formats, classification operators, and comparisons. The working group continues to deliberate on the specification of additional operations.

1.1 Typographical conventions and notation

Bold text describes the decisions and specifications of this document.

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

This document specifies 8-bit floating-point interchange formats (binary formats) and associated operations. Binary formats are parameterized by their width, the number of bits spanned in memory (here, 8); and their precision (P), the number of bits spanned by the true significand (this is one more than the bits of the significand that are stored explicitly).

The 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 a format with 3 bits of precision; one bit is an implicit leading bit and two bits are explicit.

2 Values

This section describes the set of values that a binary8 format shall represent. The universe of values in existing floating point usage encompasses some finite real numerical values, the non-finite numerical values positive and negative infinity ($-\text{Inf}$, $+\text{Inf}$), the non-numeric not-a-number values (NaN , NaN_1, \dots), and negative zero (-0). The value set for each binary8 format specifies the set of values that are available in that format.

Each binary format shall be associated with a unique encoding. An 8-bit binary encoding is a mapping from 8-bit strings to values. Some of these mappings are included in Appendix C.

Values are considered either “special” or “ordinary”. Encodings of the special values, shared by all binary8 formats, are shown in Table 1 as Zero, $+\text{Inf}$, $-\text{Inf}$, and NaN . The binary8 formats have only a single NaN , and relocate it to the -0 position, providing an increased range. The ordinary values consist of the normal and subnormal values.

Table 1: Special value encodings

| Special Value | Symbol | Hexadecimal Encoding | Bit Sequence |
|-------------------|---------------|----------------------|--------------|
| Zero | 0 | 0x00 | 0000 0000 |
| Positive Infinity | $+\text{Inf}$ | 0x7F | 0111 1111 |
| Negative Infinity | $-\text{Inf}$ | 0xFF | 1111 1111 |
| Not a Number | NaN | 0x80 | 1000 0000 |

These mappings are shared by all binary8 formats.

Table 2: Parameters for binary formats

| Symbol | Parameter Description | Derived Value | binary8pP, 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 |
| SE | all-special exponent | SE | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| 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 |
| emax | maximum exponent | $2^{W-1} - 1$ | 0 | 1 | 3 | 7 | 15 | 31 | 63 | 15 | 127 | 1023 |
| emin | minimum exponent | $SE - \text{emax}$ | 0 | -1 | -3 | -7 | -15 | -31 | -62 | -14 | -126 | -1022 |
| bias | exponent bias | $1 - \text{emin}$ | 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.

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

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

- All-special exponent flag SE, indicating the all-ones exponent encoding contains only special values

With binary8P1 and IEEE Std 754, there is one exponent symbol (all ones) that contains only special NaN and/or Inf values – this is indicated by the *all-special exponent* flag, SE. When this flag is set, the bias parameter is odd and emin is even.

All other parameters, such as the exponent of the largest finite value emax, are derived from the format-defining parameters. Unlike IEEE Std 754, the bias term in the binary8 formats is typically even. This allows a more symmetrical range, where emin = −emax, and arises in most binary8 formats (except P = 1).¹ While it is believed this symmetry will be useful, we need to find additional rationale to support this decision.

For $P \geq 3$, we note the binary8pP value sets are subsets of the IEEE Std 754 binary16 value set.

IEEE Std 754-2019 includes the radix B and the minimum exponent emin 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 emin is determined by P and emax; it cannot be varied independently, so it cannot be a format-defining parameter.

The P3109 committee has not yet fully reconciled the following notes from IEEE Std 754-2019:

- For binary formats, the precision p should be at least 3, as some numerical properties do not hold for lower precisions.
- Similarly, emax should be at least 2 to support the operations listed in 9.2.

The operations referred to in 9.2 are sin/cos/exp/log/etc.

It is unclear whether adopting $P = 8$ into the binary8 family provides any value. Strictly following Table 2, $\text{emax} = -\frac{1}{2}$ which means all ordinary values are irrational. Rounding this computation upward yields $\text{emax}(8) = 0$ and $\text{bias}(8) = 1$, where all representations within the format are subnormal, with the consequence that the value sets and encodings for binary8p7 and binary8p8 are identical except for a scaling factor of 1/2. As these different binary8p8 encodings are largely redundant with binary8p7, this interpretation does not appear to be useful. Finally, we could define $P = 8$ as a special case, e.g. a signed-magnitude Q7.0 integer representation with a range of 0 to 126 plus the special encodings for NaN, +Inf, and −Inf. This last interpretation, where binary8p8 is an 8-bit integer representation, might be useful in some training situations where the special values or saturating features are useful. However, at the moment, $P = 8$ remains outside of the scope of P3109.

2.1 Notes on emax

The choice of emax 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 emax. 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). Thus, the biased exponent of the largest finite value is $2^W - 2$, from which bias should be defined so that

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

¹The binary8p1 encoding is different because there are no bits in the significand allowing it to distinguish Inf from finite values.

Rearranging, we obtain the following for IEEE Std 754 formats

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

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

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

For $P = 1$, there are zero trailing significand bits, so all values where the exponent bits are ones are special, and again $\text{bias} = \text{emax}$.

2.2 Subnormals

Binary8 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 $0.T \times 2^{1-\text{bias}}$.

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 pulled 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 in Gaussian-like distributions values. 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.²

2.3 Not a number (NaN)

Binary8 value sets shall include exactly one NaN, encoded as 0x80, 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 binary8 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 binary8, supporting multiple NaNs 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.

²Conversely, subnormals provide a limited range increase to formats with narrow significands.

2.4 Zero

Binary8 formats shall have exactly one zero, encoded as 0x00. This zero value is nonnegative.

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 [4, 5]. 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.³

2.5 Infinities

Binary8 formats shall include positive and negative infinities, encoded as 0x7f and 0xff, respectively.

This decision causes a reduction in dynamic range (252 values rather than 254), 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 [1].
- Representation of overflow, for example, to adjust dynamic loss scaling factors [7].

As illustrated in Appendix A, both uses are facilitated by the presence of infinity.

³Sorting using comparison operations is undefined in the presence of NaNs. However, existing practice is to sort NaNs using `totalOrder`.

2.6 Extremal values

Table 3: Extremal values

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

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

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

2.7 Equivalents

The binary8p7 format is equivalent to a signed-magnitude fixed-point integer in Q1.6 format, except there are code points reserved for NaN, +Inf, and -Inf.

3 Classification operators

Conforming implementations shall provide the classification predicates defined by Table 4 and the classifier operator defined by Table 5. The classification predicates and the classifier function shall not signal exceptions.

The classification operators comprise: 1) a set of predicate functions with a boolean return value, taking a single binary8 value as input; 2) a classifier operator $\text{class}(x)$ that returns a single value of enumeration type, describing the input value's properties.

3.1 Classification predicates

Classification predicates shall behave as defined by Table 4.

Table 4: Predicate logic

| Predicate | Definition |
|-------------------------|--|
| $\text{isZero}(x)$ | iff x is Zero |
| $\text{isNaN}(x)$ | iff x is NaN |
| $\text{isInfinite}(x)$ | iff x is infinite |
| $\text{isFinite}(x)$ | iff x is zero, subnormal or normal |
| $\text{isNormal}(x)$ | iff x is normal, hence finite |
| $\text{isSubnormal}(x)$ | iff x is subnormal |
| $\text{isSignMinus}(x)$ | iff x has a negative sign ^a |
| $\text{isCanonical}(x)$ | True ^b |
| $\text{isSignaling}(x)$ | False ^c |

^a $\text{isSignMinus}(\text{NaN})$ is True: NaN is 0x80 (0b1000_0000).

^bThere are no non-canonical binary8 interchange formats.

^cAll binary8 formats have one NaN; it does not signal.

3.2 Classifier operator

The Classifier operator $\text{class}(x)$ tells which of the eight classes x falls into as defined by Table 5.

Table 5: Classifier operator

| Enumeration | Condition |
|-------------------|---|
| NaN | $\text{isNaN}(x)$ |
| negativeInfinity | $\text{isInfinite}(x)$ and $\text{isSignMinus}(x)$ |
| negativeNormal | $\text{isNormal}(x)$ and $\text{isSignMinus}(x)$ |
| negativeSubnormal | $\text{isSubnormal}(x)$ and $\text{isSignMinus}(x)$ |
| Zero | $\text{isZero}(x)$ |
| positiveSubnormal | $\text{isSubnormal}(x)$ and $\text{not}(\text{isSignMinus}(x))$ |
| positiveNormal | $\text{isNormal}(x)$ and $\text{not}(\text{isSignMinus}(x))$ |
| positiveInfinity | $\text{isInfinite}(x)$ and $\text{not}(\text{isSignMinus}(x))$ |

4 Comparison Predicates

Conforming implementations shall provide the comparison predicates defined by Table 6 and the `totalOrder(x, y)` predicate.

