Evaluation of Bfloat16, Posit, and Takum Arithmetics in Sparse Linear Solvers

Laslo Hunhold¹ James Quinlan²

¹Parallel and Distributed Systems Group, University of Cologne, Germany

²Department of Computer Science, University of Southern Maine, Portland, ME, USA

5th May 2025



Laslo Hunhold, James Quinlan

Memory wall dictates move to smaller machine number types (end of 'doubles everywhere')

- Memory wall dictates move to smaller machine number types (end of 'doubles everywhere')
- Overparametrization of DNNs makes low-precision arithmetic sufficient

- Memory wall dictates move to smaller machine number types (end of 'doubles everywhere')
- Overparametrization of DNNs makes low-precision arithmetic sufficient
- Common criticisms of IEEE 754 for low-precision arithmetic

- Memory wall dictates move to smaller machine number types (end of 'doubles everywhere')
- Overparametrization of DNNs makes low-precision arithmetic sufficient
- Common criticisms of IEEE 754 for low-precision arithmetic
 - No standardized 8-bit type

- Memory wall dictates move to smaller machine number types (end of 'doubles everywhere')
- Overparametrization of DNNs makes low-precision arithmetic sufficient
- Common criticisms of IEEE 754 for low-precision arithmetic
 - No standardized 8-bit type \rightarrow OCP OFP8 (*E4M3*, E5M2)

- Memory wall dictates move to smaller machine number types (end of 'doubles everywhere')
- Overparametrization of DNNs makes low-precision arithmetic sufficient
- Common criticisms of IEEE 754 for low-precision arithmetic
 - No standardized 8-bit type \rightarrow OCP OFP8 (*E4M3*, E5M2)
 - Redundancies dominate at low precisions, overflows due to small dynamic range

- Memory wall dictates move to smaller machine number types (end of 'doubles everywhere')
- Overparametrization of DNNs makes low-precision arithmetic sufficient
- Common criticisms of IEEE 754 for low-precision arithmetic
 - No standardized 8-bit type \rightarrow OCP OFP8 (*E4M3*, E5M2)

- Memory wall dictates move to smaller machine number types (end of 'doubles everywhere')
- Overparametrization of DNNs makes low-precision arithmetic sufficient
- Common criticisms of IEEE 754 for low-precision arithmetic
 - No standardized 8-bit type \rightarrow OCP OFP8 (*E4M3*, E5M2)

 - float16 has insufficient dynamic range

- Memory wall dictates move to smaller machine number types (end of 'doubles everywhere')
- Overparametrization of DNNs makes low-precision arithmetic sufficient
- Common criticisms of IEEE 754 for low-precision arithmetic
 - No standardized 8-bit type \rightarrow OCP OFP8 (*E4M3*, E5M2)

 - float16 has insufficient dynamic range \rightarrow bfloat16

- Memory wall dictates move to smaller machine number types (end of 'doubles everywhere')
- Overparametrization of DNNs makes low-precision arithmetic sufficient
- Common criticisms of IEEE 754 for low-precision arithmetic
 - No standardized 8-bit type \rightarrow OCP OFP8 (*E4M3*, E5M2)

 - float16 has insufficient dynamic range \rightarrow bfloat16
 - Number distribution does not reflect real-world usage

- Memory wall dictates move to smaller machine number types (end of 'doubles everywhere')
- Overparametrization of DNNs makes low-precision arithmetic sufficient
- Common criticisms of IEEE 754 for low-precision arithmetic
 - No standardized 8-bit type \rightarrow OCP OFP8 (*E4M3*, E5M2)

 - float16 has insufficient dynamic range \rightarrow bfloat16
 - \blacktriangleright Number distribution does not reflect real-world usage \rightarrow tapered precision arithmetic

- Memory wall dictates move to smaller machine number types (end of 'doubles everywhere')
- Overparametrization of DNNs makes low-precision arithmetic sufficient
- Common criticisms of IEEE 754 for low-precision arithmetic
 - No standardized 8-bit type \rightarrow OCP OFP8 (*E4M3*, E5M2)

 - float16 has insufficient dynamic range \rightarrow bfloat16
 - \blacktriangleright Number distribution does not reflect real-world usage \rightarrow tapered precision arithmetic



- Memory wall dictates move to smaller machine number types (end of 'doubles everywhere')
- Overparametrization of DNNs makes low-precision arithmetic sufficient
- Common criticisms of IEEE 754 for low-precision arithmetic
 - No standardized 8-bit type \rightarrow OCP OFP8 (*E4M3*, E5M2)

 - float16 has insufficient dynamic range \rightarrow bfloat16
 - \blacktriangleright Number distribution does not reflect real-world usage \rightarrow tapered precision arithmetic



- Memory wall dictates move to smaller machine number types (end of 'doubles everywhere')
- Overparametrization of DNNs makes low-precision arithmetic sufficient
- Common criticisms of IEEE 754 for low-precision arithmetic
 - No standardized 8-bit type \rightarrow OCP OFP8 (*E4M3*, E5M2)
 - Redundancies dominate at low precisions, overflows due to small dynamic range
 → IEEE P3109 (one NaN/zero, saturation)
 - float16 has insufficient dynamic range \rightarrow bfloat16
 - \blacktriangleright Number distribution does not reflect real-world usage \rightarrow tapered precision arithmetic



Strong focus on deep learning \rightarrow Performance in general-purpose computing?

