


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

Laslo Hunhold 

Parallel and Distributed Systems Group  
University of Cologne  
Cologne, Germany  
hunhold@uni-koeln.de

James Quinlan 

Department of Computer Science  
University of Southern Maine  
Portland, ME, USA  
james.quinlan@maine.edu

**Abstract**—Solving sparse linear systems lies at the core of numerous computational applications. Consequently, understanding the numerical performance of recently proposed alternatives to the established IEEE 754 floating-point numbers, such as bfloat16 and the tapered-precision posit and takum machine number formats, is of significant interest. This paper examines these formats in the context of widely used solvers, namely LU, QR, and GMRES, with incomplete LU preconditioning and mixed precision iterative refinement (MPIR). This contrasts with the prevailing emphasis on designing specialized algorithms tailored to new arithmetic formats.

This paper presents an extensive and unprecedented evaluation based on the SuiteSparse Matrix Collection—a dataset of real-world matrices with diverse sizes and condition numbers. A key contribution is the faithful reproduction of SuiteSparse’s UMF-PACK multifrontal LU factorization and SPQR multifrontal QR factorization for machine number formats beyond single and double-precision IEEE 754. Tapered-precision posit and takum formats show better accuracy in direct solvers and reduced iteration counts in indirect solvers. Takum arithmetic, in particular, exhibits increased stability, even at low precision.

**Index Terms**—machine numbers, IEEE 754, floating-point arithmetic, tapered precision, posit arithmetic, takum arithmetic, sparse linear systems, direct solvers, indirect solvers

## I. INTRODUCTION

The numerical solution of sparse linear systems is a cornerstone problem in scientific computing, with applications encompassing structural analysis, circuit simulation, fluid dynamics, and machine learning. Historically, such computations have relied on the IEEE 754 floating-point standard [1], which has become the default format for numerical representation. However, the landscape is shifting towards low-precision formats to mitigate processor performance outpacing memory interconnect bandwidth in modern high-performance computing (commonly referred to as the “memory wall”).

Emerging number formats such as bfloat16 [2], posit [3], and takum [4] introduce opportunities to improve computational performance and accuracy, particularly in low-precision arithmetic. Posits and takums, for instance,

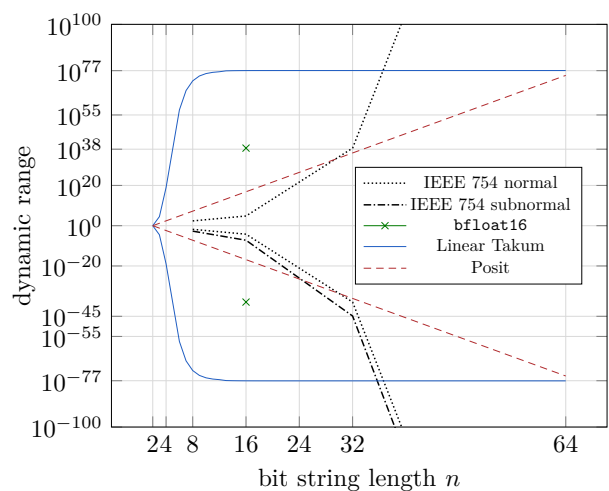


Figure 1: Dynamic range relative to the bit string length  $n$  for linear takum, posit and a selection of floating-point formats.

employ a tapered precision scheme through variable-width exponent encoding, which allocates higher precision to values near 1 while sacrificing precision for values further from 1. Takum arithmetic represents a novel advancement over posits by offering an extensive dynamic range even at very low precisions. This design is motivated by the principle that bit-string length should primarily determine precision without imposing constraints on dynamic range—a common limitation in other formats. This paper focuses on the linear takum variant, which is a floating-point format, as opposed to the logarithmic representation in the standard takum format. For the `float8` format, which lacks standardization, we adopt the E4M3 OFP8 definition provided in [5] with 4 exponent bits and 3 fraction bits. Figure 1 illustrates the dynamic ranges of the formats evaluated in this study.

Despite increasing interest in posits for numerical analysis [6], [7], [8], no prior work has examined takums in this context, given their novelty. Moreover, a systematic, large-

scale comparison of posits and takums against bfloat16, another promising low-precision format, is currently lacking. A key question in this domain is how the expansive constant dynamic range of takums compares to that of posits, which, despite their potential, have faced adoption challenges and criticism for their limited dynamic range [9].