4.1 Details of comparison predicates

Comparison operations are two-argument predicates, and their negations, that return True or False. Comparisons shall not raise exceptions. 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 6: Comparison predicates and negations

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

4.2 Details of `totalOrder` predicate

The `totalOrder(x, y)` predicate provides a total ordering over each binary8 format's value set and shall return { True, False } as defined by the logic below. It shall not raise any exceptions.

```
boolean totalOrder(x, y)
    if isNaN(x): return True
    if isNaN(y): return False
    return compareLessEqual(x, y)
end
```

The above definition is consistent with the IEEE Std 754-2019 definition of `totalOrder`. In particular, among binary8 formats, there is a single NaN and it always compares as the most-negative value.

Question: should we define *totalOrder*(*x*, *y*) between different binary8 formats, e.g. between binary8p5 and binary8p4? Should we ensure that comparison predicates are also defined between different binary8 formats?

Question: do any of our binary8 formats have redundant representations for the same number, particularly when different formats are compared? if so, `totalOrder()` in IEEE Std 754-2019 defines additional sorting by exponent magnitude that we'll have to add here:

```
if( x==y ) {  
    if isNegative(x) return Exponent(x) >= Exponent(y)  
    else              return Exponent(x) <= Exponent(y)  
}
```

A Rationales

This section is being rewritten

A.1 Exponent bias

- we chose to follow 754 in defining emax - we could alternatively have defined bias consistently
this would mean no round tripping to b16
but means conversion to b16 requires an add or subtraction of one
conversion *from* b16 is expensive anyway...

A.2 Infinity

A.2.1 Mask Values

A common use for ∞ is to create masks, for example, in Transformer models in machine learning [1].

use $\beta = 1/t$. <https://openreview.net/pdf?id=LBA2Jj5Gqn>

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:

$$\log(\text{sum}(\exp(t * (A + M))))$$

where t is a “temperature” or “base” parameter [2]. This calculation depends on the property $\exp(t * (A_{ij} - \infty)) = 0$.

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 $\log(\text{sum}(\exp(v)))$ 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 maxFinite=240 without Inf, and 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.2.2 Overflow to Infinity

A second use of infinity is to indicate overflow on conversion to the binary8 type. 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.

A.3 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.3.1 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 [6], 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 | | N | Y | N | | N | |
| Include negative zero | | N | | N | | Y | | N | |
| Max exponent emax | 15 | 7 | 3 | 15 | 8 | 15 | 7 | N/A | N/A |

C Value Tables

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

A typical entry is of the form:

HEX = BINARY = BINARY_FLOAT = DECIMAL
 0x0b = 0.0001.011 = +0b1.011×2⁻⁷ = 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 decimal 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” ≈.

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

0x05 = 0.0000.101 = +0b0.101×2⁻⁷ = 0.0048828125 Subnormal value
 0x80 = 1.0000.000 = NaN Special value (NaN, +Inf, -Inf)

C.1 Value Table: P1, P = 1, emax = 63

```
0x00 = 0.0000000_ = 0.0
0x01 = 0.0000001_ = +0b1.0×2-62 ≈ 2.1684043E-19
0x02 = 0.0000001_ = +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 ≈ 2.7755576E-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 ≈ 3.725903E-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 = 0.01953125
0x37 = 0.0110111_ = +0b1.0×2-8 = 0.0390625
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 ≈ -2.7755576E-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 = 0.0100000_ = -0b1.0×2-31 ≈ -4.6566129E-10
0xa1 = 0.0100001_ = -0b1.0×2-30 ≈ -9.3132257E-10
0xa2 = 0.0100010_ = -0b1.0×2-29 ≈ -1.8626451E-09
0xa3 = 0.0100011_ = -0b1.0×2-28 ≈ -3.725903E-09
0xa4 = 0.0100100_ = -0b1.0×2-27 ≈ -7.4505806E-09
0xa5 = 0.0100101_ = -0b1.0×2-26 ≈ -1.4901161E-08
0xa6 = 0.0100110_ = -0b1.0×2-25 ≈ -2.9802322E-08
0xa7 = 0.0100111_ = -0b1.0×2-24 ≈ -5.9604645E-08
0xa8 = 0.0101000_ = -0b1.0×2-23 ≈ -1.1920929E-07
0xa9 = 0.0101001_ = -0b1.0×2-22 ≈ -2.3841858E-07
0xaa = 0.0101010_ = -0b1.0×2-21 ≈ -4.7683716E-07
0xab = 0.0101011_ = -0b1.0×2-20 ≈ -9.5367432E-07
0xac = 0.0101100_ = -0b1.0×2-19 ≈ -1.9073486E-06
0xad = 0.0101101_ = -0b1.0×2-18 ≈ -3.8146973E-06
0xae = 0.0101110_ = -0b1.0×2-17 ≈ -7.6293945E-06
0xaf = 0.0101111_ = -0b1.0×2-16 ≈ -1.5258789E-05
0xb0 = 0.0110000_ = -0b1.0×2-15 ≈ -3.0517578E-05
0xb1 = 0.0110001_ = -0b1.0×2-14 ≈ -6.1035156E-05
0xb2 = 0.0110010_ = -0b1.0×2-13 ≈ -0.00012207031
0xb3 = 0.0110011_ = -0b1.0×2-12 ≈ -0.000244140625
0xb4 = 0.0110100_ = -0b1.0×2-11 ≈ -0.00048828125
0xb5 = 0.0110101_ = -0b1.0×2-10 ≈ -0.0009765625
0xb6 = 0.0110110_ = -0b1.0×2-9 ≈ -0.01953125
0xb7 = 0.0110111_ = -0b1.0×2-8 ≈ -0.0390625
0xb8 = 0.0111000_ = -0b1.0×2-7 ≈ -0.0078125
0xb9 = 0.0111001_ = -0b1.0×2-6 ≈ -0.015625
0xba = 0.0111010_ = -0b1.0×2-5 ≈ -0.03125
0xbb = 0.0111011_ = -0b1.0×2-4 ≈ -0.0625
0xbc = 0.0111100_ = -0b1.0×2-3 ≈ -0.125
0xbd = 0.0111101_ = -0b1.0×2-2 ≈ -0.25
0xbe = 0.0111110_ = -0b1.0×2-1 ≈ -0.5
0xbf = 0.0111111_ = -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 = 1.0100000_ = -0b1.0×217 = -131072.0
0xd1 = 1.0100001_ = -0b1.0×218 = -262144.0
0xd2 = 1.0100010_ = -0b1.0×219 = -524288.0
0xd3 = 1.0100011_ = -0b1.0×220 = -1048576.0
0xd4 = 1.0101000_ = -0b1.0×221 = -2097152.0
0xd5 = 1.0101001_ = -0b1.0×222 = -4194304.0
0xd6 = 1.0101010_ = -0b1.0×223 = -8388608.0
0xd7 = 1.0101011_ = -0b1.0×224 = -16777216.0
0xd8 = 1.0101000_ = -0b1.0×225 = -33554432.0
0xd9 = 1.0101001_ = -0b1.0×226 = -67108864.0
0xda = 1.0101010_ = -0b1.0×227 = -134217728.0
0xdb = 1.0101011_ = -0b1.0×228 = -268435456.0
0xdc = 1.0101100_ = -0b1.0×229 = -536870912.0
0xdd = 1.0101101_ = -0b1.0×230 = -1073741824.0
0xde = 1.0101110_ = -0b1.0×231 = -2147483648.0
0xdf = 1.0101111_ = -0b1.0×232 = -4294967296.0
0xe0 = 1.1000000_ = -0b1.0×233 = -8589934592.0
0xe1 = 1.1000001_ = -0b1.0×234 = -17179869184.0
0xe2 = 1.1000010_ = -0b1.0×235 = -34359738368.0
0xe3 = 1.1000011_ = -0b1.0×236 = -68719476736.0
0xe4 = 1.1000100_ = -0b1.0×237 ≈ -1.3743895e+11
0xe5 = 1.1000101_ = -0b1.0×238 ≈ -2.7487791e+11
0xe6 = 1.1000110_ = -0b1.0×239 ≈ -5.4975581e+11
0xe7 = 1.1000111_ = -0b1.0×240 ≈ -1.0995116e+12
0xe8 = 1.1001000_ = -0b1.0×241 ≈ -2.1990233e+12
0xe9 = 1.1001001_ = -0b1.0×242 ≈ -4.3980465e+12
0xea = 1.1001010_ = -0b1.0×243 ≈ -8.796093e+12
0xeb = 1.1001011_ = -0b1.0×244 ≈ -1.7592186e+13
0xec = 1.1001100_ = -0b1.0×245 ≈ -3.5184372e+13
0xed = 1.1001101_ = -0b1.0×246 ≈ -7.0368744e+13
0xee = 1.1001110_ = -0b1.0×247 ≈ -1.4073749e+14
0xef = 1.1001111_ = -0b1.0×248 ≈ -2.8147498e+14
0xf0 = 1.1010000_ = -0b1.0×249 ≈ -5.6294995e+14
0xf1 = 1.1010001_ = -0b1.0×250 ≈ -1.1258999e+15
0xf2 = 1.1010010_ = -0b1.0×251 ≈ -2.2517998e+15
0xf3 = 1.1010011_ = -0b1.0×252 ≈ -4.5035996e+15
0xf4 = 1.1010100_ = -0b1.0×253 ≈ -9.0071993e+15
0xf5 = 1.1010101_ = -0b1.0×254 ≈ -1.8014399e+16
0xf6 = 1.1010110_ = -0b1.0×255 ≈ -3.6028797e+16
0xf7 = 1.1010111_ = -0b1.0×256 ≈ -7.2057594e+16
0xf8 = 1.1011000_ = -0b1.0×257 ≈ -1.4411519e+17
0xf9 = 1.1011001_ = -0b1.0×258 ≈ -2.8823038e+17
0xfa = 1.1011010_ = -0b1.0×259 ≈ -5.7646075e+17
0xfb = 1.1011011_ = -0b1.0×260 ≈ -1.1529215e+18
0xfc = 1.1011100_ = -0b1.0×261 ≈ -2.305843e+18
0xfd = 1.1011101_ = -0b1.0×262 ≈ -4.611686e+18
0xfe = 1.1011110_ = -0b1.0×263 ≈ -9.223372e+18
0xff = 1.1011111_ = -Inf
```


