# Low- and Mixed-Precision Floating-Point Hardware: Algorithms, Error analysis, and Standardisation

Mantas Mikaitis

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

Mathematics Seminar Series, School of Computing and Mathematical Sciences,
University of Leicester, Leicester, UK
12 March, 2025
Slides available: `https://mmikaitis.github.io`

# IEEE 754: Most important standard in computing history

## Standard for binary and decimal fixed-precision arithmetic

Defines subsets of reals, their encoding in memory, conversion and arithmetic behaviour, rounding, exception handling, and more. **Concept of correct rounding**.

Released in 1985, revised in

- 2008
- 2019 (active)
- 2029 (work in progress)

## We are participating in IEEE 754-2029

- Fortnightly meetings, discussion on the mailing list, thoroughly reading the 2019 revision and raising issues.
- Working group: international, many members work in computing industry.

# Floating-point arithmetic, main tools

A floating-point system $\mathbb{F} \subset \mathbb{R}$ is described with $\beta, t, e_{min}, e_{max}$ with elements
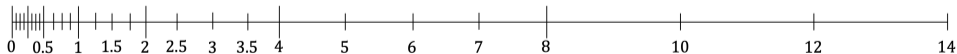
$$x = \pm m \times \beta^{e-t+1}.$$

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

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

## Toy FP system

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

# Floating-point arithmetic, main tools

## Standard model [Higham, 2002]

Given $x, y \in \mathbb{R}$ that lie in the range of $\mathbb{F}$ it can be shown that

$$\mathrm{fl}(x \text{ op } y) = (x \text{ op } y)(1 + \delta), \quad |\delta| \leq u,$$

where $u = 2^{-t}$, $\mathrm{op} \in \{+, -, \times, \div\}$ and **round-to-nearest** mode.

# Building error bounds: small example

Rounding errors $\delta$ accumulate. For example, consider computing $s = x_1 y_1 + x_2 y_2 + x_3 y_3$.

We compute $\hat{s}$ with

$$\hat{s} = \Big( \big( x_1 y_1 (1 + \delta_1) + x_2 y_2 (1 + \delta_2) \big)(1 + \delta_3) + x_3 y_3 (1 + \delta_4) \Big)(1 + \delta_5)$$
$$= x_1 y_1 (1 + \delta_1)(1 + \delta_3)(1 + \delta_5) + x_2 y_2 (1 + \delta_2)(1 + \delta_3)(1 + \delta_5) + x_3 y_3 (1 + \delta_4)(1 + \delta_5).$$

Therefore we compute a solution for the inputs *perturbed at most by* $\prod_{i=1}^{n}(1 + \delta_i)$.

## Worst case backward error bound

$\prod_{i=1}^{n}(1 + \delta_i) = 1 + \theta_n, \quad |\theta_n| \leq \gamma_n$, with $\gamma_n = \frac{nu}{1 - nu}$ and assuming $nu < 1$.

To simplify, we say worst-case error growth is $\mathcal{O}(nu)$.

## 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$K$p$P$ 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.

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

We can write down all numbers in a particular binary8p$P$ set on one page. $\rightarrow$

# IEEE P3109: format on one page (screenshot from the report)

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

# Open Compute Project Floating-Point Standards

Two standards by AMD, Arm, Google, Intel, Meta, and NVIDIA:

- OCP 8-bit Floating Point Specification (released 2023)
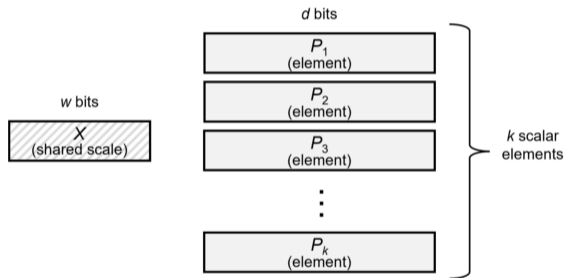- OCP Microscaling Formats (MX) Specification (released 2023)

Key points:

- 8-, 6-, and 4-bit formats.
- Retains $-0$.
- NaNs are removed in some formats, to increase the range.
- Scaled floating-point: common scale factor ($2^X$) for vectors of FP or integer numbers.

## 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.

# Open Compute Project Floating-Point Standards (screenshot from the MX standard)

# Using the TOP500 to anticipate where HPC hardware is going



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.

## Research direction 1: Mixed-precision matrix multipliers

| Architecture | Input format | Accumulation format |
|---|---|---|
| NVIDIA PTX ISA 8.5 | 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 |

# Research direction 1: Using low precision floating-point matrix arith.

Mary and Mikaitis [2024]; Fasi et al. [2023]; Higham et al. [2019].