In this paper, we evaluate the numerical performance of these alternative number formats within a selection of established solvers that underpin scientific computing software. Our analysis is based on a large, diverse, sparse matrix test set, representing an unprecedented scale in such studies. We avoid tailoring algorithms to any specific number format to ensure unbiased results. Instead, we fully reproduce the SuiteSparse library with respect to the UMFPACK LU solver and the SPQR QR solver. This approach simulates a “blind” replacement of the underlying arithmetic in a computing environment, offering a more realistic assessment than tailored implementations. Additionally, we evaluate GMRES and mixed-precision iterative refinement methods, extending the latter to 8-bit precision—an exploration that, to the best of our knowledge, has not been previously investigated in the literature. Since most of the arithmetic formats examined in this study are implemented in software, neither computation time nor power consumption is measured; the analysis focuses solely on numerical performance.

The remainder of this paper is organized as follows. Section II outlines the experimental methods used to benchmark the formats. Section III presents the main results, including detailed analyses and visualizations. Finally, Section IV summarizes our findings and offers conclusions.

## II. EXPERIMENTAL METHODS

We evaluate the numerical performance of four fundamental approaches to solving sparse linear systems across multiple numeric formats: LU decomposition and QR factorization as direct methods, and the Generalized Minimal Residual (GMRES) method with incomplete LU preconditioning alongside Mixed Precision Iterative Refinement (MPIR) as iterative methods. These approaches encompass core techniques in numerical linear algebra, each distinguished by unique trade-offs in computational efficiency, memory consumption, and numerical stability.

The benchmarking framework presented in this work, MuFoLAB (Multi-Format Linear Algebra Benchmarks) [10], is designed to facilitate systematic and reproducible evaluations. It comprises three key components: a test matrix generator, a unified experimental interface for solvers, and implementations of the four solver methods under consideration. The subsequent sections provide a detailed description of these components and their roles in the benchmarking process.

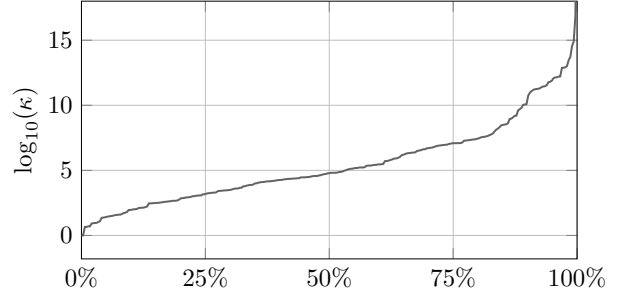


Figure 2: Cumulative distribution of test matrix  $L^1$  condition numbers.

### A. Test Matrices Generation

The first step in the experimental setup involves preparing a comprehensive set of sparse test matrices for benchmarking. To achieve this, we draw matrices from the SuiteSparse Matrix Collection [11], a well-established repository containing matrices from diverse application domains, including computational fluid dynamics, chemical simulation, materials science, optimal control, structural mechanics, and 2D/3D sequencing. Initially, we discard non-real matrices and those with more than  $10^4$  non-zero entries, resulting in a preliminary dataset of 833 matrices. Further refinement, limited to square matrices with full rank — criteria necessary for our benchmarks — reduces the dataset to 295 matrices. The  $L^1$  condition numbers of these matrices span several orders of magnitude, with a median value of approximately  $10^5$ , and around 25% of the matrices exhibit condition numbers exceeding  $10^7$ , as illustrated in Figure 2. It should be noted that the size of the matrices is far less significant than the number of nonzero entries, as the sparse linear algebra algorithms examined in this study disregard zero entries, and any differences in fill-in affect all numerical formats equally. Furthermore, the number of nonzero entries in the matrices is limited not by inherent constraints of the number formats but to maintain reasonable computation times and reproducibility, given that the number formats are implemented in software.

The established Julia package for accessing the SuiteSparse Matrix Collection (`MatrixDepot.jl`) retrieves matrices on-demand via individual internet requests. While this approach is functional in local setups with caching, it becomes unsuitable for cold-cache deployments on high-performance computing (HPC) systems and containerized environments, especially given the scale of our benchmarks, which involve hundreds of matrices. Moreover, the metadata provided by the package is frequently incomplete.