C.2 Value Table: P2, P = 2, emax = 31

```
0x00 = 0.000000.0 = 0.0
0x01 = 0.000000.1 = +0b0.1×2-31 ≈ 2.3283064E-10
0x02 = 0.000001.0 = +0b1.0×2-31 ≈ 4.6566129E-10
0x03 = 0.000001.1 = +0b1.1×2-31 ≈ 6.9849193E-10
0x04 = 0.000010.0 = +0b1.0×2-30 ≈ 9.3132257E-10
0x05 = 0.000010.1 = +0b1.1×2-30 ≈ 1.3969839E-09
0x06 = 0.000011.0 = +0b1.0×2-29 ≈ 1.8626451E-09
0x07 = 0.000011.1 = +0b1.1×2-29 ≈ 2.7939677E-09
0x08 = 0.000010.0 = +0b1.0×2-28 ≈ 3.7252903E-09
0x09 = 0.000100.1 = +0b1.1×2-28 ≈ 5.5879354E-09
0x0a = 0.000101.0 = +0b1.0×2-27 ≈ 7.4505806E-09
0x0b = 0.000101.1 = +0b1.1×2-27 ≈ 1.1175871E-08
0x0c = 0.000110.0 = +0b1.0×2-26 ≈ 1.4901161E-08
0x0d = 0.000110.1 = +0b1.1×2-26 ≈ 2.2351742E-08
0x0e = 0.000111.0 = +0b1.0×2-25 ≈ 2.9802322E-08
0x0f = 0.000111.1 = +0b1.1×2-25 ≈ 4.4703484E-08
0x10 = 0.001000.0 = +0b1.0×2-24 ≈ 5.9604645E-08
0x11 = 0.001000.1 = +0b1.1×2-24 ≈ 8.9406967E-08
0x12 = 0.001001.0 = +0b1.0×2-23 ≈ 1.1920929E-07
0x13 = 0.001001.1 = +0b1.1×2-23 ≈ 1.7881393E-07
0x14 = 0.001010.0 = +0b1.0×2-22 ≈ 2.3841858E-07
0x15 = 0.001010.1 = +0b1.1×2-22 ≈ 3.5762787E-07
0x16 = 0.001011.0 = +0b1.0×2-21 ≈ 4.7683716E-07
0x17 = 0.001011.1 = +0b1.1×2-21 ≈ 7.1525574E-07
0x18 = 0.001100.0 = +0b1.0×2-20 ≈ 9.5367432E-07
0x19 = 0.001100.1 = +0b1.1×2-20 ≈ 1.4305115E-06
0x1a = 0.001101.0 = +0b1.0×2-19 ≈ 1.9073486E-06
0x1b = 0.001101.1 = +0b1.1×2-19 ≈ 2.8610229E-06
0x1c = 0.001110.0 = +0b1.0×2-18 ≈ 3.8146973E-06
0x1d = 0.001110.1 = +0b1.1×2-18 ≈ 5.7220459E-06
0x1e = 0.001111.0 = +0b1.0×2-17 ≈ 7.6293945E-06
0x1f = 0.001111.1 = +0b1.1×2-17 ≈ 1.1444092E-05
0x20 = 0.010000.0 = +0b1.0×2-16 ≈ 1.5258789E-05
0x21 = 0.010000.1 = +0b1.1×2-16 ≈ 2.2888184E-05
0x22 = 0.010001.0 = +0b1.0×2-15 ≈ 3.0517578E-05
0x23 = 0.010001.1 = +0b1.1×2-15 ≈ 4.5776367E-05
0x24 = 0.010010.0 = +0b1.0×2-14 ≈ 6.1035156E-05
0x25 = 0.010010.1 = +0b1.1×2-14 ≈ 9.1552734E-05
0x26 = 0.010011.0 = +0b1.0×2-13 ≈ 0.00012207031
0x27 = 0.010011.1 = +0b1.1×2-13 ≈ 0.00018310547
0x28 = 0.010100.0 = +0b1.0×2-12 ≈ 0.000244140625
0x29 = 0.010100.1 = +0b1.1×2-12 ≈ 0.00036621094
0x2a = 0.010101.0 = +0b1.0×2-11 ≈ 0.00048828125
0x2b = 0.010101.1 = +0b1.1×2-11 ≈ 0.000732421875
0x2c = 0.010110.0 = +0b1.0×2-10 ≈ 0.0009765625
0x2d = 0.010110.1 = +0b1.1×2-10 ≈ 0.00146484375
0x2e = 0.010111.0 = +0b1.0×2-9 ≈ 0.001953125
0x2f = 0.010111.1 = +0b1.1×2-9 ≈ 0.0029296875
0x30 = 0.011000.0 = +0b1.0×2-8 ≈ 0.00390625
0x31 = 0.011000.1 = +0b1.1×2-8 ≈ 0.005859375
0x32 = 0.011001.0 = +0b1.0×2-7 ≈ 0.0078125
0x33 = 0.011001.1 = +0b1.1×2-7 ≈ 0.01171875
0x34 = 0.011010.0 = +0b1.0×2-6 ≈ 0.015625
0x35 = 0.011010.1 = +0b1.1×2-6 ≈ 0.0234375
0x36 = 0.011011.0 = +0b1.0×2-5 ≈ 0.03125
0x37 = 0.011011.1 = +0b1.1×2-5 ≈ 0.046875
0x38 = 0.011100.0 = +0b1.0×2-4 ≈ 0.0625
0x39 = 0.011100.1 = +0b1.1×2-4 ≈ 0.09375
0x3a = 0.011101.0 = +0b1.0×2-3 ≈ 0.125
0x3b = 0.011101.1 = +0b1.1×2-3 ≈ 0.1875
0x3c = 0.011110.0 = +0b1.0×2-2 ≈ 0.25
0x3d = 0.011110.1 = +0b1.1×2-2 ≈ 0.375
0x3e = 0.011111.0 = +0b1.0×2-1 ≈ 0.5
0x3f = 0.011111.1 = +0b1.1×2-1 ≈ 0.75

0x40 = 0.100000.0 = +0b1.0×20 = 1.0
0x41 = 0.100000.1 = +0b1.1×20 = 1.5
0x42 = 0.100001.0 = +0b1.0×21 = 2.0
0x43 = 0.100001.1 = +0b1.1×21 = 3.0
0x44 = 0.100010.0 = +0b1.0×22 = 4.0
0x45 = 0.100010.1 = +0b1.1×22 = 6.0
0x46 = 0.100011.0 = +0b1.0×23 = 8.0
0x47 = 0.100011.1 = +0b1.1×23 = 12.0
0x48 = 0.100100.0 = +0b1.0×24 = 16.0
0x49 = 0.100100.1 = +0b1.1×24 = 24.0
0x4a = 0.100101.0 = +0b1.0×25 = 32.0
0x4b = 0.100101.1 = +0b1.1×25 = 48.0
0x4c = 0.100110.0 = +0b1.0×26 = 64.0
0x4d = 0.100110.1 = +0b1.1×26 = 96.0
0x4e = 0.100111.0 = +0b1.0×27 = 128.0
0x4f = 0.100111.1 = +0b1.1×27 = 192.0
0x50 = 0.101000.0 = +0b1.0×28 = 256.0
0x51 = 0.101000.1 = +0b1.1×28 = 384.0
0x52 = 0.101001.0 = +0b1.0×29 = 512.0
0x53 = 0.101001.1 = +0b1.1×29 = 768.0
0x54 = 0.101010.0 = +0b1.0×210 = 1024.0
0x55 = 0.101010.1 = +0b1.1×210 = 1536.0
0x56 = 0.101011.0 = +0b1.0×211 = 2048.0
0x57 = 0.101011.1 = +0b1.1×211 = 3072.0
0x58 = 0.101100.0 = +0b1.0×212 = 4096.0
0x59 = 0.101100.1 = +0b1.1×212 = 6144.0
0x5a = 0.101101.0 = +0b1.0×213 = 8192.0
0x5b = 0.101101.1 = +0b1.1×213 = 12288.0
0x5c = 0.101110.0 = +0b1.0×214 = 16384.0
0x5d = 0.101110.1 = +0b1.1×214 = 24576.0
0x5e = 0.101111.0 = +0b1.0×215 = 32768.0
0x5f = 0.101111.1 = +0b1.1×215 = 49152.0
0x60 = 0.110000.0 = +0b1.0×216 = 65536.0
0x61 = 0.110000.1 = +0b1.1×216 = 98304.0
0x62 = 0.110001.0 = +0b1.0×217 = 131072.0
0x63 = 0.110001.1 = +0b1.1×217 = 196608.0
0x64 = 0.110010.0 = +0b1.0×218 = 262144.0
0x65 = 0.110010.1 = +0b1.1×218 = 393216.0
0x66 = 0.110011.0 = +0b1.0×219 = 524288.0
0x67 = 0.110011.1 = +0b1.1×219 = 786432.0
0x68 = 0.110100.0 = +0b1.0×220 = 1048576.0
0x69 = 0.110100.1 = +0b1.1×220 = 1572864.0
0x6a = 0.110101.0 = +0b1.0×221 = 2097152.0
0x6b = 0.110101.1 = +0b1.1×221 = 3145728.0
0x6c = 0.110110.0 = +0b1.0×222 = 4194304.0
0x6d = 0.110110.1 = +0b1.1×222 = 6291456.0
0x6e = 0.110111.0 = +0b1.0×223 = 8388608.0
0x6f = 0.110111.1 = +0b1.1×223 = 12582912.0
0x70 = 0.111000.0 = +0b1.0×224 = 16777216.0
0x71 = 0.111000.1 = +0b1.1×224 = 25165824.0
0x72 = 0.111001.0 = +0b1.0×225 = 33554432.0
0x73 = 0.111001.1 = +0b1.1×225 = 50331648.0
0x74 = 0.111010.0 = +0b1.0×226 = 67108864.0
0x75 = 0.111010.1 = +0b1.1×226 = 102463296.0
0x76 = 0.111011.0 = +0b1.0×227 = 134217728.0
0x77 = 0.111011.1 = +0b1.1×227 = 201326592.0
0x78 = 0.111100.0 = +0b1.0×228 = 268435456.0
0x79 = 0.111100.1 = +0b1.1×228 = 402653184.0
0x7a = 0.111101.0 = +0b1.0×229 = 536870912.0
0x7b = 0.111101.1 = +0b1.1×229 = 805306368.0
0x7c = 0.111110.0 = +0b1.0×230 = 1073741824.0
0x7d = 0.111110.1 = +0b1.1×230 = 1610612736.0
0x7e = 0.111111.0 = +0b1.0×231 = 2147483648.0
0x7f = 0.111111.1 = +Inf

0x80 = 1.000000.0 = NaN