- Memory wall dictates move to smaller machine number types (end of 'doubles everywhere')
- Overparametrization of DNNs makes low-precision arithmetic sufficient
- Common criticisms of IEEE 754 for low-precision arithmetic
 - No standardized 8-bit type \rightarrow OCP OFP8 (*E4M3*, E5M2)
 - ▶ Redundancies dominate at low precisions, overflows due to small dynamic range → IEEE P3109 (one NaN/zero, saturation)
 - float16 has insufficient dynamic range \rightarrow bfloat16
 - \blacktriangleright Number distribution does not reflect real-world usage \rightarrow tapered precision arithmetic



- Strong focus on deep learning \rightarrow Performance in general-purpose computing?
- This work: Evaluation within sparse linear solvers

Consider general floating-point format



Consider general floating-point format



Sign and fraction bits already contain maximum information

Consider general floating-point format



Sign and fraction bits already contain maximum information

Angle of approach: exponent bits

Consider general floating-point format



Sign and fraction bits already contain maximum information

Angle of approach: exponent bits

Want higher density of numbers (i.e. more fraction bits, precision) for exponent values close to zero

Consider general floating-point format



Sign and fraction bits already contain maximum information

Angle of approach: exponent bits

- Want higher density of numbers (i.e. more fraction bits, precision) for exponent values close to zero
- Solution: variable-length exponent encoding

Consider general floating-point format



Sign and fraction bits already contain maximum information

Angle of approach: exponent bits

- Want higher density of numbers (i.e. more fraction bits, precision) for exponent values close to zero
- Solution: variable-length exponent encoding

Small magnitude: shorter exponent, longer fraction, higher density



Consider general floating-point format



Sign and fraction bits already contain maximum information

Angle of approach: exponent bits

- Want higher density of numbers (i.e. more fraction bits, precision) for exponent values close to zero
- Solution: variable-length exponent encoding

Small magnitude: shorter exponent, longer fraction, higher density



Large magnitude: longer exponent, shorter fraction, lower density

| $\stackrel{\text{sign}}{\leftrightarrow}$ | exponent | + fraction |
|---|----------|------------|
| | | |

Laslo Hunhold, James Quinlan

Evaluation of Bfloat16, Posit, and Takum Arithmetics in Sparse Linear Solvers

¹John L. Gustafson et al. 'Standard for Posit Arithmetic (2022)'. Mar. 2022

²Florent de Dinechin et al. 'Posits: the good, the bad and the ugly'. Mar. 2019

State of the art tapered machine number format, in standardisation¹

¹ John L. Gustafson et al. 'Standard for Posit Arithmetic (2022)'. Mar. 2022

 $^{^2}$ Florent de Dinechin et al. 'Posits: the good, the bad and the ugly'. Mar. 2019

- State of the art tapered machine number format, in standardisation¹
- Active research field (hundreds of publications since 2017)

¹John L. Gustafson et al. 'Standard for Posit Arithmetic (2022)'. Mar. 2022

²Florent de Dinechin et al. 'Posits: the good, the bad and the ugly'. Mar. 2019

- State of the art tapered machine number format, in standardisation¹
- Active research field (hundreds of publications since 2017)
- Numerous implementations, latest by Calligo with RISC-V SoC 'TUNGA' (2024)

¹John L. Gustafson et al. 'Standard for Posit Arithmetic (2022)'. Mar. 2022

²Florent de Dinechin et al. 'Posits: the good, the bad and the ugly'. Mar. 2019

- State of the art tapered machine number format, in standardisation¹
- Active research field (hundreds of publications since 2017)
- Numerous implementations, latest by Calligo with RISC-V SoC 'TUNGA' (2024)
- Exponent coding

 $^{^{1}}$ John L. Gustafson et al. 'Standard for Posit Arithmetic (2022)'. Mar. 2022

²Florent de Dinechin et al. 'Posits: the good, the bad and the ugly'. Mar. 2019

- State of the art tapered machine number format, in standardisation¹
- Active research field (hundreds of publications since 2017)
- Numerous implementations, latest by Calligo with RISC-V SoC 'TUNGA' (2024)
- Exponent coding



Evaluation of Bfloat16, Posit, and Takum Arithmetics in Sparse Linear Solvers

¹John L. Gustafson et al. 'Standard for Posit Arithmetic (2022)'. Mar. 2022

²Florent de Dinechin et al. 'Posits: the good, the bad and the ugly'. Mar. 2019

- State of the art tapered machine number format, in standardisation¹
- Active research field (hundreds of publications since 2017)
- Numerous implementations, latest by Calligo with RISC-V SoC 'TUNGA' (2024)
- Exponent coding



▶ *R* is run of *k* zeros or ones, followed by one or zero $(\overline{R_0})$ (prefix code)

¹John L. Gustafson et al. 'Standard for Posit Arithmetic (2022)'. Mar. 2022

²Florent de Dinechin et al. 'Posits: the good, the bad and the ugly'. Mar. 2019

- State of the art tapered machine number format, in standardisation¹
- Active research field (hundreds of publications since 2017)
- Numerous implementations, latest by Calligo with RISC-V SoC 'TUNGA' (2024)
- Exponent coding



▶ *R* is run of *k* zeros or ones, followed by one or zero $(\overline{R_0})$ (prefix code)