To address these challenges, we introduce a second processing step to streamline and enhance matrix preparation. In this step, all filtered matrices are converted to sparse `float64` format, and complete metadata—including the

number of non-zero entries, absolute minimum and maximum values,  $L^1$  condition number, rank, symmetry, and positive definiteness—is computed. The processed matrices and their metadata are stored in a single compressed Julia Data File (JLD2). This self-contained file significantly improves deployment efficiency by enabling rapid access to test matrices with consistent metadata without requiring an internet connection [10, src/TestMatrices-Generator.jl, src/TestMatrices.jl].

### B. Common Solver Experiment Interface

With the test matrices prepared, the next step is establishing a standardized interface for evaluating different numeric formats within a selection of solvers. Given only the matrix  $A$ , constructing a linear test system  $Ax = b$  requires generating the solution vector  $x \in \mathbb{R}^n$  and the corresponding right-hand side  $b \in \mathbb{R}^n$ .

Although a common approach is to set  $x = (1, \dots, 1)^T$ , a more representative method, as outlined in [12, Section 5], involves generating  $b$  randomly such that  $\|b\|_\infty = 1$  and computing the reference solution by solving the system  $Ax = b$  in `float128` using `Quadmath.jl` and a custom type-agnostic sparse QR solver [10, src/QR.jl]. Random generation is performed using a Xoshiro pseudorandom number generator (PRNG) with a fixed seed to ensure reproducibility.

For each numeric format under evaluation, the matrix  $A$  and right-hand side  $b$  are converted to the target type, denoted as  $\tilde{A}$  and  $\tilde{b}$ , respectively. If any entry in  $\tilde{A}$  or  $\tilde{b}$  underflows or overflows during conversion, the experiment is aborted, and the failure is recorded in the results. For further discussion on dynamic range and its implications in low-precision arithmetic, refer to [13]. For the target types we use the Julia packages `Float8s.jl`, `BFloat16s.jl`, `Posits.jl` and `Takums.jl`.

The respective solver is applied to the linear system  $\tilde{A}x = \tilde{b}$ , yielding an approximate solution  $\tilde{x}$ . This solution, cast back to `float128`, is compared to the exact solution  $x$ , and the absolute and relative 2-norm errors are calculated. While it may seem counterintuitive, generating only one random sample for each matrix is sufficient because the ensemble diversity is provided by the large number of matrices in the test set. The same random seed also ensures each number format is tested with the same pseudorandom outcome [10, src/Experiment.jl].

### C. LU Solver

Although the LU decomposition algorithm is relatively straightforward to implement, the primary challenge lies in determining effective row and column permutations to minimize fill-in and enhance numerical performance. UMFPACK [14] is widely used and highly optimized for LU decompositions. It employs a sophisticated rule set to select an optimal pivoting strategy, producing row and column permutations, as well as a row scaling matrix [14, Section 1].

However, a significant limitation is that UMFPACK is implemented exclusively for `float32` and `float64` data types. To extend its capabilities to other numeric formats, we emulate UMFPACK’s behavior by precomputing an LU decomposition in `float64` using UMFPACK for each test matrix. This computation yields the row and column permutations, which depend solely on the structural properties of the matrix and not on the specific type of its elements. The UMFPACK row scaling, however, must be computed separately for each number format, as it depends on the numerical properties of the data type. This is relatively straightforward, as it simply involves summing the absolute values in each row, yielding the vector of row 1-norms (cf. [14, p. 49]).

Once the precomputed permutations and scaling factors are determined, they are applied to the matrix in the target numeric format. The system is then solved using a simple non-pivoting LU solver, effectively replicating UMFPACK’s behavior. This approach ensures consistency in the decomposition process across all tested numeric formats while maintaining parity with UMFPACK’s sophisticated pivoting strategy.

### D. QR Solver

Similar to the case of LU decomposition with UMFPACK, the core algorithm for QR decomposition, which employs HOUSEHOLDER rotations, is relatively straightforward. However, determining optimal row and column permutations to minimize fill-in during the decomposition is the primary challenge. This challenge is addressed by SPQR, a highly optimized implementation that, like UMFPACK, is part of the SuiteSparse library [15]. SPQR is designed to maximize numerical efficiency and reduce memory requirements through sophisticated permutation strategies.