0x81 = 1.000000.1 = -0b0.1×2-31 ≈ -2.3283064E-10
0x82 = 1.000001.0 = -0b1.0×2-31 ≈ -4.6566129E-10
0x83 = 1.000001.1 = -0b1.1×2-31 ≈ -6.9849193E-10
0x84 = 1.000010.0 = -0b1.0×2-30 ≈ -9.3132257E-10
0x85 = 1.000010.1 = -0b1.1×2-30 ≈ -1.3969839E-09
0x86 = 1.000011.0 = -0b1.0×2-29 ≈ -1.8626451E-09
0x87 = 1.000011.1 = -0b1.1×2-29 ≈ -2.7939677E-09
0x88 = 1.000100.0 = -0b1.0×2-28 ≈ -3.7252903E-09
0x89 = 1.000100.1 = -0b1.1×2-28 ≈ -5.5879354E-09
0x8a = 1.000101.0 = -0b1.0×2-27 ≈ -7.4505806E-09
0x8b = 1.000101.1 = -0b1.1×2-27 ≈ -1.1175871E-08
0x8c = 1.000110.0 = -0b1.0×2-26 ≈ -1.4901161E-08
0x8d = 1.000110.1 = -0b1.1×2-26 ≈ -2.2351742E-08
0x8e = 1.000111.0 = -0b1.0×2-25 ≈ -2.9802322E-08
0x8f = 1.000111.1 = -0b1.1×2-25 ≈ -4.4703484E-08
0x90 = 1.001000.0 = -0b1.0×2-24 ≈ -5.9604645E-08
0x91 = 1.001000.1 = -0b1.1×2-24 ≈ -8.9406967E-08
0x92 = 1.001001.0 = -0b1.0×2-23 ≈ -1.1920929E-07
0x93 = 1.001001.1 = -0b1.1×2-23 ≈ -1.7881393E-07
0x94 = 1.001010.0 = -0b1.0×2-22 ≈ -2.3841858E-07
0x95 = 1.001010.1 = -0b1.1×2-22 ≈ -3.5762787E-07
0x96 = 1.001011.0 = -0b1.0×2-21 ≈ -4.7683716E-07
0x97 = 1.001011.1 = -0b1.1×2-21 ≈ -7.1525574E-07
0x98 = 1.001100.0 = -0b1.0×2-20 ≈ -9.5367432E-07
0x99 = 1.001100.1 = -0b1.1×2-20 ≈ -1.4305115E-06
0x9a = 1.001101.0 = -0b1.0×2-19 ≈ -1.9073486E-06
0x9b = 1.001101.1 = -0b1.1×2-19 ≈ -2.8610229E-06
0x9c = 1.001110.0 = -0b1.0×2-18 ≈ -3.8146973E-06
0x9d = 1.001110.1 = -0b1.1×2-18 ≈ -5.7220459E-06
0x9e = 1.001111.0 = -0b1.0×2-17 ≈ -7.6293945E-06
0x9f = 1.001111.1 = -0b1.1×2-17 ≈ -1.1444092E-05
0xa0 = 1.010000.0 = -0b1.0×2-16 ≈ -1.5258789E-05
0xa1 = 1.010000.1 = -0b1.1×2-16 ≈ -2.2888184E-05
0xa2 = 1.010001.0 = -0b1.0×2-15 ≈ -3.0517578E-05
0xa3 = 1.010001.1 = -0b1.1×2-15 ≈ -4.5776367E-05
0xa4 = 1.010010.0 = -0b1.0×2-14 ≈ -6.1035156E-05
0xa5 = 1.010010.1 = -0b1.1×2-14 ≈ -9.1552734E-05
0xa6 = 1.010011.0 = -0b1.0×2-13 ≈ -0.00012207031
0xa7 = 1.010011.1 = -0b1.1×2-13 ≈ -0.00018310547
0xa8 = 1.010100.0 = -0b1.0×2-12 ≈ -0.000244140625
0xa9 = 1.010100.1 = -0b1.1×2-12 ≈ -0.00036621094
0xaa = 1.010101.0 = -0b1.0×2-11 ≈ -0.00048828125
0xab = 1.010101.1 = -0b1.1×2-11 ≈ -0.000732421875
0xac = 1.010110.0 = -0b1.0×2-10 ≈ -0.0009765625
0xad = 1.010110.1 = -0b1.1×2-10 ≈ -0.00146484375
0xae = 1.010111.0 = -0b1.0×2-9 ≈ -0.001953125
0xaf = 1.010111.1 = -0b1.1×2-9 ≈ -0.0029296875
0xb0 = 1.011000.0 = -0b1.0×2-8 ≈ -0.00390625
0xb1 = 1.011000.1 = -0b1.1×2-8 ≈ -0.005859375
0xb2 = 1.011001.0 = -0b1.0×2-7 ≈ -0.0078125
0xb3 = 1.011001.1 = -0b1.1×2-7 ≈ -0.01171875
0xb4 = 1.011010.0 = -0b1.0×2-6 ≈ -0.015625
0xb5 = 1.011010.1 = -0b1.1×2-6 ≈ -0.0234375
0xb6 = 1.011011.0 = -0b1.0×2-5 ≈ -0.03125
0xb7 = 1.011011.1 = -0b1.1×2-5 ≈ -0.046875
0xb8 = 1.011100.0 = -0b1.0×2-4 ≈ -0.0625
0xb9 = 1.011100.1 = -0b1.1×2-4 ≈ -0.09375
0xba = 1.011101.0 = -0b1.0×2-3 ≈ -0.125
0xbb = 1.011101.1 = -0b1.1×2-3 ≈ -0.1875
0xbc = 1.011110.0 = -0b1.0×2-2 ≈ -0.25
0xbd = 1.011110.1 = -0b1.1×2-2 ≈ -0.375
0xbe = 1.011111.0 = -0b1.0×2-1 ≈ -0.5
0xbf = 1.011111.1 = -0b1.1×2-1 ≈ -0.75

0xc0 = 1.100000.0 = -0b1.0×20 = -1.0
0xc1 = 1.100000.1 = -0b1.1×20 = -1.5
0xc2 = 1.100001.0 = -0b1.0×21 = -2.0
0xc3 = 1.100001.1 = -0b1.1×21 = -3.0
0xc4 = 1.100010.0 = -0b1.0×22 = -4.0
0xc5 = 1.100010.1 = -0b1.1×22 = -6.0
0xc6 = 1.100011.0 = -0b1.0×23 = -8.0
0xc7 = 1.100011.1 = -0b1.1×23 = -12.0
0xc8 = 1.100100.0 = -0b1.0×24 = -16.0
0xc9 = 1.100100.1 = -0b1.1×24 = -24.0
0xca = 1.100101.0 = -0b1.0×25 = -32.0
0xcb = 1.100101.1 = -0b1.1×25 = -48.0
0xcc = 1.100110.0 = -0b1.0×26 = -64.0
0xcd = 1.100110.1 = -0b1.1×26 = -96.0
0xce = 1.100111.0 = -0b1.0×27 = -128.0
0xcf = 1.100111.1 = -0b1.1×27 = -192.0
0xd0 = 1.101000.0 = -0b1.0×28 = -256.0
0xd1 = 1.101000.1 = -0b1.1×28 = -384.0
0xd2 = 1.101001.0 = -0b1.0×29 = -512.0
0xd3 = 1.101001.1 = -0b1.1×29 = -768.0
0xd4 = 1.101010.0 = -0b1.0×210 = -1024.0
0xd5 = 1.101010.1 = -0b1.1×210 = -1536.0
0xd6 = 1.101011.0 = -0b1.0×211 = -2048.0
0xd7 = 1.101011.1 = -0b1.1×211 = -3072.0
0xd8 = 1.101100.0 = -0b1.0×212 = -4096.0
0xd9 = 1.101100.1 = -0b1.1×212 = -6144.0
0xda = 1.101101.0 = -0b1.0×213 = -8192.0
0xdb = 1.101101.1 = -0b1.1×213 = -12288.0
0xdc = 1.101110.0 = -0b1.0×214 = -16384.0
0xdd = 1.101110.1 = -0b1.1×214 = -24576.0
0xde = 1.101111.0 = -0b1.0×215 = -32768.0
0xdf = 1.101111.1 = -0b1.1×215 = -49152.0
0xe0 = 1.110000.0 = -0b1.0×216 = -65536.0
0xe1 = 1.110000.1 = -0b1.1×216 = -98304.0
0xe2 = 1.110001.0 = -0b1.0×217 = -131072.0
0xe3 = 1.110001.1 = -0b1.1×217 = -196608.0
0xe4 = 1.110010.0 = -0b1.0×218 = -262144.0
0xe5 = 1.110010.1 = -0b1.1×218 = -393216.0
0xe6 = 1.110011.0 = -0b1.0×219 = -524288.0
0xe7 = 1.110011.1 = -0b1.1×219 = -786432.0
0xe8 = 1.110100.0 = -0b1.0×220 = -1048576.0
0xe9 = 1.110100.1 = -0b1.1×220 = -1572864.0
0xea = 1.110101.0 = -0b1.0×221 = -2097152.0
0xeb = 1.110101.1 = -0b1.1×221 = -3145728.0
0xec = 1.110110.0 = -0b1.0×222 = -4194304.0
0xed = 1.110110.1 = -0b1.1×222 = -6291456.0
0xee = 1.110111.0 = -0b1.0×223 = -8388608.0
0xef = 1.110111.1 = -0b1.1×223 = -12582912.0
0xf0 = 1.111000.0 = -0b1.0×224 = -16777216.0
0xf1 = 1.111000.1 = -0b1.1×224 = -25165824.0
0xf2 = 1.111001.0 = -0b1.0×225 = -33554432.0
0xf3 = 1.111001.1 = -0b1.1×225 = -50331648.0
0xf4 = 1.111010.0 = -0b1.0×226 = -67108864.0
0xf5 = 1.111010.1 = -0b1.1×226 = -102463296.0
0xf6 = 1.111011.0 = -0b1.0×227 = -134217728.0
0xf7 = 1.111011.1 = -0b1.1×227 = -201326592.0
0xf8 = 1.111100.0 = -0b1.0×228 = -268435456.0
0xf9 = 1.111100.1 = -0b1.1×228 = -402653184.0
0xfa = 1.111101.0 = -0b1.0×229 = -536870912.0
0xfb = 1.111101.1 = -0b1.1×229 = -805306368.0
0xfc = 1.111110.0 = -0b1.0×230 = -1073741824.0
0xfd = 1.111110.1 = -0b1.1×230 = -1610612736.0
0xfe = 1.111111.0 = -0b1.0×231 = -2147483648.0
0xff = 1.111111.1 = -Inf
```

