

Overview of New Floating-Point Standards: IEEE P3109 (May 2024 interim report) and Open Compute Project OCP8/MX 1.0

Mantas Mikaitis

School of Computer Science, University of Leeds, Leeds, UK

Scientific Computation Group Seminar

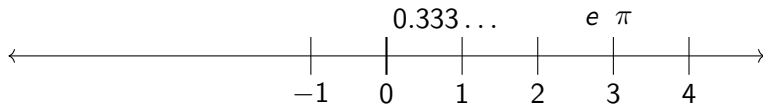
22 May, 2025

Slides: mmikaitis.github.io

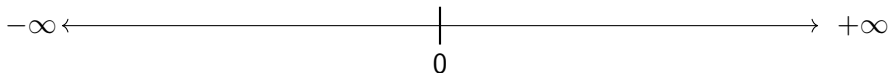


Working with real numbers on computers

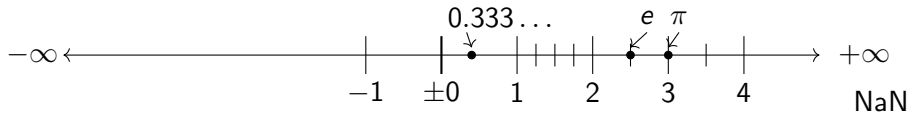
\mathbb{R}



$\mathbb{R} \cup \{-\infty, +\infty\}$



\mathbb{F}



What is floating point? Picking subsets of reals

A floating-point system $\mathbb{F} \subset \mathbb{R} \cup \{\pm\infty, -0, \text{NaN}\}$ is described with $\beta, t, e_{\min}, e_{\max}$ with elements

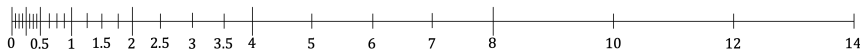
$$\pm m \times \beta^{e-t+1}.$$

Virtually all computers have $\beta = 2$ (binary FP).

Here t is precision (binary digits), $e_{\min} \leq e \leq e_{\max}$ an exponent, $m \leq \beta^t - 1$ a significand ($m, t, e \in \mathbb{Z}$).

Toy FP system

Below: the positive numbers in $\mathbb{F}(\beta = 2, t = 3, e_{\min} = -2, e_{\max} = 3)$.



IEEE Standard for Floating-Point Arithmetic

Part 1: IEEE 754 Standard for Floating-Point Arithmetic

IEEE Computer Society

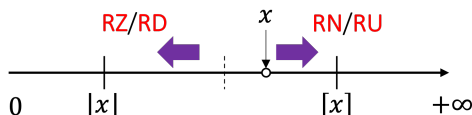
Developed by the
Microprocessor Standards Committee

IEEE
3 Park Avenue
New York, NY 10016-5997
USA

IEEE Std 754™-2019

IEEE 754 standard FP arithmetic: rounding

- Round-to-nearest (**RN**) (ties even)
- Round-toward-zero (**RZ**)
- Round-down (**RD**)
- Round-up (**RU**)



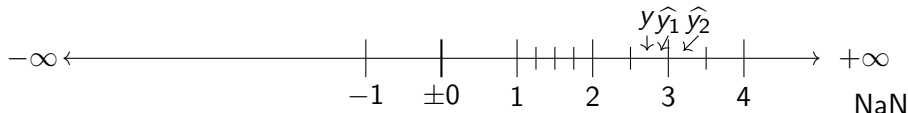
Use of rounding modes

RN is a default. *Directed modes* used for special cases, such as **interval arithmetic**.

Why a standard is needed?

- Rounding modes: result lies between two FP numbers—return one of them consistently.
- *Correct rounding* (directed): exact result y and approximation \hat{y} are surrounded by the same two FP numbers, consistently.
- *Correct rounding* (nearest): nearest FP number to y is also a nearest FP number to \hat{y} .

\mathbb{F}

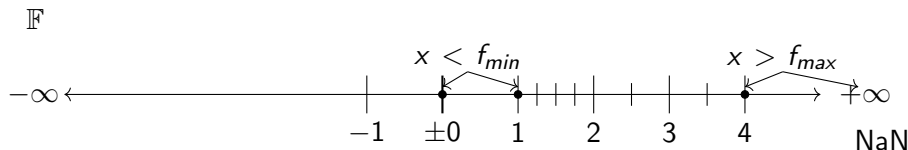


Accuracy is therefore not a consideration in IEEE 754—compliant operations are as accurate as precision allows.

Why a standard is needed? Overflows and underflows

The standard sets out behaviour for edge cases:

- Overflow is when rounding produces $\pm\infty$.
- Underflow occurs when the rounded result is between zero and smallest magnitude normal FP value.



Why a standard is needed? Operations

IEEE 754 requires *correctly rounded*

- $+$, $-$, \times , \div
- Square root
- FMA: $a \times b + c$

However, it does not require x^y and mathematical functions \exp , \sin , \cos , \log , and others.

Nevertheless, it recommends providing them, and they *shall be* correctly rounded—most mathematical function libraries do not assure this today.

Why a standard is needed? Invalid operations

A NaN is generated by

- $0/0$
- $0 \times \infty$
- ∞/∞
- $(+\infty) - (-\infty)$
- $\sqrt{-1}$

Once a NaN is generated, it propagates so the user sees it in the result, or causes an exception from the floating-point unit.

Why a standard is needed? Infinities

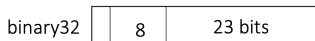
∞ obeys mathematical rules

- $\infty + \infty = \infty$
- $(-1) \times \infty = -\infty$
- $x/\infty = 0$ where x is a finite number

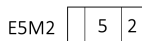
Floating-point format encoding

Numbers are held in memory using bits (convenient when $\beta = 2$).