A significant limitation of SPQR is its exclusive implementation for the `float64` and `float32` data types. To enable the use of other number formats, we leverage the fact that row and column permutations are solely dependent on the structural properties of the matrix and not on the specific numeric type of its elements. Accordingly, we precompute a QR decomposition in `float64` using SPQR to extract the optimal row and column permutations.

Once the matrix is permuted according to these precomputed permutations, we collect all non-zero entries below the diagonal and apply one HOUSEHOLDER rotation per column. We store both the rotation vector and the vector of indices affected by each rotation. This process effectively emulates the behavior of SPQR while allowing for a broader range of number formats, ensuring consistency and efficiency in the decomposition across all tested data types.

### E. Mixed Precision Iterative Refinement (MPIR) Solver

The iterative refinement technique [16] is a classical approach to improving an approximate solution  $\tilde{x}$  to a

linear system  $Ax = b$ . The method iteratively refines  $\tilde{x}$  by solving a correction equation  $Ac = r$ , where  $r = b - A\tilde{x}$  is the residual, and updating  $\tilde{x} \leftarrow \tilde{x} + c$ . This process is repeated until a convergence criterion is met, which, in our implementation, is based on the normwise backward error [17]. However, iterative refinement may fail to converge if the low-precision arithmetic causes  $A$  to appear singular, thereby preventing the accurate computation of  $\tilde{x}$ .

Mixed-precision iterative refinement (MPIR) extends this method by employing different levels of precision to optimize computational efficiency and solution accuracy [18]. Specifically, MPIR uses three distinct precision levels:

- 1) *Working precision* ( $W$ ), where  $A$ ,  $b$ , and  $\tilde{x}$  are stored.
- 2) *Low precision* ( $L$ ) for the factorization of  $A$ , typically to reduce computational cost.
- 3) *High precision* ( $H$ ) for the residual calculation, ensuring accurate error correction.

These precision levels are collectively represented as a triple  $(L, W, H)$ .

We evaluated several precision configurations, including a novel (8, 16, 32) setup, alongside established configurations such as (16, 16, 32), (16, 32, 32), and (16, 32, 64). These configurations were tested across multiple number formats, including IEEE floating-point, **bfloat16**, linear takums, and posits.

For each configuration, the error tolerance was adjusted to align with the precision levels:  $10^{-3}$  for (8, 16, 32) and (16, 16, 32),  $10^{-6}$  for (16, 32, 32), and  $10^{-9}$  for (16, 32, 64). The maximum number of iterations was set to 100 for all experiments.

#### F. Incomplete LU Preconditioned GMRES Solver

Our implementation employs left-looking level 0 incomplete LU factorization ( $\text{ILU}(0)$ ) via **IncompleteLU.jl** to reduce fill-in and use the resulting factors as a preconditioner in the Generalized Minimal Residual (GMRES) method. This combination balances computational efficiency and accelerated convergence for sparse linear systems while keeping memory requirements manageable.

We leverage the **IterativeSolvers.jl** package, adhering largely to the GMRES default parameters. Specifically, we use a restart value of  $\min(20, n)$ , a maximum iteration count of  $n$ , and the modified GRAM-SCHMIDT process for orthogonalization. For the relative tolerance, which defaults to the square root of the machine epsilon of the working precision, we instead use the square root of the machine precision of the corresponding reference float type. This approach ensures fairness across all numeric formats. For example, **float8** is used for all 8-bit types, **float16** for all 16-bit types, **float32** for all 32-bit types, and **float64** for all 64-bit types. Without this adjustment, numeric types with smaller machine epsilons than their IEEE 754 counterparts would be disproportionately disadvantaged in terms of convergence criteria.

### III. RESULTS

The results for the LU solver are presented in Figure 3. As shown, both **posit8** and, to an even greater extent, **takum\_linear8** significantly outperform **float8** in terms of solution accuracy. This pattern persists across 16, 32, and 64 bits, with posits and takums consistently surpassing or at least matching the corresponding IEEE 754 floating-point types.

An especially noteworthy observation is that **takum\_linear16** consistently outperforms **bfloat16**, whereas **posit16** exhibits reduced accuracy for the lower quartile of matrices. Interestingly, while **bfloat16** generally achieves higher accuracy than **float16**, it is less accurate for approximately 25% of the test matrices. This behavior underscores the nuanced trade-offs among these numeric formats and suggests that takums offer greater dependability in challenging cases.