C.3 Value Table: P3, P = 3, emax = 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 = +0b1.10×2 ⁻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 = -0.875 | 0xff = 1. |

C.4 Value Table: P4, P = 4, 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 ⁻¹ = 0.5 | 0x78 = 0.1111.000 = +0b1.000×2 ⁷⁷ = 128.0 | 0xb8 = 1.0111.000 = -0b1.000×2 ⁻¹ = -0.5 | 0xf8 = 1.1111.000 = -0b1.000×2 ⁷⁷ = -128.0 |
| 0x39 = 0.0111.001 = +0b1.001×2 ⁻¹ = 0.5625 | 0x79 = 0.1111.001 = +0b1.001×2 ⁷⁷ = 144.0 | 0xb9 = 1.0111.001 = -0b1.001×2 ⁻¹ = -0.5625 | 0xf9 = 1.1111.001 = -0b1.001×2 ⁷⁷ = -144.0 |
| 0x3a = 0.0111.010 = +0b1.010×2 ⁻¹ = 0.625 | 0x7a = 0.1111.010 = +0b1.010×2 ⁷⁷ = 160.0 | 0xba = 1.0111.010 = -0b1.010×2 ⁻¹ = -0.625 | 0xfa = 1.1111.010 = -0b1.010×2 ⁷⁷ = -160.0 |
| 0x3b = 0.0111.011 = +0b1.011×2 ⁻¹ = 0.6875 | 0x7b = 0.1111.011 = +0b1.011×2 ⁷⁷ = 176.0 | 0xbb = 1.0111.011 = -0b1.011×2 ⁻¹ = -0.6875 | 0xfb = 1.1111.011 = -0b1.011×2 ⁷⁷ = -176.0 |
| 0x3c = 0.0111.100 = +0b1.100×2 ⁻¹ = 0.75 | 0x7c = 0.1111.100 = +0b1.100×2 ⁷⁷ = 192.0 | 0xbc = 1.0111.100 = -0b1.100×2 ⁻¹ = -0.75 | 0xfc = 1.1111.100 = -0b1.100×2 ⁷⁷ = -192.0 |
| 0x3d = 0.0111.101 = +0b1.101×2 ⁻¹ = 0.8125 | 0x7d = 0.1111.101 = +0b1.101×2 ⁷⁷ = 208.0 | 0xbd = 1.0111.101 = -0b1.101×2 ⁻¹ = -0.8125 | 0xfd = 1.1111.101 = -0b1.101×2 ⁷⁷ = -208.0 |
| 0x3e = 0.0111.110 = +0b1.110×2 ⁻¹ = 0.875 | 0x7e = 0.1111.110 = +0b1.110×2 ⁷⁷ = 224.0 | 0xbe = 1.0111.110 = -0b1.110×2 ⁻¹ = -0.875 | 0xfe = 1.1111.110 = -0b1.110×2 ⁷⁷ = -224.0 |
| 0x3f = 0.0111.111 = +0b1.111×2 ⁻¹ = 0.9375 | 0x7f = 0.1111.111 = +Inf | 0xbf = 1.0111.111 = -0b1.111×2 ⁻¹ = -0.9375 | 0xff = 1.1111.111 = -Inf |