Main IEEE 754 formats (**double**, **single**, **half**):



Some ubiquitous non-standard formats:



Floating point formats and their ranges

Format	precision	min pos. (f_{min})	max pos. (f_{max})	u
binary64	53	2^{-1022}	$\sim 1.798 \times 10^{308}$	2^{-53}
binary32	24	2^{-126}	$\sim 3.403 \times 10^{38}$	2^{-24}
tf32 (19-bit)	11	2^{-126}	$\sim 3.401 \times 10^{38}$	2^{-11}
bfloat16	8	2^{-126}	$\sim 3.389 \times 10^{38}$	2^{-8}
binary16	11	2^{-14}	65504	2^{-11}
fp8-E4M3	4	2^{-6}	448	2^{-4}
fp8-E5M2	3	2^{-14}	57344	2^{-3}
fp6-E2M3	4	2^0	7.5	2^{-4}
fp6-E3M2	3	2^{-2}	28	2^{-3}
fp4-E2M1	2	2^0	6	2^{-2}

Limitations of standardisation: ordering of operations

Some of the known properties of floating point:

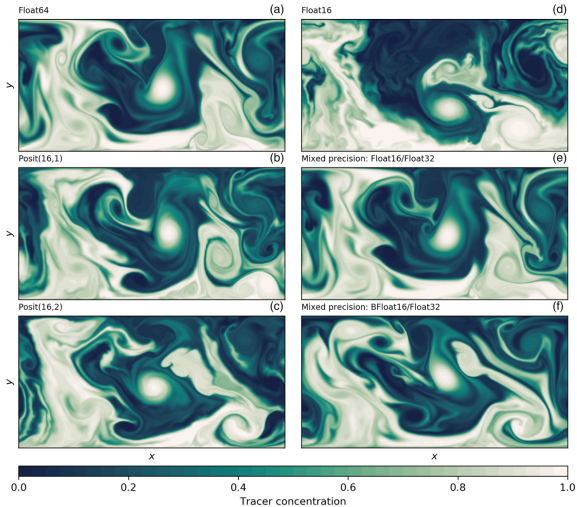
- ✓ Commutativity: $a \times b = b \times a$.
- × Associativity: $(a + (b + c)) \neq (a + b) + c$.
- × Distributivity: $(a \times (b + c)) \neq (a \times b) + (a \times c)$.

Even though $+$ is correctly rounded, IEEE 754 does not mandate the order in $\sum_{i=0}^n x_i$ —results implementation dependent.

But that is OK: ordering is in programmer's control, and two systems can match their implementations if needed.

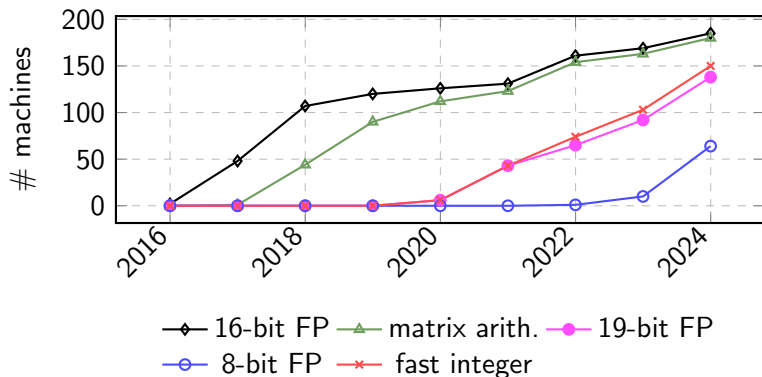
(except when race conditions and thread scheduling change ordering)

Part 2: Low-Precision Floating Point (made for AI, but people use it for scientific computation)



Shallow water model simulation [Klower et al. 2020].

Low-precision floating point on the TOP500



Devices counted: P100, V100, A100, H100, MI210, MI250X, MI300X, Intel Data Center GPU, from <https://www.top500.org>.

With NVIDIA Blackwell 4/6-bit FP will appear.

Mixed-precision matrix multipliers

Non-standard low-precision formats are available in matrix multiply form.

$$\begin{array}{ccccccc} D & = & C & + & A & \times & B, \\ \underbrace{\begin{bmatrix} \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \end{bmatrix}}_{\text{binary16 or binary32}} & = & \underbrace{\begin{bmatrix} \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \end{bmatrix}}_{\text{binary16 or binary32}} & + & \underbrace{\begin{bmatrix} \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \end{bmatrix}}_{\text{8-bit FP}} & \times & \underbrace{\begin{bmatrix} \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \end{bmatrix}}_{\text{8-bit FP}} \end{array}$$

Example above is 4×4 , but dimensions differ across architectures.

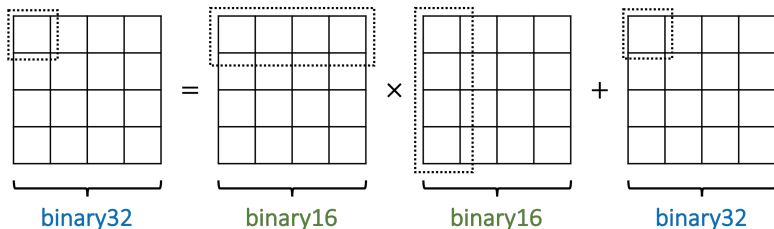
Mixed-precision matrix multipliers

Architecture	Input format	Accumulation format
NVIDIA PTX ISA	fp8-E5M2	binary32
	fp8-E4M3	binary32
	binary16	binary16
	binary16	binary32
	bfloat16	binary32
	19-bit FP	binary32
AMD MI300 ISA	fp8-E5M2	binary32
	fp8-E4M3	binary32
	binary16	binary32
	bfloat16	binary32
	19-bit FP	binary32

NVIDIA Blackwell throughputs (FLOPS)

fp8 (9×10^{15}) fp16 (4.5×10^{15}) fp64 (0.04×10^{15}).

Nonstandard features of matrix multipliers



$$a_{11} \times b_{11} + a_{12} \times b_{21} + a_{13} \times b_{31} + a_{14} \times b_{41}$$

- We are used to **normalize-round after each op.**
- In hardware it is not necessarily the case.
- Normalize at the end? Savings in circuit area and latency.
- Round or drop bits? Savings.
- Accumulate in higher precision? Can do with 1-bit granularity.