Approach 1: Use 8-bit FP directly (scale to avoid overflow)

$$C = \Lambda^{-1}\Big(\mathrm{fl}(\Lambda A)\mathrm{fl}(BM)\Big)M^{-1}$$

Approach 2: Use multiple 8-bit FP directly

$$A^{(i)} = \mathrm{fl}\bigg(\Big(\Lambda A - \sum_{k=0}^{i-1} u^k A^{(k)}\Big)/u^i\bigg)$$

and we approximate $C = AB$ as

$$C \approx \Lambda^{-1}\bigg(\sum_{i+j<p} u^{i+j} A^{(i)} B^{(j)}\bigg)M^{-1}.$$

T. Mary and M. Mikaitis [2024].

Matrix multiply; data in $[-10^{10}, -10^{-10}] \cup [10^{-10}, 10^{10}]$.



fp8-E4M3 input
binary32 accumulation
subnormals off

fp8-E4M3 input
binary32 accumulation
subnormals on

A. Abdelfattah, J. Dongarra, M. Fasi, M. Mikaitis, and F. Tisseur, [2025, in prep.].

Ootomo et al. [2024] discovered algorithms for simulating FP matrix multiply with integer matrix multiply. Small example split:

$$\begin{bmatrix} 2^0 \cdot 1.1001000 \\ 2^3 \cdot 1.0000000 \\ -2^1 \cdot 1.1101100 \end{bmatrix} \Rightarrow 2^4 \cdot \begin{bmatrix} \emptyset.\underline{000}\ \underline{110}\ \underline{010}\ \underline{000} \\ \emptyset.\underline{100}\ \underline{000}\ \underline{000}\ \underline{000} \\ -\emptyset.\underline{001}\ \underline{110}\ \underline{110}\ \underline{000} \end{bmatrix} \Rightarrow 2^1 \cdot \begin{bmatrix} 000 \\ 100 \\ -001 \end{bmatrix} + 2^{-2} \cdot \begin{bmatrix} 110 \\ 000 \\ -110 \end{bmatrix} + 2^{-5} \cdot \begin{bmatrix} 010 \\ 000 \\ -110 \end{bmatrix} + 2^{-8} \cdot \begin{bmatrix} 000 \\ 000 \\ 000 \end{bmatrix}$$

$$\quad\ A^T \qquad\qquad\qquad \text{Block fixed-point} \qquad\qquad\quad A^T_{(1)} \qquad\quad A^T_{(2)} \qquad\quad A^T_{(3)} \qquad\quad A^T_{(4)}$$

$$\begin{bmatrix} 2^0 \cdot 1.0110001 \\ -2^2 \cdot 1.1110100 \\ 2^1 \cdot 1.1101000 \end{bmatrix} \Rightarrow 2^3 \cdot \begin{bmatrix} \emptyset.\underline{001}\ \underline{011}\ \underline{000}\ \underline{100} \\ -\emptyset.\underline{111}\ \underline{101}\ \underline{000}\ \underline{000} \\ \emptyset.\underline{011}\ \underline{101}\ \underline{000}\ \underline{000} \end{bmatrix} \Rightarrow 2^0 \cdot \begin{bmatrix} 001 \\ -111 \\ 011 \end{bmatrix} + 2^{-3} \cdot \begin{bmatrix} 011 \\ -101 \\ 101 \end{bmatrix} + 2^{-6} \cdot \begin{bmatrix} 000 \\ 000 \\ 000 \end{bmatrix} + 2^{-9} \cdot \begin{bmatrix} 100 \\ 000 \\ 000 \end{bmatrix}$$

$$\qquad\ B \qquad\qquad\qquad\ \text{Block fixed-point} \qquad\qquad\quad B^{(1)} \qquad\quad B^{(2)} \qquad\quad B^{(3)} \qquad\quad B^{(4)}$$

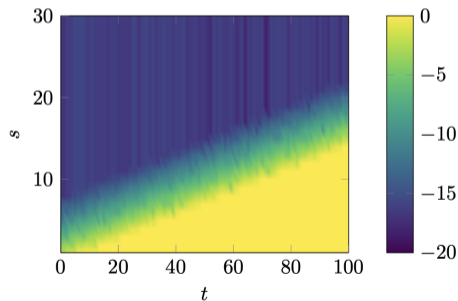# Research direction 1: Using low precision integer matrix arith.

A. Abdelfattah, J. Dongarra, M. Fasi, M. Mikaitis, and F. Tisseur, [2025, in prep.].

We are developing theoretical and experimental analysis of the algorithms. Analysis is general, to cover any future hardware changes.