C.5 Value Table: P5, P = 5, 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 = +0b0.1000×2 ⁻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, P = 6, 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 (Linear), P = 7, emax = 0

0x00 = 0.0.000000 = 0.0
0x01 = 0.0.000001 = +0b0.000001×2⁰ = 0.015625
0x02 = 0.0.000010 = +0b0.000010×2⁰ = 0.03125
0x03 = 0.0.000011 = +0b0.000011×2⁰ = 0.046875
0x04 = 0.0.000100 = +0b0.000100×2⁰ = 0.0625
0x05 = 0.0.000101 = +0b0.000101×2⁰ = 0.078125
0x06 = 0.0.000110 = +0b0.000110×2⁰ = 0.09375
0x07 = 0.0.000111 = +0b0.000111×2⁰ = 0.109375
0x08 = 0.0.001000 = +0b0.001000×2⁰ = 0.125
0x09 = 0.0.001001 = +0b0.001001×2⁰ = 0.140625
0x0a = 0.0.001010 = +0b0.001010×2⁰ = 0.15625
0x0b = 0.0.001011 = +0b0.001011×2⁰ = 0.171875
0x0c = 0.0.001100 = +0b0.001100×2⁰ = 0.1875
0x0d = 0.0.001101 = +0b0.001101×2⁰ = 0.203125
0x0e = 0.0.001110 = +0b0.001110×2⁰ = 0.21875
0x0f = 0.0.001111 = +0b0.001111×2⁰ = 0.234375
0x10 = 0.0.010000 = +0b0.010000×2⁰ = 0.25
0x11 = 0.0.010001 = +0b0.010001×2⁰ = 0.265625
0x12 = 0.0.010010 = +0b0.010010×2⁰ = 0.28125
0x13 = 0.0.010011 = +0b0.010011×2⁰ = 0.296875
0x14 = 0.0.010100 = +0b0.010100×2⁰ = 0.3125
0x15 = 0.0.010101 = +0b0.010101×2⁰ = 0.328125
0x16 = 0.0.010110 = +0b0.010110×2⁰ = 0.34375
0x17 = 0.0.010111 = +0b0.010111×2⁰ = 0.359375
0x18 = 0.0.011000 = +0b0.011000×2⁰ = 0.375
0x19 = 0.0.011001 = +0b0.011001×2⁰ = 0.390625
0x1a = 0.0.011010 = +0b0.011010×2⁰ = 0.40625
0x1b = 0.0.011011 = +0b0.011011×2⁰ = 0.421875
0x1c = 0.0.011100 = +0b0.011100×2⁰ = 0.4375
0x1d = 0.0.011101 = +0b0.011101×2⁰ = 0.453125
0x1e = 0.0.011110 = +0b0.011110×2⁰ = 0.46875
0x1f = 0.0.011111 = +0b0.011111×2⁰ = 0.484375
0x20 = 0.0.100000 = +0b0.100000×2⁰ = 0.5
0x21 = 0.0.100001 = +0b0.100001×2⁰ = 0.515625
0x22 = 0.0.100010 = +0b0.100010×2⁰ = 0.53125
0x23 = 0.0.100011 = +0b0.100011×2⁰ = 0.546875
0x24 = 0.0.100100 = +0b0.100100×2⁰ = 0.5625
0x25 = 0.0.100101 = +0b0.100101×2⁰ = 0.578125
0x26 = 0.0.100110 = +0b0.100110×2⁰ = 0.59375
0x27 = 0.0.100111 = +0b0.100111×2⁰ = 0.609375
0x28 = 0.0.101000 = +0b0.101000×2⁰ = 0.625
0x29 = 0.0.101001 = +0b0.101001×2⁰ = 0.640625
0x2a = 0.0.101010 = +0b0.101010×2⁰ = 0.65625
0x2b = 0.0.101011 = +0b0.101011×2⁰ = 0.671875
0x2c = 0.0.101100 = +0b0.101100×2⁰ = 0.6875
0x2d = 0.0.101101 = +0b0.101101×2⁰ = 0.703125
0x2e = 0.0.101110 = +0b0.101110×2⁰ = 0.71875
0x2f = 0.0.101111 = +0b0.101111×2⁰ = 0.734375
0x30 = 0.0.110000 = +0b0.110000×2⁰ = 0.75
0x31 = 0.0.110001 = +0b0.110001×2⁰ = 0.765625
0x32 = 0.0.110010 = +0b0.110010×2⁰ = 0.78125
0x33 = 0.0.110011 = +0b0.110011×2⁰ = 0.796875
0x34 = 0.0.110100 = +0b0.110100×2⁰ = 0.8125
0x35 = 0.0.110101 = +0b0.110101×2⁰ = 0.828125
0x36 = 0.0.110110 = +0b0.110110×2⁰ = 0.84375
0x37 = 0.0.110111 = +0b0.110111×2⁰ = 0.859375
0x38 = 0.0.111000 = +0b0.111000×2⁰ = 0.875
0x39 = 0.0.111001 = +0b0.111001×2⁰ = 0.890625
0x3a = 0.0.111010 = +0b0.111010×2⁰ = 0.90625
0x3b = 0.0.111011 = +0b0.111011×2⁰ = 0.921875
0x3c = 0.0.111100 = +0b0.111100×2⁰ = 0.9375
0x3d = 0.0.111101 = +0b0.111101×2⁰ = 0.953125
0x3e = 0.0.111110 = +0b0.111110×2⁰ = 0.96875
0x3f = 0.0.111111 = +0b0.111111×2⁰ = 0.984375