Why test the mathematical hardware

IEEE 754 does not define strict rules on dot products:

“Implementations may associate in any order or evaluate in any wider format.”

- Implementations might differ.
- Not documented in detail by vendors.
- Massive speed means we are using MMAs in other areas.
- No control over ordering means we may not be able to match two systems.

For now we test to determine features.

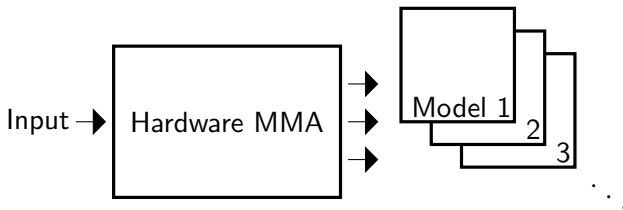
How do we test mathematical hardware

[Fasi, Higham, Mikaitis, Pranesh, 2021]

Technique goes back to software called Paranoia from the 1980s.

MMAs more complicated than $+$, \times , \div

Find floating-point **inputs that will yield different outputs on different hardware**.



Our main findings from testing

$$a_{11} \times b_{11} + a_{12} \times b_{21} + a_{13} \times b_{31} + a_{14} \times b_{41}$$

1-bit difference in accumulators of NVIDIA V100 (2018) and T4/A100 (2020).

Normalization at the end, not at intermediate adds.

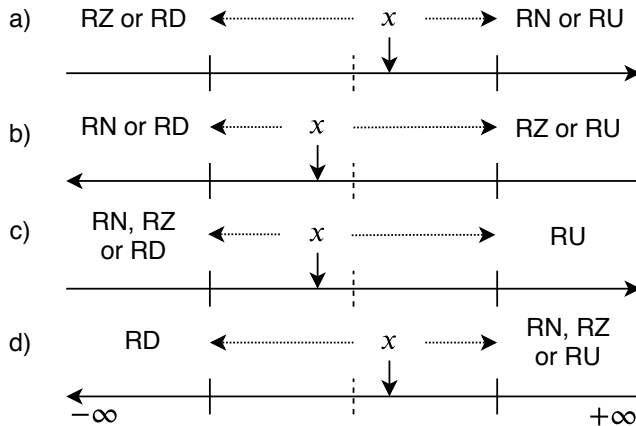
Rounding is not to nearest.

Nonmonotonicity: increase one input, do not change order, dot product decreases. See [\[Mikaitis, 2024\]](#).

No fixed order.

Currently a manual process

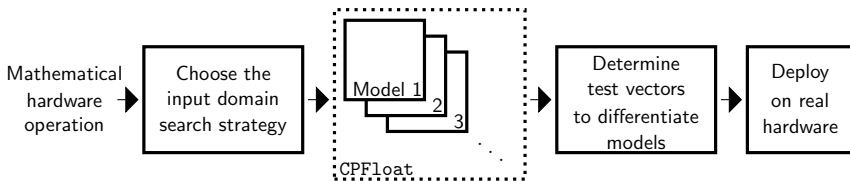
Consider rounding:



New, more systematic approach

3-year EPSRC project: 2025-2028. Intel and Argonne National Lab as partners.

- Make the search for testing vectors automatic, or at least partially.
- Remove the need for very specialized floating-point knowledge.
- Make a library of behaviours and maintain as new hardware comes.



Simulating custom precision with CPFLOAT in MATLAB/Octave

Paper in ACM TOMS provides details. [\[Fasi and Mikaitis, 2023\]](#)

<https://github.com/north-numerical-computing/cpfloat>

Example use in MATLAB:

```
>> input.format = 'q43';
>> input.subnormal = 0;
>> accum.format = 'binary16';
>> accum.subnormal = 0;
>> cpfloat(pi, input)
ans =
    3.2500000000000000
>> cpfloat(cpfloat(
    cpfloat(pi,input)*cpfloat(pi,input),accum)+0.5,accum)
ans =
   11.0625000000000000
```


Part 3: Open Compute Project (OCP) standards

Open Compute Project floating-point standards

Two standards:

- Standard 1: OCP 8-bit Floating Point Specification
- Standard 2: OCP Microscaling Formats (MX) Specification

Version 1.0 released June and September 2023.



OPEN
Compute Project

OCP 8-bit Floating Point Specification (OFP8)

Revision 1.0

Date Submitted: May 26, 2023

Date Approved: June 20, 2023

Author: Paulius Mickevicius and Stuart Oberman, NVIDIA

Author: Pradeep Dubey, Marius Cornea and Andres Rodriguez, Intel

Author: Ian Bratt and Richard Grisenthwaite, Arm

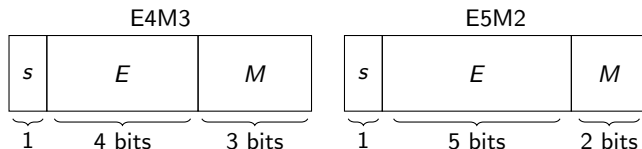
Author: Norm Jouppi, Chiachen Chou and Amber Huffman, Google

Author: Michael Schulte and Ralph Wittig, AMD

Author: Dharmesh Jani and Summer Deng, Meta

OCP floating-point standard 1

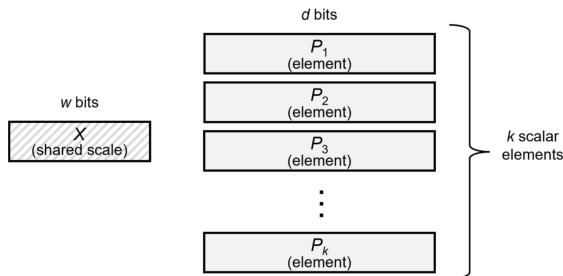
Key aspects in the Standard 1



- Defines two formats: OFP8. E4M3 no $\pm\infty$, one NaN.
- Defines conversion from higher precision formats (binary32, binary16, bfloat16) to OFP8.
- Conversion includes
 - Round to nearest (no other rounding modes required)
 - Saturation mode: after rounding, if $> f_{max}$, return f_{max} instead of ∞ .