► Coded exponent is
$$e = \begin{cases} -4k + \text{uint}(E) & \overline{R_0} = 0\\ 4(k-1) + \text{uint}(E) & \overline{R_0} = 1 \end{cases}$$

Evaluation of Bfloat16, Posit, and Takum Arithmetics in Sparse Linear Solvers

 $¹_{\mbox{John L. Gustafson et al.}}$ 'Standard for Posit Arithmetic (2022)'. Mar. 2022

²Florent de Dinechin et al. 'Posits: the good, the bad and the ugly'. Mar. 2019

- State of the art tapered machine number format, in standardisation¹
- Active research field (hundreds of publications since 2017)
- Numerous implementations, latest by Calligo with RISC-V SoC 'TUNGA' (2024)
- Exponent coding



▶ *R* is run of *k* zeros or ones, followed by one or zero $(\overline{R_0})$ (prefix code)

► Coded exponent is $e = \begin{cases} -4k + \text{uint}(E) & \overline{R_0} = 0\\ 4(k-1) + \text{uint}(E) & \overline{R_0} = 1 \end{cases}$

• Efficient for small exponents (e.g. $1010 \equiv 2, 0111 \equiv -1$)

¹ John L. Gustafson et al. 'Standard for Posit Arithmetic (2022)'. Mar. 2022

²Florent de Dinechin et al. 'Posits: the good, the bad and the ugly'. Mar. 2019

- State of the art tapered machine number format, in standardisation¹
- Active research field (hundreds of publications since 2017)
- Numerous implementations, latest by Calligo with RISC-V SoC 'TUNGA' (2024)
- Exponent coding



R is run of k zeros or ones, followed by one or zero $(\overline{R_0})$ (prefix code)

- ► Coded exponent is $e = \begin{cases} -4k + \text{uint}(E) & \overline{R_0} = 0\\ 4(k-1) + \text{uint}(E) & \overline{R_0} = 1 \end{cases}$
- Efficient for small exponents (e.g. $1010 \equiv 2, 0111 \equiv -1$)
- Inefficient for large exponents

Evaluation of Bfloat16, Posit, and Takum Arithmetics in Sparse Linear Solvers

¹John L. Gustafson et al. 'Standard for Posit Arithmetic (2022)'. Mar. 2022

²Florent de Dinechin et al. 'Posits: the good, the bad and the ugly'. Mar. 2019

- State of the art tapered machine number format, in standardisation¹
- Active research field (hundreds of publications since 2017)
- Numerous implementations, latest by Calligo with RISC-V SoC 'TUNGA' (2024)
- Exponent coding



- R is run of k zeros or ones, followed by one or zero $(\overline{R_0})$ (prefix code)
- Coded exponent is $e = \begin{cases} -4k + \operatorname{uint}(E) & \overline{R_0} = 0\\ 4(k-1) + \operatorname{uint}(E) & \overline{R_0} = 1 \end{cases}$
- Efficient for small exponents (e.g. $1010 \equiv 2, 0111 \equiv -1$)

Evaluation of Bfloat16, Posit, and Takum Arithmetics in Sparse Linear Solvers

 $^{^{1}}$ John L. Gustafson et al. 'Standard for Posit Arithmetic (2022)'. Mar. 2022

²Florent de Dinechin et al. 'Posits: the good, the bad and the ugly'. Mar. 2019
- State of the art tapered machine number format, in standardisation¹
- Active research field (hundreds of publications since 2017)
- Numerous implementations, latest by Calligo with RISC-V SoC 'TUNGA' (2024)
- Exponent coding



R is run of k zeros or ones, followed by one or zero $(\overline{R_0})$ (prefix code)

- ► Coded exponent is $e = \begin{cases} -4k + \text{uint}(E) & \overline{R_0} = 0\\ 4(k-1) + \text{uint}(E) & \overline{R_0} = 1 \end{cases}$
- Efficient for small exponents (e.g. $1010 \equiv 2, 0111 \equiv -1$)

Properties

¹ John L. Gustafson et al. 'Standard for Posit Arithmetic (2022)'. Mar. 2022

²Florent de Dinechin et al. 'Posits: the good, the bad and the ugly'. Mar. 2019

- State of the art tapered machine number format, in standardisation¹
- Active research field (hundreds of publications since 2017)
- Numerous implementations, latest by Calligo with RISC-V SoC 'TUNGA' (2024)
- Exponent coding



- R is run of k zeros or ones, followed by one or zero $(\overline{R_0})$ (prefix code)
- Coded exponent is $e = \begin{cases} -4k + \operatorname{uint}(E) & \overline{R_0} = 0\\ 4(k-1) + \operatorname{uint}(E) & \overline{R_0} = 1 \end{cases}$
- Efficient for small exponents (e.g. $1010 \equiv 2, 0111 \equiv -1$)

Properties

No redundant representations

¹John L. Gustafson et al. 'Standard for Posit Arithmetic (2022)'. Mar. 2022

²Florent de Dinechin et al. 'Posits: the good, the bad and the ugly'. Mar. 2019

- State of the art tapered machine number format, in standardisation¹
- Active research field (hundreds of publications since 2017)
- Numerous implementations, latest by Calligo with RISC-V SoC 'TUNGA' (2024)
- Exponent coding



- R is run of k zeros or ones, followed by one or zero $(\overline{R_0})$ (prefix code)
- ► Coded exponent is $e = \begin{cases} -4k + \operatorname{uint}(E) & \overline{R_0} = 0\\ 4(k-1) + \operatorname{uint}(E) & \overline{R_0} = 1 \end{cases}$
- Efficient for small exponents (e.g. $1010 \equiv 2, 0111 \equiv -1$)