0x40 = 0.1.000000 = +0b1.000000×2⁰ = 1.0
0x41 = 0.1.000001 = +0b1.000001×2⁰ = 1.015625
0x42 = 0.1.000010 = +0b1.000010×2⁰ = 1.03125
0x43 = 0.1.000011 = +0b1.000011×2⁰ = 1.046875
0x44 = 0.1.000100 = +0b1.000100×2⁰ = 1.0625
0x45 = 0.1.000101 = +0b1.000101×2⁰ = 1.078125
0x46 = 0.1.000110 = +0b1.000110×2⁰ = 1.09375
0x47 = 0.1.000111 = +0b1.000111×2⁰ = 1.109375
0x48 = 0.1.001000 = +0b1.001000×2⁰ = 1.125
0x49 = 0.1.001001 = +0b1.001001×2⁰ = 1.140625
0x4a = 0.1.001010 = +0b1.001010×2⁰ = 1.15625
0x4b = 0.1.001011 = +0b1.001011×2⁰ = 1.171875
0x4c = 0.1.001100 = +0b1.001100×2⁰ = 1.1875
0x4d = 0.1.001101 = +0b1.001101×2⁰ = 1.203125
0x4e = 0.1.001110 = +0b1.001110×2⁰ = 1.21875
0x4f = 0.1.001111 = +0b1.001111×2⁰ = 1.234375
0x50 = 0.1.010000 = +0b1.010000×2⁰ = 1.25
0x51 = 0.1.010001 = +0b1.010001×2⁰ = 1.265625
0x52 = 0.1.010010 = +0b1.010010×2⁰ = 1.28125
0x53 = 0.1.010011 = +0b1.010011×2⁰ = 1.296875
0x54 = 0.1.010100 = +0b1.010100×2⁰ = 1.3125
0x55 = 0.1.010101 = +0b1.010101×2⁰ = 1.328125
0x56 = 0.1.010110 = +0b1.010110×2⁰ = 1.34375
0x57 = 0.1.010111 = +0b1.010111×2⁰ = 1.359375
0x58 = 0.1.011000 = +0b1.011000×2⁰ = 1.375
0x59 = 0.1.011001 = +0b1.011001×2⁰ = 1.390625
0x5a = 0.1.011010 = +0b1.011010×2⁰ = 1.40625
0x5b = 0.1.011011 = +0b1.011011×2⁰ = 1.421875
0x5c = 0.1.011100 = +0b1.011100×2⁰ = 1.4375
0x5d = 0.1.011101 = +0b1.011101×2⁰ = 1.453125
0x5e = 0.1.011110 = +0b1.011110×2⁰ = 1.46875
0x5f = 0.1.011111 = +0b1.011111×2⁰ = 1.484375
0x60 = 0.1.100000 = +0b1.100000×2⁰ = 1.5
0x61 = 0.1.100001 = +0b1.100001×2⁰ = 1.515625
0x62 = 0.1.100010 = +0b1.100010×2⁰ = 1.53125
0x63 = 0.1.100011 = +0b1.100011×2⁰ = 1.546875
0x64 = 0.1.100100 = +0b1.100100×2⁰ = 1.5625
0x65 = 0.1.100101 = +0b1.100101×2⁰ = 1.578125
0x66 = 0.1.100110 = +0b1.100110×2⁰ = 1.59375
0x67 = 0.1.100111 = +0b1.100111×2⁰ = 1.609375
0x68 = 0.1.101000 = +0b1.101000×2⁰ = 1.625
0x69 = 0.1.101001 = +0b1.101001×2⁰ = 1.640625
0x6a = 0.1.101010 = +0b1.101010×2⁰ = 1.65625
0x6b = 0.1.101011 = +0b1.101011×2⁰ = 1.671875
0x6c = 0.1.101100 = +0b1.101100×2⁰ = 1.6875
0x6d = 0.1.101101 = +0b1.101101×2⁰ = 1.703125
0x6e = 0.1.101110 = +0b1.101110×2⁰ = 1.71875
0x6f = 0.1.101111 = +0b1.101111×2⁰ = 1.734375
0x70 = 0.1.110000 = +0b1.110000×2⁰ = 1.75
0x71 = 0.1.110001 = +0b1.110001×2⁰ = 1.765625
0x72 = 0.1.110010 = +0b1.110010×2⁰ = 1.78125
0x73 = 0.1.110011 = +0b1.110011×2⁰ = 1.796875
0x74 = 0.1.110100 = +0b1.110100×2⁰ = 1.8125
0x75 = 0.1.110101 = +0b1.110101×2⁰ = 1.828125
0x76 = 0.1.110110 = +0b1.110110×2⁰ = 1.84375
0x77 = 0.1.110111 = +0b1.110111×2⁰ = 1.859375
0x78 = 0.1.111000 = +0b1.111000×2⁰ = 1.875
0x79 = 0.1.111001 = +0b1.111001×2⁰ = 1.890625
0x7a = 0.1.111010 = +0b1.111010×2⁰ = 1.90625
0x7b = 0.1.111011 = +0b1.111011×2⁰ = 1.921875
0x7c = 0.1.111100 = +0b1.111100×2⁰ = 1.9375
0x7d = 0.1.111101 = +0b1.111101×2⁰ = 1.953125
0x7e = 0.1.111110 = +0b1.111110×2⁰ = 1.96875
0x7f = 0.1.111111 = +Inf

0x80 = 1.0.000000 = NaN
0x81 = 1.0.000001 = -0b0.000001×2⁰ = -0.015625
0x82 = 1.0.000010 = -0b0.000010×2⁰ = -0.03125
0x83 = 1.0.000011 = -0b0.000011×2⁰ = -0.046875
0x84 = 1.0.000100 = -0b0.000100×2⁰ = -0.0625
0x85 = 1.0.000101 = -0b0.000101×2⁰ = -0.078125
0x86 = 1.0.000110 = -0b0.000110×2⁰ = -0.09375
0x87 = 1.0.000111 = -0b0.000111×2⁰ = -0.109375
0x88 = 1.0.001000 = -0b0.001000×2⁰ = -0.125
0x89 = 1.0.001001 = -0b0.001001×2⁰ = -0.140625
0x8a = 1.0.001010 = -0b0.001010×2⁰ = -0.15625
0x8b = 1.0.001011 = -0b0.001011×2⁰ = -0.171875
0x8c = 1.0.001100 = -0b0.001100×2⁰ = -0.1875
0x8d = 1.0.001101 = -0b0.001101×2⁰ = -0.203125
0x8e = 1.0.001110 = -0b0.001110×2⁰ = -0.21875
0x8f = 1.0.001111 = -0b0.001111×2⁰ = -0.234375
0x90 = 1.0.010000 = -0b0.010000×2⁰ = -0.25
0x91 = 1.0.010001 = -0b0.010001×2⁰ = -0.265625
0x92 = 1.0.010010 = -0b0.010010×2⁰ = -0.28125
0x93 = 1.0.010011 = -0b0.010011×2⁰ = -0.296875
0x94 = 1.0.010100 = -0b0.010100×2⁰ = -0.3125
0x95 = 1.0.010101 = -0b0.010101×2⁰ = -0.328125
0x96 = 1.0.010110 = -0b0.010110×2⁰ = -0.34375
0x97 = 1.0.010111 = -0b0.010111×2⁰ = -0.359375
0x98 = 1.0.011000 = -0b0.011000×2⁰ = -0.375
0x99 = 1.0.011001 = -0b0.011001×2⁰ = -0.390625
0x9a = 1.0.011010 = -0b0.011010×2⁰ = -0.40625
0x9b = 1.0.011011 = -0b0.011011×2⁰ = -0.421875
0x9c = 1.0.011100 = -0b0.011100×2⁰ = -0.4375
0x9d = 1.0.011101 = -0b0.011101×2⁰ = -0.453125
0x9e = 1.0.011110 = -0b0.011110×2⁰ = -0.46875
0x9f = 1.0.011111 = -0b0.011111×2⁰ = -0.484375
0xa0 = 1.0.100000 = -0b0.100000×2⁰ = -0.5
0xa1 = 1.0.100001 = -0b0.100001×2⁰ = -0.515625
0xa2 = 1.0.100010 = -0b0.100010×2⁰ = -0.53125
0xa3 = 1.0.100011 = -0b0.100011×2⁰ = -0.546875
0xa4 = 1.0.100100 = -0b0.100100×2⁰ = -0.5625
0xa5 = 1.0.100101 = -0b0.100101×2⁰ = -0.578125
0xa6 = 1.0.100110 = -0b0.100110×2⁰ = -0.59375
0xa7 = 1.0.100111 = -0b0.100111×2⁰ = -0.609375
0xa8 = 1.0.101000 = -0b0.101000×2⁰ = -0.625
0xa9 = 1.0.101001 = -0b0.101001×2⁰ = -0.640625
0xaa = 1.0.101010 = -0b0.101010×2⁰ = -0.65625
0xab = 1.0.101011 = -0b0.101011×2⁰ = -0.671875
0xac = 1.0.101100 = -0b0.101100×2⁰ = -0.6875
0xad = 1.0.101101 = -0b0.101101×2⁰ = -0.703125
0xae = 1.0.101110 = -0b0.101110×2⁰ = -0.71875
0xaf = 1.0.101111 = -0b0.101111×2⁰ = -0.734375
0xb0 = 1.0.110000 = -0b0.110000×2⁰ = -0.75
0xb1 = 1.0.110001 = -0b0.110001×2⁰ = -0.765625
0xb2 = 1.0.110010 = -0b0.110010×2⁰ = -0.78125
0xb3 = 1.0.110011 = -0b0.110011×2⁰ = -0.796875
0xb4 = 1.0.110100 = -0b0.110100×2⁰ = -0.8125
0xb5 = 1.0.110101 = -0b0.110101×2⁰ = -0.828125
0xb6 = 1.0.110110 = -0b0.110110×2⁰ = -0.84375
0xb7 = 1.0.110111 = -0b0.110111×2⁰ = -0.859375
0xb8 = 1.0.111000 = -0b0.111000×2⁰ = -0.875
0xb9 = 1.0.111001 = -0b0.111001×2⁰ = -0.890625
0xba = 1.0.111010 = -0b0.111010×2⁰ = -0.90625
0xbb = 1.0.111011 = -0b0.111011×2⁰ = -0.921875
0xbc = 1.0.111100 = -0b0.111100×2⁰ = -0.9375
0xbd = 1.0.111101 = -0b0.111101×2⁰ = -0.953125
0xbe = 1.0.111110 = -0b0.111110×2⁰ = -0.96875
0xbf = 1.0.111111 = -0b0.111111×2⁰ = -0.984375