Arithmetic operations are not in scope of standard 1.

OCP standard 2 (screenshot from the MX standard)



- Characterized by scale (X) type, data (P_i) type, and block size k
- Layout in memory not prescribed.
- $w + kd$ bits required. $w = 8$, $k = 32$ for all configurations.
- Value i encodes $X \times P_i$.

Benefits of shared scale: cheaply extend dynamic range (extra w/k bits per element).

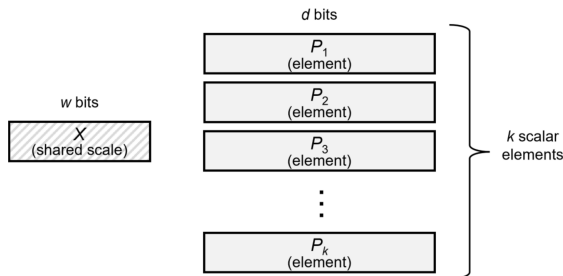
OCP standard 2 (screenshot from the MX standard)

Format Name	Element Data Type	Element Bits (d)	Scaling Block Size (k)	Scale Data Type	Scale Bits (w)
MXFP8	FP8 (E5M2)	8	32	E8M0	8
	FP8 (E4M3)				
MXFP6	FP6 (E3M2)	6	32	E8M0	8
	FP6 (E2M3)				
MXFP4	FP4 (E2M1)	4	32	E8M0	8
MXINT8	INT8	8	32	E8M0	8

INT8 is a fixed-point format encoding values $(-2, 2)$ in steps of 2^{-6} .

Conversion from standard vectors of size 32, to the MX formats, must be provided.

OCP standard 2 (screenshot from the MX standard)



Dot product of two MX vectors must be provided

$$C = \text{Dot}(A, B) = X^A X^B \sum_{i=1}^{32} (P_i^A \times P_i^B)$$

(*“internal precision of the dot product and order of operations is implementation-defined”*)

Part 4: IEEE P3109 Standard for Arithmetic Formats for Machine Learning (Oct. 2024 state)

Standard for Arithmetic Formats for Machine Learning

New IEEE standard for computer arithmetic for AI is in progress. We are also actively participating in it: meeting fortnightly.

Standardises:

- Small floating-point subsets of reals,
- encoding in 8-bit words,
- rounding behaviour,
- arithmetic operations needed in AI workloads,
- conversion to/from 754,
- exception handling.

Interim report available: <http://bit.ly/42gPWcy>

IEEE floating-point standardisation work: IEEE P3109

Current draft outlines these key aspects:

- Defines formats as binary KpP with $K = 8$ and $P = \{1, 2, 3, 4, 5, 6, 7\}$.
- No -0 (would have been 0x80)
- Only one NaN (0x80) - note IEEE 754 numbers have many bit patterns for NaNs.
- $\pm\infty$ 0x7F and 0xFF.
- Saturation mode: on overflow, return maximum finite value.
- Several rounding modes, operations.

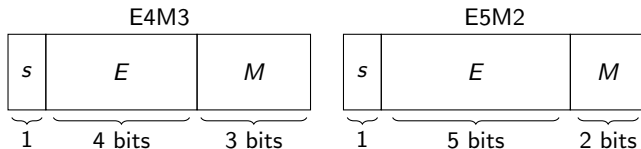
Format	minSubnormal	maxSubnormal	minNormal	maxNormal
binary8p1	N/A	N/A	1×2^{-62}	1×2^{63}
binary8p2	1×2^{-32}	1×2^{-32}	1×2^{-31}	1×2^{31}
binary8p3	1×2^{-17}	$3/2 \times 2^{-16}$	1×2^{-15}	$3/2 \times 2^{15}$
binary8p4	1×2^{-10}	$7/4 \times 2^{-8}$	1×2^{-7}	$7/4 \times 2^7$
binary8p5	1×2^{-7}	$15/8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