Here $s$ is the number of splits into 8-bit chunks; colour denotes forward error of (pow 10) $\frac{|a^T b - c|}{|c|}$, where
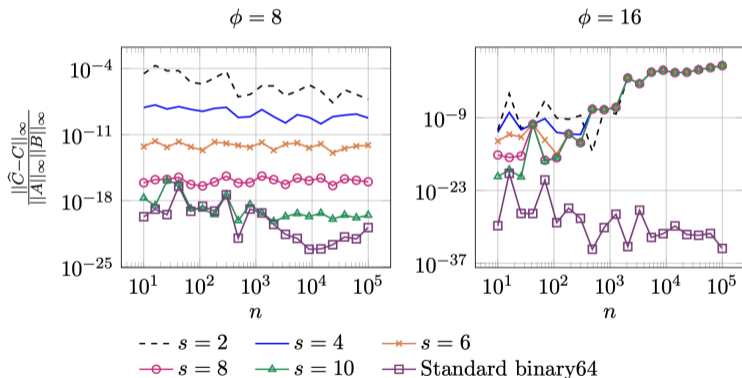
$$a = \begin{bmatrix} 2^{-t}x \\ 1 \end{bmatrix}, \qquad b = \begin{bmatrix} 2^t y \\ 1 \end{bmatrix},$$

and $x$, $y$ are drawn from a uniform distribution and $c$ is a reference solution.

$$A_{ij}, B_{ij} = \text{uniform}(-0.5, 0.5) \times e^{\phi \times \text{normal}(0,1)}.$$

# Research direction 2: Stochastic rounding

With **stochastic rounding** (**SR**), we are not rounding a number to the same direction, but to either direction with probability.

Given some $x$ and FP neighbours $\lfloor x \rfloor$, $\lceil x \rceil$, we round to $\lceil x \rceil$ with prob. $p$ and $\lfloor x \rfloor$ with $p - 1$.



**Mode 1 SR** (nearness): $p = \frac{x - \lfloor x \rfloor}{\lceil x \rceil - \lfloor x \rfloor}$     **Mode 2 SR**: $p = 0.5$

## Mode 1

With **Mode 1 SR** we round $x$ depending on its distances to the nearest two FP numbers, **cancelling out errors of different signs**.

## Standard error model for SR

With SR we replace $u$ by $2u$ since it can round to the second nearest neighbour in $\mathbb{F}$.

## Rounding error analysis

Worst-case error analysis determines the **upper bounds of errors**, while probabilistic error analysis describes **more realistic bounds**.

- Worst-case b-err bound with **RN**: $\frac{nu}{1-nu}$.
- Probabilistic bound with **RN**: $\lambda\sqrt{n}u + \mathcal{O}(u^2)$ w. p. $1 - 2ne^{-\lambda^2/2}$. Requires an assumption that $\delta_i$ are *mean independent zero-mean* quantities—do not always hold [Connolly, Higham, Mary, 2021].

## Wilkinson rule of thumb

$\mathcal{O}(\sqrt{n}u)$ error growth is a rule of thumb with **RN**, but always holds with **SR**.

Backward error in $y = Ax$ where $A \in \mathbb{R}^{100 \times n}$ with entries from uniform dist over $[0, 10^{-3}]$ and $x \in \mathbb{R}^n$ over $[0, 1]$: $\max_i \frac{|\hat{y} - y|_i}{(|A||x|)_i}$.

(a) binary16 arithmetic    (b) bfloat16 arithmetic

# Research direction 2: Consider implementation of SR

Take $m_t$ to be a high precision unrounded significand from an operation.

Take $t$ to be source precision and $k$ the precision of random numbers.

E.-M. El Arar et al. [2024]: we derived a new bound $\mathcal{O}(\sqrt{n}u_p + nu_{p+r})$

Error when adding $n$ random values in $(0, 1)$ in 16-bit floating point:



Relative forward error | Bounds with $1 - \lambda = 0.9$

Legend:
- RN
- $\text{SR}_{11,7}$
- $\text{SR}_{11,3}$
- $\text{SR}_{11,8}$
- $\text{SR}_{11,6}$
- $\text{SR}_{11,10}$

Given a small matrix multiplier not necessarily using IEEE 754 $+, \times$:

$$D \qquad = \qquad C \qquad + \qquad A \qquad \times \qquad B,$$

$$\begin{bmatrix} \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \end{bmatrix} = \begin{bmatrix} \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \end{bmatrix} + \begin{bmatrix} \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \end{bmatrix} \times \begin{bmatrix} \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \end{bmatrix}$$

$$\underbrace{\phantom{xxxxxxxx}}_{\substack{\text{binary16 or} \\ \text{binary32}}} \qquad \underbrace{\phantom{xxxxxxxx}}_{\substack{\text{binary16 or} \\ \text{binary32}}} \qquad \underbrace{\phantom{xxxxxxxx}}_{\text{binary16}} \qquad \underbrace{\phantom{xxxxxxxx}}_{\text{binary16}}$$

We are building theory and tools that allow to *report the internal precision, order of computations, rounding mode, and more*. Fasi, Higham, Pranesh, Mikaitis [2021].