0xc0 = 1.1.000000 = -0b1.000000×2⁰ = -1.0
0xc1 = 1.1.000001 = -0b1.000001×2⁰ = -1.015625
0xc2 = 1.1.000010 = -0b1.000010×2⁰ = -1.03125
0xc3 = 1.1.000011 = -0b1.000011×2⁰ = -1.046875
0xc4 = 1.1.000100 = -0b1.000100×2⁰ = -1.0625
0xc5 = 1.1.000101 = -0b1.000101×2⁰ = -1.078125
0xc6 = 1.1.000110 = -0b1.000110×2⁰ = -1.09375
0xc7 = 1.1.000111 = -0b1.000111×2⁰ = -1.109375
0xc8 = 1.1.001000 = -0b1.001000×2⁰ = -1.125
0xc9 = 1.1.001001 = -0b1.001001×2⁰ = -1.140625
0xca = 1.1.001010 = -0b1.001010×2⁰ = -1.15625
0xcb = 1.1.001011 = -0b1.001011×2⁰ = -1.171875
0xcc = 1.1.001100 = -0b1.001100×2⁰ = -1.1875
0xcd = 1.1.001101 = -0b1.001101×2⁰ = -1.203125
0xce = 1.1.001110 = -0b1.001110×2⁰ = -1.21875
0xcf = 1.1.001111 = -0b1.001111×2⁰ = -1.234375
0xd0 = 1.1.010000 = -0b1.010000×2⁰ = -1.25
0xd1 = 1.1.010001 = -0b1.010001×2⁰ = -1.265625
0xd2 = 1.1.010010 = -0b1.010010×2⁰ = -1.28125
0xd3 = 1.1.010011 = -0b1.010011×2⁰ = -1.296875
0xd4 = 1.1.010100 = -0b1.010100×2⁰ = -1.3125
0xd5 = 1.1.010101 = -0b1.010101×2⁰ = -1.328125
0xd6 = 1.1.010110 = -0b1.010110×2⁰ = -1.34375
0xd7 = 1.1.010111 = -0b1.010111×2⁰ = -1.359375
0xd8 = 1.1.011000 = -0b1.011000×2⁰ = -1.375
0xd9 = 1.1.011001 = -0b1.011001×2⁰ = -1.390625
0xda = 1.1.011010 = -0b1.011010×2⁰ = -1.40625
0xdb = 1.1.011011 = -0b1.011011×2⁰ = -1.421875
0xdc = 1.1.011100 = -0b1.011100×2⁰ = -1.4375
0xdd = 1.1.011101 = -0b1.011101×2⁰ = -1.453125
0xde = 1.1.011110 = -0b1.011110×2⁰ = -1.46875
0xdf = 1.1.011111 = -0b1.011111×2⁰ = -1.484375
0xe0 = 1.1.100000 = -0b1.100000×2⁰ = -1.5
0xe1 = 1.1.100001 = -0b1.100001×2⁰ = -1.515625
0xe2 = 1.1.100010 = -0b1.100010×2⁰ = -1.53125
0xe3 = 1.1.100011 = -0b1.100011×2⁰ = -1.546875
0xe4 = 1.1.100100 = -0b1.100100×2⁰ = -1.5625
0xe5 = 1.1.100101 = -0b1.100101×2⁰ = -1.578125
0xe6 = 1.1.100110 = -0b1.100110×2⁰ = -1.59375
0xe7 = 1.1.100111 = -0b1.100111×2⁰ = -1.609375
0xe8 = 1.1.101000 = -0b1.101000×2⁰ = -1.625
0xe9 = 1.1.101001 = -0b1.101001×2⁰ = -1.640625
0xea = 1.1.101010 = -0b1.101010×2⁰ = -1.65625
0xeb = 1.1.101011 = -0b1.101011×2⁰ = -1.671875
0xec = 1.1.101100 = -0b1.101100×2⁰ = -1.6875
0xed = 1.1.101101 = -0b1.101101×2⁰ = -1.703125
0xee = 1.1.101110 = -0b1.101110×2⁰ = -1.71875
0xef = 1.1.101111 = -0b1.101111×2⁰ = -1.734375
0xf0 = 1.1.110000 = -0b1.110000×2⁰ = -1.75
0xf1 = 1.1.110001 = -0b1.110001×2⁰ = -1.765625
0xf2 = 1.1.110010 = -0b1.110010×2⁰ = -1.78125
0xf3 = 1.1.110011 = -0b1.110011×2⁰ = -1.796875
0xf4 = 1.1.110100 = -0b1.110100×2⁰ = -1.8125
0xf5 = 1.1.110101 = -0b1.110101×2⁰ = -1.828125
0xf6 = 1.1.110110 = -0b1.110110×2⁰ = -1.84375
0xf7 = 1.1.110111 = -0b1.110111×2⁰ = -1.859375
0xf8 = 1.1.111000 = -0b1.111000×2⁰ = -1.875
0xf9 = 1.1.111001 = -0b1.111001×2⁰ = -1.890625
0xfa = 1.1.111010 = -0b1.111010×2⁰ = -1.90625
0xfb = 1.1.111011 = -0b1.111011×2⁰ = -1.921875
0xfc = 1.1.111100 = -0b1.111100×2⁰ = -1.9375
0xfd = 1.1.111101 = -0b1.111101×2⁰ = -1.953125
0xfe = 1.1.111110 = -0b1.111110×2⁰ = -1.96875
0xff = 1.1.111111 = -Inf

References

- [1] PyTorch authors. Pytorch torchtext package: `_t5_multi_head_attention_forward` . <https://github.com/pytorch/text/blob/a933cbe5a008bc2cb61d985cf5864069194157eb/torchtext/prototype/models/t5/modules.py#L236>.
- [2] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*, chapter 6.2.2.3 Softmax Units for Multinoulli Output Distributions, pages 180–184. MIT Press, 2016.
- [3] Google. Jax lax package: `_float_to_int_for_sort` . https://github.com/google/jax/blob/fc5960f2b8b7a0ef74dbae4e27c5c08ff1564cff/jax/_src/lax/lax.py#L3934.
- [4] W. Kahan. Branch cuts for complex elementary functions or much ado about nothing’s sign bit. *Inst. Math. Appl. Conf. Ser. New Ser.*, 1987.
- [5] W. Kahan and J. W. Thomas. Augmenting a programming language with complex arithmetic. Technical report, EECS Department, University of California, Berkeley, 1991.
- [6] 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). Technical report, opencompute.org, 2023. <https://www.opencompute.org/documents/ocp-8-bit-floating-point-specification-ofp8-revision-1-0-2023-06-20-pdf>.
- [7] B. Nouné, P. Jones, D. Justus, D. Masters, and C. Luschi. Adaptive loss scaling for mixed precision training. Technical report, arXiv cs.LG, 2019. <https://arxiv.org/abs/1910.12385>.
- [8] B. Nouné, P. Jones, D. Justus, D. Masters, and C. Luschi. 8-bit numerical formats for deep neural networks. Technical report, 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.