IEEE P33300 8bit format on one slide

C.4 Value Table: P4, $e_{\min} = -7$, $e_{\max} = 7$

$0x00 = 0.0000.000 = 0.0$	$0x40 = 0.1000.000 = +0b1.000 \times 2^0 = 1.0$	$0x80 = 1.0000.000 = +Inf$	$0xc0 = 1.0000.000 = -0b1.000 \times 2^0 = -1.0$
$0x01 = 0.0000.001 = +0b0.001 \times 2^{-7} = 0.0008765625$	$0x41 = 0.1000.001 = +0b1.001 \times 2^0 = 1.125$	$0x81 = 1.0000.001 = -0b0.001 \times 2^{-7} = -0.0008765625$	$0xc1 = 1.0000.001 = -0b1.001 \times 2^0 = -1.125$
$0x02 = 0.0000.010 = +0b0.010 \times 2^{-7} = 0.00173125$	$0x42 = 0.1000.010 = +0b1.010 \times 2^0 = 1.25$	$0x82 = 1.0000.010 = -0b0.010 \times 2^{-7} = -0.00173125$	$0xc2 = 1.0000.010 = -0b1.010 \times 2^0 = -1.25$
$0x03 = 0.0000.011 = +0b0.011 \times 2^{-7} = 0.002596875$	$0x43 = 0.1000.011 = +0b1.011 \times 2^0 = 1.375$	$0x83 = 1.0000.011 = -0b0.011 \times 2^{-7} = -0.002596875$	$0xc3 = 1.0000.011 = -0b1.011 \times 2^0 = -1.375$
$0x04 = 0.0000.100 = +0b0.100 \times 2^{-7} = 0.00390625$	$0x44 = 0.1000.100 = +0b1.100 \times 2^0 = 1.5$	$0x84 = 1.0000.100 = -0b0.100 \times 2^{-7} = -0.00390625$	$0xc4 = 1.0000.100 = -0b1.100 \times 2^0 = -1.5$
$0x05 = 0.0000.101 = +0b0.101 \times 2^{-7} = 0.0048828125$	$0x45 = 0.1000.101 = +0b1.101 \times 2^0 = 1.625$	$0x85 = 1.0000.101 = -0b0.101 \times 2^{-7} = -0.0048828125$	$0xc5 = 1.0000.101 = -0b1.101 \times 2^0 = -1.625$
$0x06 = 0.0000.110 = +0b0.110 \times 2^{-7} = 0.005859375$	$0x46 = 0.1000.110 = +0b1.110 \times 2^0 = 1.75$	$0x86 = 1.0000.110 = -0b0.110 \times 2^{-7} = -0.005859375$	$0xc6 = 1.0000.110 = -0b1.110 \times 2^0 = -1.75$
$0x07 = 0.0000.111 = +0b0.111 \times 2^{-7} = 0.0068359375$	$0x47 = 0.1000.111 = +0b1.111 \times 2^0 = 1.875$	$0x87 = 1.0000.111 = -0b0.111 \times 2^{-7} = -0.0068359375$	$0xc7 = 1.0000.111 = -0b1.111 \times 2^0 = -1.875$
$0x08 = 0.0001.000 = +0b1.000 \times 2^{-7} = 0.0078125$	$0x48 = 0.1001.000 = +0b1.000 \times 2^1 = 2.0$	$0x88 = 1.0001.000 = -0b1.000 \times 2^{-7} = -0.0078125$	$0xc8 = 1.0001.000 = -0b1.000 \times 2^1 = -2.0$
$0x09 = 0.0001.001 = +0b1.001 \times 2^{-7} = 0.007890625$	$0x49 = 0.1001.001 = +0b1.001 \times 2^1 = 2.125$	$0x89 = 1.0001.001 = -0b1.001 \times 2^{-7} = -0.007890625$	$0xc9 = 1.0001.001 = -0b1.001 \times 2^1 = -2.125$
$0x0a = 0.0001.010 = +0b1.010 \times 2^{-7} = 0.00796625$	$0x4a = 0.1001.010 = +0b1.010 \times 2^1 = 2.25$	$0x8a = 1.0001.010 = -0b1.010 \times 2^{-7} = -0.00796625$	$0xca = 1.0001.010 = -0b1.010 \times 2^1 = -2.25$
$0x0b = 0.0001.011 = +0b1.011 \times 2^{-7} = 0.010741875$	$0x4b = 0.1001.011 = +0b1.011 \times 2^1 = 2.375$	$0x8b = 1.0001.011 = -0b1.011 \times 2^{-7} = -0.010741875$	$0xcb = 1.0001.011 = -0b1.011 \times 2^1 = -2.375$
$0x0c = 0.0001.100 = +0b1.100 \times 2^{-7} = 0.01171875$	$0x4c = 0.1001.100 = +0b1.100 \times 2^1 = 2.5$	$0x8c = 1.0001.100 = -0b1.100 \times 2^{-7} = -0.01171875$	$0xcc = 1.0001.100 = -0b1.100 \times 2^1 = -2.5$
$0x0d = 0.0001.101 = +0b1.101 \times 2^{-7} = 0.0126953125$	$0x4d = 0.1001.101 = +0b1.101 \times 2^1 = 3.125$	$0x8d = 1.0001.101 = -0b1.101 \times 2^{-7} = -0.0126953125$	$0xcd = 1.0001.101 = -0b1.101 \times 2^1 = -3.125$
$0x0e = 0.0001.110 = +0b1.110 \times 2^{-7} = 0.013671875$	$0x4e = 0.1001.110 = +0b1.110 \times 2^1 = 3.25$	$0x8e = 1.0001.110 = -0b1.110 \times 2^{-7} = -0.013671875$	$0xce = 1.0001.110 = -0b1.110 \times 2^1 = -3.25$
$0x0f = 0.0001.111 = +0b1.111 \times 2^{-7} = 0.0146484375$	$0x4f = 0.1001.111 = +0b1.111 \times 2^1 = 3.75$	$0x8f = 1.0001.111 = -0b1.111 \times 2^{-7} = -0.0146484375$	$0xcf = 1.0001.111 = -0b1.111 \times 2^1 = -3.75$
$0x10 = 0.0010.000 = +0b1.000 \times 2^{-6} = 0.015625$	$0x50 = 0.1010.000 = +0b1.000 \times 2^2 = 4.0$	$0x90 = 1.0010.000 = -0b1.000 \times 2^{-6} = -0.015625$	$0xd0 = 1.0010.000 = -0b1.000 \times 2^2 = -4.0$
$0x11 = 0.0010.001 = +0b1.001 \times 2^{-6} = 0.01578125$	$0x51 = 0.1010.001 = +0b1.001 \times 2^2 = 4.5$	$0x91 = 1.0010.001 = -0b1.001 \times 2^{-6} = -0.01578125$	$0xd1 = 1.0010.001 = -0b1.001 \times 2^2 = -4.5$
$0x12 = 0.0010.010 = +0b1.010 \times 2^{-6} = 0.01963125$	$0x52 = 0.1010.010 = +0b1.010 \times 2^2 = 5.0$	$0x92 = 1.0010.010 = -0b1.010 \times 2^{-6} = -0.01963125$	$0xd2 = 1.0010.010 = -0b1.010 \times 2^2 = -5.0$
$0x13 = 0.0010.011 = +0b1.011 \times 2^{-6} = 0.021484375$	$0x53 = 0.1010.011 = +0b1.011 \times 2^2 = 5.5$	$0x93 = 1.0010.011 = -0b1.011 \times 2^{-6} = -0.021484375$	$0xd3 = 1.0010.011 = -0b1.011 \times 2^2 = -5.5$