Properties

- No redundant representations
- Symmetry with two's complement integers

¹John L. Gustafson et al. 'Standard for Posit Arithmetic (2022)'. Mar. 2022

²Florent de Dinechin et al. 'Posits: the good, the bad and the ugly'. Mar. 2019

- State of the art tapered machine number format, in standardisation¹
- Active research field (hundreds of publications since 2017)
- Numerous implementations, latest by Calligo with RISC-V SoC 'TUNGA' (2024)
- Exponent coding



- R is run of k zeros or ones, followed by one or zero $(\overline{R_0})$ (prefix code)
- ► Coded exponent is $e = \begin{cases} -4k + \operatorname{uint}(E) & \overline{R_0} = 0\\ 4(k-1) + \operatorname{uint}(E) & \overline{R_0} = 1 \end{cases}$
- Efficient for small exponents (e.g. $1010 \equiv 2, 0111 \equiv -1$)

Properties

- No redundant representations
- Symmetry with two's complement integers (unsigned zero of all zero bits

¹John L. Gustafson et al. 'Standard for Posit Arithmetic (2022)'. Mar. 2022

²Florent de Dinechin et al. 'Posits: the good, the bad and the ugly'. Mar. 2019

- State of the art tapered machine number format, in standardisation¹
- Active research field (hundreds of publications since 2017)
- Numerous implementations, latest by Calligo with RISC-V SoC 'TUNGA' (2024)
- Exponent coding



- R is run of k zeros or ones, followed by one or zero $(\overline{R_0})$ (prefix code)
- ► Coded exponent is $e = \begin{cases} -4k + \operatorname{uint}(E) & \overline{R_0} = 0\\ 4(k-1) + \operatorname{uint}(E) & \overline{R_0} = 1 \end{cases}$
- Efficient for small exponents (e.g. $1010 \equiv 2, 0111 \equiv -1$)

Properties

- No redundant representations
- Symmetry with two's complement integers (unsigned zero of all zero bits, negation

¹John L. Gustafson et al. 'Standard for Posit Arithmetic (2022)'. Mar. 2022

²Florent de Dinechin et al. 'Posits: the good, the bad and the ugly'. Mar. 2019

- State of the art tapered machine number format, in standardisation¹
- Active research field (hundreds of publications since 2017)
- Numerous implementations, latest by Calligo with RISC-V SoC 'TUNGA' (2024)
- Exponent coding



- R is run of k zeros or ones, followed by one or zero $(\overline{R_0})$ (prefix code)
- ► Coded exponent is $e = \begin{cases} -4k + \operatorname{uint}(E) & \overline{R_0} = 0\\ 4(k-1) + \operatorname{uint}(E) & \overline{R_0} = 1 \end{cases}$
- Efficient for small exponents (e.g. $1010 \equiv 2, 0111 \equiv -1$)

Properties

- No redundant representations
- Symmetry with two's complement integers (unsigned zero of all zero bits, negation, ordering

¹John L. Gustafson et al. 'Standard for Posit Arithmetic (2022)'. Mar. 2022

²Florent de Dinechin et al. 'Posits: the good, the bad and the ugly'. Mar. 2019

- State of the art tapered machine number format, in standardisation¹
- Active research field (hundreds of publications since 2017)
- Numerous implementations, latest by Calligo with RISC-V SoC 'TUNGA' (2024)
- Exponent coding



- R is run of k zeros or ones, followed by one or zero $(\overline{R_0})$ (prefix code)
- ► Coded exponent is $e = \begin{cases} -4k + \operatorname{uint}(E) & \overline{R_0} = 0\\ 4(k-1) + \operatorname{uint}(E) & \overline{R_0} = 1 \end{cases}$
- Efficient for small exponents (e.g. $1010 \equiv 2, 0111 \equiv -1$)

Properties

- No redundant representations
- Symmetry with two's complement integers (unsigned zero of all zero bits, negation, ordering, sign)

¹John L. Gustafson et al. 'Standard for Posit Arithmetic (2022)'. Mar. 2022

²Florent de Dinechin et al. 'Posits: the good, the bad and the ugly'. Mar. 2019

- State of the art tapered machine number format, in standardisation¹
- Active research field (hundreds of publications since 2017)
- Numerous implementations, latest by Calligo with RISC-V SoC 'TUNGA' (2024)
- Exponent coding



- R is run of k zeros or ones, followed by one or zero $(\overline{R_0})$ (prefix code)
- Coded exponent is $e = \begin{cases} -4k + \operatorname{uint}(E) & \overline{R_0} = 0\\ 4(k-1) + \operatorname{uint}(E) & \overline{R_0} = 1 \end{cases}$
- Efficient for small exponents (e.g. $1010 \equiv 2, 0111 \equiv -1$)

Properties

- No redundant representations
- Symmetry with two's complement integers (unsigned zero of all zero bits, negation, ordering, sign)
- Defined for all n, conversion between lengths simple rounding or expansion

¹John L. Gustafson et al. 'Standard for Posit Arithmetic (2022)'. Mar. 2022

²Florent de Dinechin et al. 'Posits: the good, the bad and the ugly'. Mar. 2019

Laslo Hunhold, James Quinlan









Goals

More efficient exponent code



- More efficient exponent code
- Preserve useful posit properties



- More efficient exponent code
- Preserve useful posit properties
- Design Approach