The results for the QR solver, shown in Figure 4, exhibit a similar overall trend. Both posits and takums consistently outperform or match their corresponding IEEE 754 floating-point counterparts. Notably, **takum\_linear16** consistently surpasses **bfloat16**, showing increased accuracy across all test cases.

The mixed precision iterative refinement results are presented in Figure 5. Overall, both posits and takums demonstrate significantly lower iteration counts and fewer occurrences of singularities or maximum iteration limit exceedances compared to their respective IEEE 754 floating-point counterparts. When comparing posits and takums, no clear advantage of one format emerges, as their numerical performance appears comparable across the evaluated test cases.

For GMRES preconditioned with incomplete LU, the results displayed in Figure 6 highlight significant differences in performance across numeric formats. While **float8** frequently experiences overflows or requires a high number of iterations, both **posit8** and **takum\_linear8** exhibit much greater numerical stability. Notably, **takum\_linear8** avoids overflow entirely for all test matrices.

This trend persists at 16 bits, where **takum\_linear16** consistently achieves lower iteration counts than **bfloat16**, in contrast to **posit16**, which occasionally lags behind. At 32 and 64 bits, posits and takums demonstrate very similar performance, both achieving significantly fewer iterations compared to **float32** and **float64**. These results underscore the advantages of tapered-precision formats in reducing computational effort and enhancing stability.

### IV. CONCLUSION

We evaluated IEEE 754 floating-point numbers, **bfloat16**, posits, and takums across four widely used direct and iterative solving algorithms. Our experiments demonstrate that tapered-precision arithmetic consistently outperforms IEEE 754 floating-point numbers in all tested scenarios. Among the tapered-precision formats,