$0x14 = 0.0010.100 = +0b1.100 \times 2^{-6} = 0.0234375$	$0x54 = 0.1010.100 = +0b1.100 \times 2^2 = 6.0$	$0x94 = 1.0010.100 = -0b1.100 \times 2^{-6} = -0.0234375$	$0xd4 = 1.0010.100 = -0b1.100 \times 2^2 = -6.0$
$0x15 = 0.0010.101 = +0b1.101 \times 2^{-6} = 0.025390625$	$0x55 = 0.1010.101 = +0b1.101 \times 2^2 = 6.5$	$0x95 = 1.0010.101 = -0b1.101 \times 2^{-6} = -0.025390625$	$0xd5 = 1.0010.101 = -0b1.101 \times 2^2 = -6.5$
$0x16 = 0.0010.110 = +0b1.110 \times 2^{-6} = 0.02734375$	$0x56 = 0.1010.110 = +0b1.110 \times 2^2 = 7.0$	$0x96 = 1.0010.110 = -0b1.110 \times 2^{-6} = -0.02734375$	$0xd6 = 1.0010.110 = -0b1.110 \times 2^2 = -7.0$
$0x17 = 0.0010.111 = +0b1.111 \times 2^{-6} = 0.029296875$	$0x57 = 0.1010.111 = +0b1.111 \times 2^2 = 7.5$	$0x97 = 1.0010.111 = -0b1.111 \times 2^{-6} = -0.029296875$	$0xd7 = 1.0010.111 = -0b1.111 \times 2^2 = -7.5$
$0x18 = 0.0011.000 = +0b1.000 \times 2^{-5} = 0.03125$	$0x58 = 0.1011.000 = +0b1.000 \times 2^3 = 8.0$	$0x98 = 1.0011.000 = -0b1.000 \times 2^{-5} = -0.03125$	$0xd8 = 1.0011.000 = -0b1.000 \times 2^3 = -8.0$
$0x19 = 0.0011.001 = +0b1.001 \times 2^{-5} = 0.0315625$	$0x59 = 0.1011.001 = +0b1.001 \times 2^3 = 8.5$	$0x99 = 1.0011.001 = -0b1.001 \times 2^{-5} = -0.0315625$	$0xd9 = 1.0011.001 = -0b1.001 \times 2^3 = -8.5$
$0x1a = 0.0011.010 = +0b1.010 \times 2^{-5} = 0.0390625$	$0x5a = 0.1011.010 = +0b1.010 \times 2^3 = 10.0$	$0x9a = 1.0011.010 = -0b1.010 \times 2^{-5} = -0.0390625$	$0xda = 1.0011.010 = -0b1.010 \times 2^3 = -10.0$
$0x1b = 0.0011.011 = +0b1.011 \times 2^{-5} = 0.04296875$	$0x5b = 0.1011.011 = +0b1.011 \times 2^3 = 11.0$	$0x9b = 1.0011.011 = -0b1.011 \times 2^{-5} = -0.04296875$	$0xdb = 1.0011.011 = -0b1.011 \times 2^3 = -11.0$
$0x1c = 0.0011.100 = +0b1.100 \times 2^{-5} = 0.046875$	$0x5c = 0.1011.100 = +0b1.100 \times 2^3 = 12.0$	$0x9c = 1.0011.100 = -0b1.100 \times 2^{-5} = -0.046875$	$0xdc = 1.0011.100 = -0b1.100 \times 2^3 = -12.0$
$0x1d = 0.0011.101 = +0b1.101 \times 2^{-5} = 0.05078125$	$0x5d = 0.1011.101 = +0b1.101 \times 2^3 = 13.0$	$0x9d = 1.0011.101 = -0b1.101 \times 2^{-5} = -0.05078125$	$0xdd = 1.0011.101 = -0b1.101 \times 2^3 = -13.0$
$0x1e = 0.0011.110 = +0b1.110 \times 2^{-5} = 0.0546875$	$0x5e = 0.1011.110 = +0b1.110 \times 2^3 = 14.0$	$0x9e = 1.0011.110 = -0b1.110 \times 2^{-5} = -0.0546875$	$0xde = 1.0011.110 = -0b1.110 \times 2^3 = -14.0$
$0x1f = 0.0011.111 = +0b1.111 \times 2^{-5} = 0.05859375$	$0x5f = 0.1011.111 = +0b1.111 \times 2^3 = 15.0$	$0x9f = 1.0011.111 = -0b1.111 \times 2^{-5} = -0.05859375$	$0xdf = 1.0011.111 = -0b1.111 \times 2^3 = -15.0$
$0x20 = 0.0100.000 = +0b1.000 \times 2^{-4} = 0.0625$	$0x60 = 0.1100.000 = +0b1.000 \times 2^4 = 16.0$	$0xa0 = 1.0100.000 = -0b1.000 \times 2^{-4} = -0.0625$	$0xe0 = 1.0100.000 = -0b1.000 \times 2^4 = -16.0$
$0x21 = 0.0100.001 = +0b1.001 \times 2^{-4} = 0.0703125$	$0x61 = 0.1100.001 = +0b1.001 \times 2^4 = 16.5$	$0xa1 = 1.0100.001 = -0b1.001 \times 2^{-4} = -0.0703125$	$0xe1 = 1.0100.001 = -0b1.001 \times 2^4 = -16.5$
$0x22 = 0.0100.010 = +0b1.010 \times 2^{-4} = 0.078125$	$0x62 = 0.1100.010 = +0b1.010 \times 2^4 = 20.0$	$0xa2 = 1.0100.010 = -0b1.010 \times 2^{-4} = -0.078125$	$0xe2 = 1.0100.010 = -0b1.010 \times 2^4 = -20.0$
$0x23 = 0.0100.011 = +0b1.011 \times 2^{-4} = 0.0859375$	$0x63 = 0.1100.011 = +0b1.011 \times 2^4 = 22.0$	$0xa3 = 1.0100.011 = -0b1.011 \times 2^{-4} = -0.0859375$	$0xe3 = 1.0100.011 = -0b1.011 \times 2^4 = -22.0$
$0x24 = 0.0100.100 = +0b1.100 \times 2^{-4} = 0.09375$	$0x64 = 0.1100.100 = +0b1.100 \times 2^4 = 24.0$	$0xa4 = 1.0100.100 = -0b1.100 \times 2^{-4} = -0.09375$	$0xe4 = 1.0100.100 = -0b1.100 \times 2^4 = -24.0$
$0x25 = 0.0100.101 = +0b1.101 \times 2^{-4} = 0.101625$	$0x65 = 0.1100.101 = +0b1.101 \times 2^4 = 26.0$	$0xa5 = 1.0100.101 = -0b1.101 \times 2^{-4} = -0.101625$	$0xe5 = 1.0100.101 = -0b1.101 \times 2^4 = -26.0$
$0x26 = 0.0100.110 = +0b1.110 \times 2^{-4} = 0.109375$	$0x66 = 0.1100.110 = +0b1.110 \times 2^4 = 28.0$	$0xa6 = 1.0100.110 = -0b1.110 \times 2^{-4} = -0.109375$	$0xe6 = 1.0100.110 = -0b1.110 \times 2^4 = -28.0$
$0x27 = 0.0100.111 = +0b1.111 \times 2^{-4} = 0.1171875$	$0x67 = 0.1100.111 = +0b1.111 \times 2^4 = 30.0$	$0xa7 = 1.0100.111 = -0b1.111 \times 2^{-4} = -0.1171875$	$0xe7 = 1.0100.111 = -0b1.111 \times 2^4 = -30.0$
$0x28 = 0.0101.000 = +0b1.000 \times 2^{-3} = 0.125$	$0x68 = 0.1101.000 = +0b1.000 \times 2^5 = 32.0$	$0xa8 = 1.0101.000 = -0b1.000 \times 2^{-3} = -0.125$	$0xe8 = 1.0101.000 = -0b1.000 \times 2^5 = -32.0$