Goals

- More efficient exponent code
- Preserve useful posit properties

Design Approach

Candidate sequence of 'saturated' integers $2^k - 1$ with $k \in \mathbb{N}_1$ (1, 3, 7, 15, 31, 63, 127, 255, 511, 1023, . . .) for largest exponent value



► Goals

- More efficient exponent code
- Preserve useful posit properties

Design Approach

- Candidate sequence of 'saturated' integers $2^k 1$ with $k \in \mathbb{N}_1$ (1, 3, 7, 15, 31, 63, 127, 255, 511, 1023, . . .) for largest exponent value
- Tapered format: Maximum exponent length must also be saturated



- More efficient exponent code
- Preserve useful posit properties
- Design Approach
 - Candidate sequence of 'saturated' integers $2^k 1$ with $k \in \mathbb{N}_1$ (1, 3, 7, 15, 31, 63, 127, 255, 511, 1023, . . .) for largest exponent value
 - ▶ Tapered format: Maximum exponent length must also be saturated
 - Take the subsequence 2^{2^k-1} 1 of integers whose bit length is an incremented saturated integer (3, 15, 255, 65535, 4294967295, ...)



- More efficient exponent code
- Preserve useful posit properties
- Design Approach
 - Candidate sequence of 'saturated' integers $2^k 1$ with $k \in \mathbb{N}_1$ (1, 3, 7, 15, 31, 63, 127, 255, 511, 1023, . . .) for largest exponent value
 - ▶ Tapered format: Maximum exponent length must also be saturated
 - Take the subsequence 2^{2^k-1} 1 of integers whose bit length is an incremented saturated integer (3, 15, 255, 65535, 4294967295, ...)
 - Maximum exponent value 255 (with bit length 8) follows naturally



- More efficient exponent code
- Preserve useful posit properties
- Design Approach
 - Candidate sequence of 'saturated' integers $2^k 1$ with $k \in \mathbb{N}_1$ (1, 3, 7, 15, 31, 63, 127, 255, 511, 1023, . . .) for largest exponent value
 - ▶ Tapered format: Maximum exponent length must also be saturated
 - ► Take the subsequence 2^{2^k-1} 1 of integers whose bit length is an incremented saturated integer (3, 15, 255, 65535, 4294967295, ...)
 - Maximum exponent value 255 (with bit length 8) follows naturally
- > 3-bit regime (0-7) encodes exponent bit count, followed by 0 to 7 exponent bits.



- More efficient exponent code
- Preserve useful posit properties
- Design Approach
 - Candidate sequence of 'saturated' integers $2^k 1$ with $k \in \mathbb{N}_1$ (1, 3, 7, 15, 31, 63, 127, 255, 511, 1023, . . .) for largest exponent value
 - ▶ Tapered format: Maximum exponent length must also be saturated
 - Take the subsequence 2^{2^k-1} 1 of integers whose bit length is an incremented saturated integer (3, 15, 255, 65535, 4294967295, ...)
 - Maximum exponent value 255 (with bit length 8) follows naturally
- ▶ 3-bit regime (0-7) encodes exponent bit count, followed by 0 to 7 exponent bits.
- Exponent value has implicit leading 1-bit (1-8 bits), subtract 1 for range 0-254



- More efficient exponent code
- Preserve useful posit properties
- Design Approach
 - Candidate sequence of 'saturated' integers $2^k 1$ with $k \in \mathbb{N}_1$ (1, 3, 7, 15, 31, 63, 127, 255, 511, 1023, . . .) for largest exponent value
 - Tapered format: Maximum exponent length must also be saturated
 - Take the subsequence 2^{2^k-1} 1 of integers whose bit length is an incremented saturated integer (3, 15, 255, 65535, 4294967295, ...)
 - Maximum exponent value 255 (with bit length 8) follows naturally
- ▶ 3-bit regime (0-7) encodes exponent bit count, followed by 0 to 7 exponent bits.
- Exponent value has implicit leading 1-bit (1-8 bits), subtract 1 for range 0-254

| value | 0 | 1 | 2 | 3 | 4 | 254 |
|------------|-----|------|------|-------|-------|----------------|
| value bits | 1 | 10 | 11 | 100 | 101 | 11111111 |
| encoding | 000 | 0010 | 0011 | 01000 | 01001 | 1111111111 |



Goals

- More efficient exponent code
- Preserve useful posit properties
- Design Approach
 - Candidate sequence of 'saturated' integers $2^k 1$ with $k \in \mathbb{N}_1$ (1, 3, 7, 15, 31, 63, 127, 255, 511, 1023, . . .) for largest exponent value
 - ▶ Tapered format: Maximum exponent length must also be saturated
 - Take the subsequence 2^{2^k-1} 1 of integers whose bit length is an incremented saturated integer (3, 15, 255, 65535, 4294967295, ...)
 - Maximum exponent value 255 (with bit length 8) follows naturally
- ▶ 3-bit regime (0-7) encodes exponent bit count, followed by 0 to 7 exponent bits.
- Exponent value has implicit leading 1-bit (1-8 bits), subtract 1 for range 0-254

| value | 0 | 1 | 2 | 3 | 4 | 254 |
|------------|-----|------|------|-------|-------|----------------|
| value bits | 1 | 10 | 11 | 100 | 101 | 11111111 |
| encoding | 000 | 0010 | 0011 | 01000 | 01001 | 1111111111 |