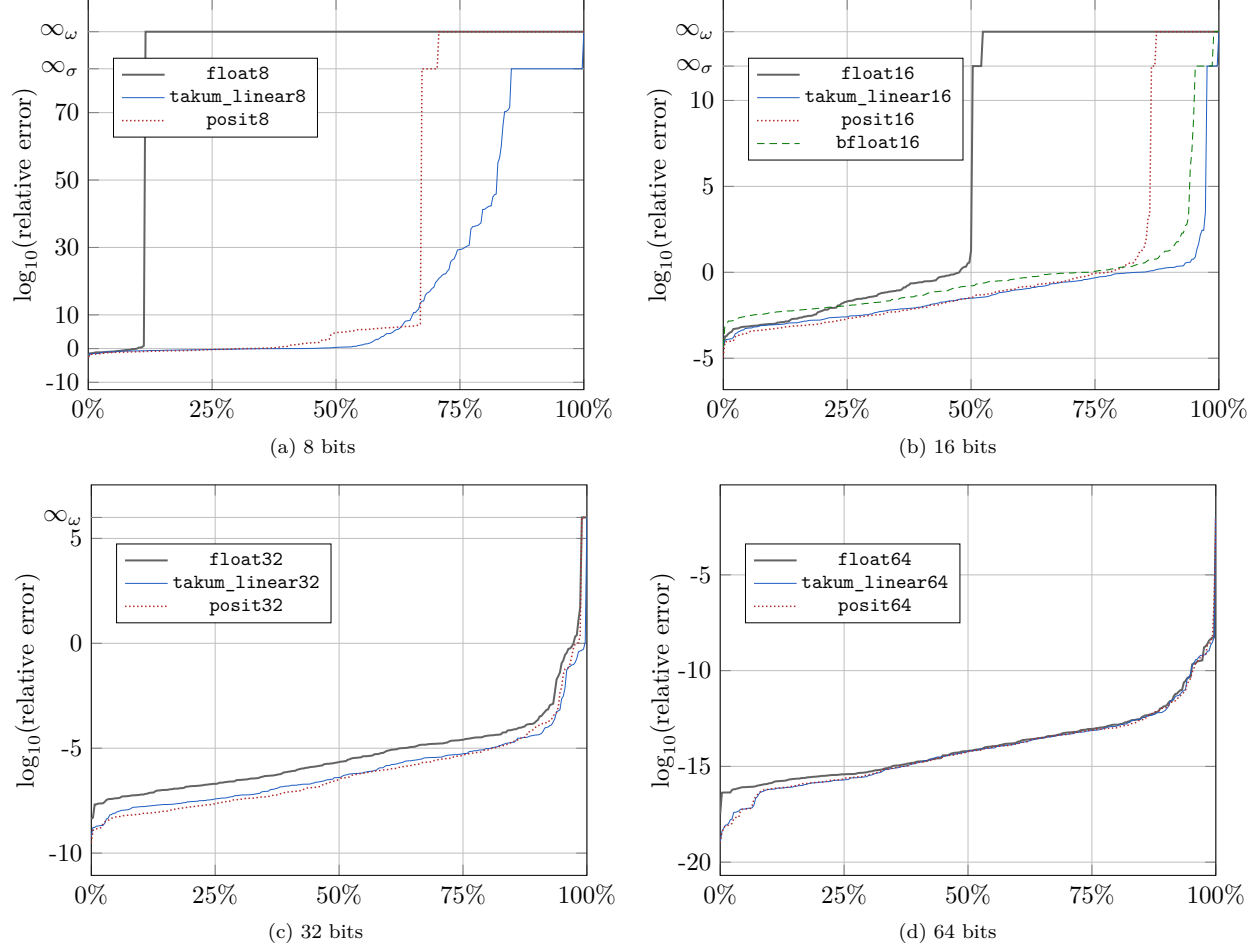


Figure 3: Cumulative error distribution of the relative errors of the solutions of the linear systems via fully pivoted LU decomposition using a range of machine number types. The symbol  $\infty_\sigma$  denotes where the conversion of the matrix to the target number type turned it singular,  $\infty_\omega$  denotes where the dynamic range of the matrix entries exceeded the target number type.

takums exhibited exceptional performance, outperforming **bfloat16** in every case. While occasionally marginally less accurate than posits, takums delivered comparable results overall and demonstrated superior numerical stability. Notably, we successfully introduced the application of 8-bit posits and takums in mixed-precision iterative refinement, marking a possibly significant milestone in numerical computing. Additionally, GMRES exhibited particular benefit from tapered-precision formats, with takums delivering outstanding results, surpassing posits in all cases.

These findings are especially relevant as they position takums as a strong candidate to replace **bfloat16** as the state-of-the-art in 16-bit arithmetic, which posits were unable to given their limited dynamic range. Furthermore, the results address a critical question: despite having a much larger dynamic range than other number formats,

including posits (see Figure 1), takums exhibit comparable and often more favorable results. This property is possibly transformative for mixed-precision workflows, as the choice of the precision level  $n$  with takums becomes purely about precision, decoupled from concerns about dynamic range.

Future research can explore further optimizations for MPIR using equilibrated matrices and investigate GMRES-based iterative refinement employing more than three precision levels [19].

#### AUTHOR CONTRIBUTIONS

**Laslo Hunhold:** Conceptualization, Data curation, Formal Analysis, Funding acquisition, Investigation, Methodology, Project administration, Resources, Software, Supervision, Validation, Visualization, Writing –