$0x29 = 0.0101.001 = +0b1.001 \times 2^{-3} = 0.140625$	$0x69 = 0.1101.001 = +0b1.001 \times 2^5 = 36.0$	$0xa9 = 1.0101.001 = -0b1.001 \times 2^{-3} = -0.140625$	$0xe9 = 1.0101.001 = -0b1.001 \times 2^5 = -36.0$
$0x2a = 0.0101.010 = +0b1.010 \times 2^{-3} = 0.15625$	$0x6a = 0.1101.010 = +0b1.010 \times 2^5 = 40.0$	$0xaa = 1.0101.010 = -0b1.010 \times 2^{-3} = -0.15625$	$0xea = 1.0101.010 = -0b1.010 \times 2^5 = -40.0$
$0x2b = 0.0101.011 = +0b1.011 \times 2^{-3} = 0.171875$	$0x6b = 0.1101.011 = +0b1.011 \times 2^5 = 46.0$	$0xab = 1.0101.011 = -0b1.011 \times 2^{-3} = -0.171875$	$0xeb = 1.0101.011 = -0b1.011 \times 2^5 = -46.0$
$0x2c = 0.0101.100 = +0b1.100 \times 2^{-3} = 0.1875$	$0x6c = 0.1101.100 = +0b1.100 \times 2^5 = 48.0$	$0xac = 1.0101.100 = -0b1.100 \times 2^{-3} = -0.1875$	$0xec = 1.0101.100 = -0b1.100 \times 2^5 = -48.0$
$0x2d = 0.0101.101 = +0b1.101 \times 2^{-3} = 0.203125$	$0x6d = 0.1101.101 = +0b1.101 \times 2^5 = 52.0$	$0xad = 1.0101.101 = -0b1.101 \times 2^{-3} = -0.203125$	$0xed = 1.0101.101 = -0b1.101 \times 2^5 = -52.0$
$0x2e = 0.0101.110 = +0b1.110 \times 2^{-3} = 0.21875$	$0x6e = 0.1101.110 = +0b1.110 \times 2^5 = 60.0$	$0xae = 1.0101.110 = -0b1.110 \times 2^{-3} = -0.21875$	$0xee = 1.0101.110 = -0b1.110 \times 2^5 = -60.0$
$0x2f = 0.0101.111 = +0b1.111 \times 2^{-3} = 0.234375$	$0x6f = 0.1101.111 = +0b1.111 \times 2^5 = 70.0$	$0xaf = 1.0101.111 = -0b1.111 \times 2^{-3} = -0.234375$	$0xef = 1.0101.111 = -0b1.111 \times 2^5 = -70.0$
$0x30 = 0.0110.000 = +0b1.000 \times 2^{-2} = 0.25$	$0x70 = 0.1110.000 = +0b1.000 \times 2^6 = 64.0$	$0xb0 = 1.0110.000 = -0b1.000 \times 2^{-2} = -0.25$	$0xf0 = 1.0110.000 = -0b1.000 \times 2^6 = -64.0$
$0x31 = 0.0110.001 = +0b1.001 \times 2^{-2} = 0.261875$	$0x71 = 0.1110.001 = +0b1.001 \times 2^6 = 72.0$	$0xb1 = 1.0110.001 = -0b1.001 \times 2^{-2} = -0.261875$	$0xf1 = 1.0110.001 = -0b1.001 \times 2^6 = -72.0$
$0x32 = 0.0110.010 = +0b1.010 \times 2^{-2} = 0.3125$	$0x72 = 0.1110.010 = +0b1.010 \times 2^6 = 80.0$	$0xb2 = 1.0110.010 = -0b1.010 \times 2^{-2} = -0.3125$	$0xf2 = 1.0110.010 = -0b1.010 \times 2^6 = -80.0$
$0x33 = 0.0110.011 = +0b1.011 \times 2^{-2} = 0.34375$	$0x73 = 0.1110.011 = +0b1.011 \times 2^6 = 88.0$	$0xb3 = 1.0110.011 = -0b1.011 \times 2^{-2} = -0.34375$	$0xf3 = 1.0110.011 = -0b1.011 \times 2^6 = -88.0$
$0x34 = 0.0110.100 = +0b1.100 \times 2^{-2} = 0.375$	$0x74 = 0.1110.100 = +0b1.100 \times 2^6 = 96.0$	$0xb4 = 1.0110.100 = -0b1.100 \times 2^{-2} = -0.375$	$0xf4 = 1.0110.100 = -0b1.100 \times 2^6 = -96.0$
$0x35 = 0.0110.101 = +0b1.101 \times 2^{-2} = 0.40625$	$0x75 = 0.1110.101 = +0b1.101 \times 2^6 = 104.0$	$0xb5 = 1.0110.101 = -0b1.101 \times 2^{-2} = -0.40625$	$0xf5 = 1.0110.101 = -0b1.101 \times 2^6 = -104.0$
$0x36 = 0.0110.110 = +0b1.110 \times 2^{-2} = 0.4375$	$0x76 = 0.1110.110 = +0b1.110 \times 2^6 = 112.0$	$0xb6 = 1.0110.110 = -0b1.110 \times 2^{-2} = -0.4375$	$0xf6 = 1.0110.110 = -0b1.110 \times 2^6 = -112.0$
$0x37 = 0.0110.111 = +0b1.111 \times 2^{-2} = 0.46875$	$0x77 = 0.1110.111 = +0b1.111 \times 2^6 = 120.0$	$0xb7 = 1.0110.111 = -0b1.111 \times 2^{-2} = -0.46875$	$0xf7 = 1.0110.111 = -0b1.111 \times 2^6 = -120.0$
$0x38 = 0.0111.000 = +0b1.000 \times 2^{-1} = 0.5$	$0x78 = 0.1111.000 = +0b1.000 \times 2^7 = 128.0$	$0xb8 = 1.0111.000 = -0b1.000 \times 2^{-1} = -0.5$	$0xf8 = 1.0111.000 = -0b1.000 \times 2^7 = -128.0$
$0x39 = 0.0111.001 = +0b1.001 \times 2^{-1} = 0.5625$	$0x79 = 0.1111.001 = +0b1.001 \times 2^7 = 144.0$	$0xb9 = 1.0111.001 = -0b1.001 \times 2^{-1} = -0.5625$	$0xf9 = 1.0111.001 = -0b1.001 \times 2^7 = -144.0$
$0x3a = 0.0111.010 = +0b1.010 \times 2^{-1} = 0.625$	$0x7a = 0.1111.010 = +0b1.010 \times 2^7 = 160.0$	$0xba = 1.0111.010 = -0b1.010 \times 2^{-1} = -0.625$	$0xfa = 1.0111.010 = -0b1.010 \times 2^7 = -160.0$
$0x3b = 0.0111.011 = +0b1.011 \times 2^{-1} = 0.6875$	$0x7b = 0.1111.011 = +0b1.011 \times 2^7 = 176.0$	$0xbb = 1.0111.011 = -0b1.011 \times 2^{-1} = -0.6875$	$0xfb = 1.0111.011 = -0b1.011 \times 2^7 = -176.0$
$0x3c = 0.0111.100 = +0b1.100 \times 2^{-1} = 0.75$	$0x7c = 0.1111.100 = +0b1.100 \times 2^7 = 192.0$	$0xbc = 1.0111.100 = -0b1.100 \times 2^{-1} = -0.75$	$0xfc = 1.0111.100 = -0b1.100 \times 2^7 = -192.0$
$0x3d = 0.0111.101 = +0b1.101 \times 2^{-1} = 0.8125$	$0x7d = 0.1111.101 = +0b1.101 \times 2^7 = 208.0$	$0xbd = 1.0111.101 = -0b1.101 \times 2^{-1} = -0.8125$	$0xfd = 1.0111.101 = -0b1.101 \times 2^7 = -208.0$
$0x3e = 0.0111.110 = +0b1.110 \times 2^{-1} = 0.875$	$0x7e = 0.1111.110$		