Icelandic 'takmarkað umfang', meaning 'limited range'



- More efficient exponent code
- Preserve useful posit properties
- Design Approach
 - Candidate sequence of 'saturated' integers $2^k 1$ with $k \in \mathbb{N}_1$ (1, 3, 7, 15, 31, 63, 127, 255, 511, 1023, . . .) for largest exponent value
 - Tapered format: Maximum exponent length must also be saturated
 - Take the subsequence 2^{2^k-1} 1 of integers whose bit length is an incremented saturated integer (3, 15, 255, 65535, 4294967295, ...)
 - Maximum exponent value 255 (with bit length 8) follows naturally
- ▶ 3-bit regime (0-7) encodes exponent bit count, followed by 0 to 7 exponent bits.
- Exponent value has implicit leading 1-bit (1-8 bits), subtract 1 for range 0-254

| value | 0 | 1 | 2 | 3 | 4 | 254 |
|------------|-----|------|------|-------|-------|----------------|
| value bits | 1 | 10 | 11 | 100 | 101 | 11111111 |
| encoding | 000 | 0010 | 0011 | 01000 | 01001 | 1111111111 |

- Icelandic 'takmarkað umfang', meaning 'limited range'
- Separate 'direction' bit D for exponent sign



- More efficient exponent code
- Preserve useful posit properties
- Design Approach
 - Candidate sequence of 'saturated' integers $2^k 1$ with $k \in \mathbb{N}_1$ (1, 3, 7, 15, 31, 63, 127, 255, 511, 1023, . . .) for largest exponent value
 - Tapered format: Maximum exponent length must also be saturated
 - Take the subsequence 2^{2^k-1} 1 of integers whose bit length is an incremented saturated integer (3, 15, 255, 65535, 4294967295, ...)
 - Maximum exponent value 255 (with bit length 8) follows naturally
- ▶ 3-bit regime (0-7) encodes exponent bit count, followed by 0 to 7 exponent bits.
- Exponent value has implicit leading 1-bit (1-8 bits), subtract 1 for range 0-254

| value | 0 | 1 | 2 | 3 | 4 | 254 |
|------------|-----|------|------|-------|-------|----------------|
| value bits | 1 | 10 | 11 | 100 | 101 | 11111111 |
| encoding | 000 | 0010 | 0011 | 01000 | 01001 | 1111111111 |

- Icelandic 'takmarkað umfang', meaning 'limited range'
- Separate 'direction' bit D for exponent sign
- ▶ Interlude: Comparison for value 254 (posit \rightarrow takum):



Goals

- More efficient exponent code
- Preserve useful posit properties
- Design Approach
 - Candidate sequence of 'saturated' integers $2^k 1$ with $k \in \mathbb{N}_1$ (1, 3, 7, 15, 31, 63, 127, 255, 511, 1023, . . .) for largest exponent value
 - ▶ Tapered format: Maximum exponent length must also be saturated
 - ► Take the subsequence 2^{2^k-1} − 1 of integers whose bit length is an incremented saturated integer (3, 15, 255, 65535, 4294967295, ...)
 - Maximum exponent value 255 (with bit length 8) follows naturally
- ▶ 3-bit regime (0-7) encodes exponent bit count, followed by 0 to 7 exponent bits.
- Exponent value has implicit leading 1-bit (1-8 bits), subtract 1 for range 0-254

| value | 0 | 1 | 2 | 3 | 4 | 254 |
|------------|-----|------|------|-------|-------|----------------|
| value bits | 1 | 10 | 11 | 100 | 101 | 11111111 |
| encoding | 000 | 0010 | 0011 | 01000 | 01001 | 1111111111 |

- Icelandic 'takmarkað umfang', meaning 'limited range'
- Separate 'direction' bit D for exponent sign
- ▶ Interlude: Comparison for value 254 (posit \rightarrow takum):



Goals

- More efficient exponent code
- Preserve useful posit properties
- Design Approach
 - Candidate sequence of 'saturated' integers $2^k 1$ with $k \in \mathbb{N}_1$ (1, 3, 7, 15, 31, 63, 127, 255, 511, 1023, . . .) for largest exponent value
 - ▶ Tapered format: Maximum exponent length must also be saturated
 - ► Take the subsequence 2^{2^k-1} − 1 of integers whose bit length is an incremented saturated integer (3, 15, 255, 65535, 4294967295, ...)
 - Maximum exponent value 255 (with bit length 8) follows naturally
- ▶ 3-bit regime (0-7) encodes exponent bit count, followed by 0 to 7 exponent bits.
- Exponent value has implicit leading 1-bit (1-8 bits), subtract 1 for range 0-254

| value | 0 | 1 | 2 | 3 | 4 | 254 |
|------------|-----|------|------|-------|-------|----------------|
| value bits | 1 | 10 | 11 | 100 | 101 | 11111111 |
| encoding | 000 | 0010 | 0011 | 01000 | 01001 | 1111111111 |

- Icelandic 'takmarkað umfang', meaning 'limited range'
- Separate 'direction' bit D for exponent sign
- ▶ Interlude: Comparison for value 254 (posit \rightarrow takum):

Precision



Dynamic Range



Laslo Hunhold, James Quinlan

Four solvers

▶ LU decomposition (simulate UMFPACK → pregenerated permutations)