We reported bit level differences between NVIDIA V100 and A100, not documented.

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

- 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.

## Research direction 3: 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 throughput means we are using MMAs in other areas that traditionally use standardized FP.
- Can be said we are back to pre-IEEE-754 with mixed-precision.
- Eventually need to standardise (IEEE working group **P3109**).

**For now we create tests that determine features and differences between architectures.**

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**.

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

1-bit difference in accumulators of V100 and T4/A100.

Normalization at the end, not at intermediate adds.

Rounding is not to nearest.

Monotonicity: increase one input, do not change order, dot product decreases. See [Mikaitis, 2024].

No fixed order.

Consider determining the rounding direction:

## Research direction 3: Looking for numerical features by testing hardware

3-year EPSRC-funded project 2025–2028.

Collaboration with Argonne National Lab and Intel (US).



WP4: Develop automated open-source software to report hardware features.

# Research tool development: custom precision simulators

- Various packages available: `chop`, `FLOATP`, `QPyTorch`.
- Usual approach is to perform ops in binary32/64 HW.
- Round down to sub-32-bit precision: careful with double rounding.
- We believe ours is most customizable and fastest: `CPFloat` [Fasi & Mikaitis, 2023].
- Can be used in MATLAB, Octave or C.

```
>> options.format = 'binary16';
>> [~,options] = cpfloat(0, options)
options = struct with fields:
       format: 'binary16'
       params: [11 15]
    subnormal: 1
        round: 5
         flip: 0
            p: 0.5000
       explim: 1
>> cpfloat(pi, options)
ans =
    3.1406
>> options.params(1) = options.params(1) + 1;
>> cpfloat(pi, options)
ans =
    3.1426
```

# Summary

- We are building theory and tools to explain computer arithmetic hardware and analyse algorithms.
- People have been there before: pre-IEEE-754 1985.
- More complicated now: vector and matrix hardware.
- Standardisation will impact future low-precision hardware.

Slides and more info at http://mmikaitis.github.io

# Leeds Mathematical Software and Hardware Lab

Informal group, within Scientific Computation group, in the School of Computing, Univ. Leeds.



Massimiliano Fasi
Lecturer

Research and teaching



Mantas Mikaitis
Lecturer

Research and teaching



- Focusing on computer arithmetic, numerical linear algebra, high-performance computing.
- Working with IEEE P3109 and IEEE 754-2029.
- Serving on PC committees of ARITH.
- Planning MSc module on computer arithmetic.
- PhD studentships available.

# References I

📄 N. J. Higham
Accuracy and Stability of Numerical Algorithms. 2nd edition
SIAM. 2002

📄 T. Mary and M. Mikaitis
Error Analysis of Matrix Multiplication with Narrow Range Floating-Point Arithmetic
hal-04671474. Aug. 2024

📄 M. Fasi, N. J. Higham, F. Lopez, T. Mary, and M. Mikaitis
Matrix Multiplication in Multiword Arithmetic: Error Analysis and Application to GPU
Tensor Cores
SIAM J. Sci. Comput., 45:1. Feb. 2023.

📄 N. J. Higham, S. Pranesh, and M. Zounon
Squeezing a Matrix into Half Precision, with an Application to Solving Linear Systems
SIAM J. Sci. Comput., 41:4. Aug. 2019.

# References II

📄 A. Abdelfattah, J. Dongarra, M. Fasi, M. Mikaitis, and F. Tisseur
Error Analysis of Floating-Point Matrix Multiplication Computed via Low-Precision Integer
Arithmetic
In preparation. 2025.

📄 H. Ootomo, K. Ozaki, and R. Yokota
DGEMM on integer matrix multiplication unit
Int. J. High Perform. Comput. Appl., 38:4. Mar. 2024.

📄 M. P. Connolly, N. J. Higham, T. Mary
Stochastic rounding and its probabilistic backward error analysis.
SIAM J. Sci. Comput., 43. 2021.

📄 E.-M. El Arar, M. Fasi, S.-I. Filip, and M. Mikaitis
Probabilistic Error Analysis of Limited-Precision Stochastic Rounding
hal-04665809. Jul. 2024.

📄 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