Standard 8-bit floating-point formats: summary



- OCP standard (<https://bit.ly/3G9EsyG>). E5M2 similar to IEEE 754 formats. E4M3: no infinities, one signed NaN.
- IEEE p3109 (interim report <https://bit.ly/42gPWcy>). No -0 . One unsigned NaN. More rigorous and complete.
- NVIDIA (PTX ISA 8.7 <https://bit.ly/3RNE1ve>) follows the OCP spec.
- AMD does not fully follow the OCP spec. for NaN and inf.
- **Standardisation work ongoing (P3109 is changing significantly).**

Summary


- Standardisation work ongoing until approx 2029.
- We are developing tools and methodologies to probe hardware.
- We maintain a base of hardware knowledge, which we input into standardisation committees.

Key takeaway

OCP is different from IEEE P3109 in terms of defining subsets of reals (formats). For the foreseeable future, two 8-bit standards will be driving the progress in hardware.

Paper that started this work

M. Fasi, N. J. Higham, M. Mikaitis, S. Pranesh. *Numerical Behavior of NVIDIA Tensor Cores*. **Peer J. Comput. Sci.**, 7. Feb. 2021.

 <https://bit.ly/442IIGT>.

References I



M. Klöwer , P. D. Düben , and T. N. Palmer

Number formats, error mitigation, and scope for 16-bit arithmetics in weather and climate modeling analyzed with a shallow water model
J. Adv. Model. Earth Systems, 12. 2020.



M. Fasi, N. J. Higham, S. Pranesh, and M. Mikaitis

Numerical behavior of NVIDIA tensor cores
PeerJ Comput. Sci., 7. 2021



M. Mikaitis

Monotonicity of Multi-Term Floating-Point Adders
IEEE Trans. Comput., 73. 2024.



M. Fasi and M. Mikaitis

CPFloat: A C library for Simulating Low-Precision Arithmetic
ACM Trans. Math. Software, 49. 2023