- LU decomposition (simulate UMFPACK \rightarrow pregenerated permutations)
- QR factorization (simulate SPQR \rightarrow pregenerated permutations)

- LU decomposition (simulate UMFPACK \rightarrow pregenerated permutations)
- QR factorization (simulate SPQR \rightarrow pregenerated permutations)
- Mixed Precision Iterative Refinement (MPIR) with ILU(0) preconditioning

- LU decomposition (simulate UMFPACK \rightarrow pregenerated permutations)
- QR factorization (simulate SPQR \rightarrow pregenerated permutations)
- Mixed Precision Iterative Refinement (MPIR) with ILU(0) preconditioning
- Generalized Minimal Residual (GMRES) method

Four solvers

- LU decomposition (simulate UMFPACK \rightarrow pregenerated permutations)
- QR factorization (simulate SPQR \rightarrow pregenerated permutations)
- Mixed Precision Iterative Refinement (MPIR) with ILU(0) preconditioning
- Generalized Minimal Residual (GMRES) method

Dataset

Four solvers

- LU decomposition (simulate UMFPACK \rightarrow pregenerated permutations)
- QR factorization (simulate SPQR \rightarrow pregenerated permutations)
- Mixed Precision Iterative Refinement (MPIR) with ILU(0) preconditioning
- Generalized Minimal Residual (GMRES) method

Dataset

SuiteSparse Matrix Collection
Four solvers

- LU decomposition (simulate UMFPACK \rightarrow pregenerated permutations)
- ▶ QR factorization (simulate SPQR → pregenerated permutations)
- Mixed Precision Iterative Refinement (MPIR) with ILU(0) preconditioning
- Generalized Minimal Residual (GMRES) method

Dataset

- SuiteSparse Matrix Collection
- Subset of 295 square, full rank matrices with less than 10⁴ non-zero entries

Four solvers

- ▶ LU decomposition (simulate UMFPACK → pregenerated permutations)
- ▶ QR factorization (simulate SPQR → pregenerated permutations)
- Mixed Precision Iterative Refinement (MPIR) with ILU(0) preconditioning
- Generalized Minimal Residual (GMRES) method

Dataset

- SuiteSparse Matrix Collection
- Subset of 295 square, full rank matrices with less than 10⁴ non-zero entries
- Custom metadata postprocessing and packing

Four solvers

- ▶ LU decomposition (simulate UMFPACK → pregenerated permutations)
- ▶ QR factorization (simulate SPQR → pregenerated permutations)
- Mixed Precision Iterative Refinement (MPIR) with ILU(0) preconditioning
- Generalized Minimal Residual (GMRES) method

Dataset

- SuiteSparse Matrix Collection
- Subset of 295 square, full rank matrices with less than 10⁴ non-zero entries
- Custom metadata postprocessing and packing

Four solvers

- LU decomposition (simulate UMFPACK \rightarrow pregenerated permutations)
- ▶ QR factorization (simulate SPQR → pregenerated permutations)
- Mixed Precision Iterative Refinement (MPIR) with ILU(0) preconditioning
- Generalized Minimal Residual (GMRES) method

Dataset

- SuiteSparse Matrix Collection
- Subset of 295 square, full rank matrices with less than 10⁴ non-zero entries
- Custom metadata postprocessing and packing

Procedure for each matrix A and type T

1. Generate random b such that $\left\|b\right\|_{\infty}=1$ (seeded Xoshiro PRNG)

Four solvers

- LU decomposition (simulate UMFPACK \rightarrow pregenerated permutations)
- ▶ QR factorization (simulate SPQR → pregenerated permutations)
- Mixed Precision Iterative Refinement (MPIR) with ILU(0) preconditioning
- Generalized Minimal Residual (GMRES) method

Dataset

- SuiteSparse Matrix Collection
- Subset of 295 square, full rank matrices with less than 10⁴ non-zero entries
- Custom metadata postprocessing and packing

- 1. Generate random b such that $\left\|b\right\|_{\infty} = 1$ (seeded Xoshiro PRNG)
- 2. Solve system Ax = b in float128 using custom sparse QR solver

Four solvers

- LU decomposition (simulate UMFPACK \rightarrow pregenerated permutations)
- ▶ QR factorization (simulate SPQR → pregenerated permutations)
- Mixed Precision Iterative Refinement (MPIR) with ILU(0) preconditioning
- Generalized Minimal Residual (GMRES) method

Dataset

- SuiteSparse Matrix Collection
- Subset of 295 square, full rank matrices with less than 10⁴ non-zero entries
- Custom metadata postprocessing and packing

- 1. Generate random b such that $\left\|b\right\|_{\infty} = 1$ (seeded Xoshiro PRNG)
- 2. Solve system Ax = b in float128 using custom sparse QR solver
- 3. Convert (A, b) to T, yielding (\tilde{A}, \tilde{b})

Four solvers

- ▶ LU decomposition (simulate UMFPACK → pregenerated permutations)
- QR factorization (simulate SPQR \rightarrow pregenerated permutations)
- Mixed Precision Iterative Refinement (MPIR) with ILU(0) preconditioning
- Generalized Minimal Residual (GMRES) method

Dataset

- SuiteSparse Matrix Collection
- Subset of 295 square, full rank matrices with less than 10⁴ non-zero entries
- Custom metadata postprocessing and packing

- 1. Generate random b such that $\left\|b\right\|_{\infty} = 1$ (seeded Xoshiro PRNG)
- 2. Solve system Ax = b in float128 using custom sparse QR solver
- 3. Convert (A, b) to T, yielding (\tilde{A}, \tilde{b})
- 4. Solve $\tilde{A}\tilde{x} = \tilde{b}$

Four solvers

- ▶ LU decomposition (simulate UMFPACK → pregenerated permutations)
- QR factorization (simulate SPQR \rightarrow pregenerated permutations)
- Mixed Precision Iterative Refinement (MPIR) with ILU(0) preconditioning
- Generalized Minimal Residual (GMRES) method

Dataset

- SuiteSparse Matrix Collection
- Subset of 295 square, full rank matrices with less than 10⁴ non-zero entries
- Custom metadata postprocessing and packing

- 1. Generate random b such that $\left\|b\right\|_{\infty} = 1$ (seeded Xoshiro PRNG)
- 2. Solve system Ax = b in float128 using custom sparse QR solver
- 3. Convert (A, b) to T, yielding (\tilde{A}, \tilde{b})
- 4. Solve $\tilde{A}\tilde{x} = \tilde{b}$
- 5. Determine $\|x \tilde{x}\|_2$ in float128

Four solvers

- LU decomposition (simulate UMFPACK \rightarrow pregenerated permutations)
- QR factorization (simulate SPQR \rightarrow pregenerated permutations)
- Mixed Precision Iterative Refinement (MPIR) with ILU(0) preconditioning
- Generalized Minimal Residual (GMRES) method

Dataset

- SuiteSparse Matrix Collection
- Subset of 295 square, full rank matrices with less than 10⁴ non-zero entries
- Custom metadata postprocessing and packing

Procedure for each matrix A and type T

- 1. Generate random b such that $\left\|b\right\|_{\infty} = 1$ (seeded Xoshiro PRNG)
- 2. Solve system Ax = b in float128 using custom sparse QR solver
- 3. Convert (A, b) to T, yielding (\tilde{A}, \tilde{b})
- 4. Solve $\tilde{A}\tilde{x} = \tilde{b}$
- 5. Determine $||x \tilde{x}||_2$ in float128

Sort all errors to obtain cumulative error distribution

LU (1/2) 8 and 16 bits



Laslo Hunhold, James Quinlan

LU (2/2) 32 and 64 bits



$\mathsf{QR}(1/2)$ 8 and 16 bits



QR(2/2) 32 and 64 bits



MPIR (1/2) (L, W, H) = (8, 16, 32) and (L, W, H) = (16, 16, 32), relative tolerance 10^{-3} and 10^{-3}



MPIR (2/2) (L, W, H) = (16, 32, 32) and (L, W, H) = (16, 32, 64), relative tolerance 10^{-6} and 10^{-9}



GMRES (1/2) 8 and 16 bits, restart value 20, relative tolerance $\sqrt{\varepsilon_{\text{float8}}}$ and $\sqrt{\varepsilon_{\text{float16}}}$



GMRES (2/2)

32 and 64 bits, restart value 20, relative tolerance $\sqrt{\varepsilon_{float32}}$ and $\sqrt{\varepsilon_{float64}}$



Laslo Hunhold, James Quinlan

Summary of Results

bfloat16 overall better than float16, but sometimes worse

- bfloat16 overall better than float16, but sometimes worse
- Posits and takums overall superior than IEEE 754 floats

- bfloat16 overall better than float16, but sometimes worse
- Posits and takums overall superior than IEEE 754 floats
- takum16 always better than bfloat16 (unlike posit16)

- bfloat16 overall better than float16, but sometimes worse
- Posits and takums overall superior than IEEE 754 floats
- takum16 always better than bfloat16 (unlike posit16)
- Takums significantly outperform posits in some benchmarks, especially GMRES (against intuition)

Summary of Results

- bfloat16 overall better than float16, but sometimes worse
- Posits and takums overall superior than IEEE 754 floats
- takum16 always better than bfloat16 (unlike posit16)
- Takums significantly outperform posits in some benchmarks, especially GMRES (against intuition)

Discussion

Summary of Results

- bfloat16 overall better than float16, but sometimes worse
- Posits and takums overall superior than IEEE 754 floats
- takum16 always better than bfloat16 (unlike posit16)
- Takums significantly outperform posits in some benchmarks, especially GMRES (against intuition)

Discussion

bfloat16 is a better general-purpose format than float16

Summary of Results

- bfloat16 overall better than float16, but sometimes worse
- Posits and takums overall superior than IEEE 754 floats
- takum16 always better than bfloat16 (unlike posit16)
- Takums significantly outperform posits in some benchmarks, especially GMRES (against intuition)

Discussion

- bfloat16 is a better general-purpose format than float16
- Posits have some shortcomings (precision loss further from 1, limited dynamic range), but overall better than IEEE 754 floats

Summary of Results

- bfloat16 overall better than float16, but sometimes worse
- Posits and takums overall superior than IEEE 754 floats
- takum16 always better than bfloat16 (unlike posit16)
- Takums significantly outperform posits in some benchmarks, especially GMRES (against intuition)

Discussion

- bfloat16 is a better general-purpose format than float16
- Posits have some shortcomings (precision loss further from 1, limited dynamic range), but overall better than IEEE 754 floats
- Takums suggest new mixed-precision workflow: Reducing n only affects precision, not dynamic range