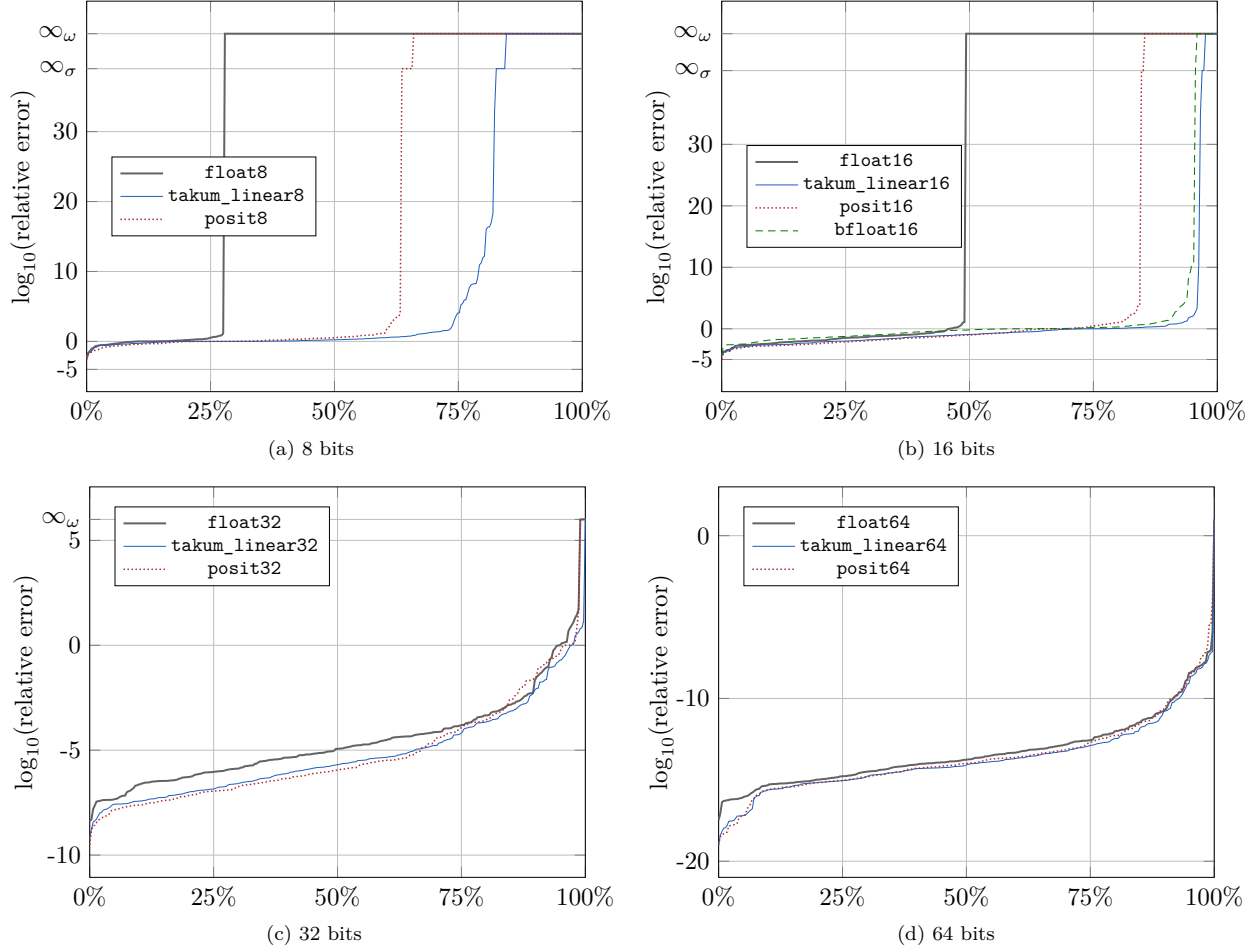


Figure 4: Cumulative error distribution of the relative errors of the solutions of the linear systems via QR decomposition using a range of machine number types. The symbol  $\infty_\sigma$  denotes where the conversion of the matrix to the target number type turned it singular,  $\infty_\omega$  denotes where the dynamic range of the matrix entries exceeded the target number type.

original draft, Writing – review & editing; **James Quinlan**: Conceptualization, Formal Analysis, Funding acquisition, Investigation, Methodology, Resources, Software, Writing – original draft, Writing – review & editing

#### REFERENCES

- [1] “IEEE Standard for Floating-Point Arithmetic,” *IEEE Std 754-2019 (Revision of IEEE 754-2008)*, pp. 1–84, 2019. DOI: 10.1109/ieeestd.2019.8766229.
- [2] S. Wang and P. Kanwar, “BFloat16: The secret to high performance on cloud TPUs,” Aug. 2019. [Online]. Available: <https://web.archive.org/web/20190826170119/https://cloud.google.com/blog/products/ai-machine-learning/bfloat16-the-secret-to-high-performance-on-cloud-tpus>.
- [3] J. L. Gustafson and I. Yonemoto, “Beating Floating Point at Its Own Game: Posit Arithmetic,” *Supercomputing Frontiers and Innovations*, vol. 4, no. 2, pp. 71–86, Jun. 2017. DOI: 10.14529/jsfi170206.
- [4] L. Hunhold, “Beating Posits at Their Own Game: Takum Arithmetic,” in *Next Generation Arithmetic, 5th International Conference, CoNGA 2024, Sydney, NSW, Australia, February 20–21, 2024, Proceedings*, ser. Lecture Notes in Computer Science, vol. 14666, Sydney, NSW, Australia: Springer Nature Switzerland, Oct. 2024. DOI: 10.1007/978-3-031-72709-2\_1.
- [5] P. Micikevicius *et al.*, “OCP 8-bit Floating Point Specification (OFP8),” Jun. 2023. [Online]. Available: <https://web.archive.org/web/20231017223546/https://www.opencompute.org/>

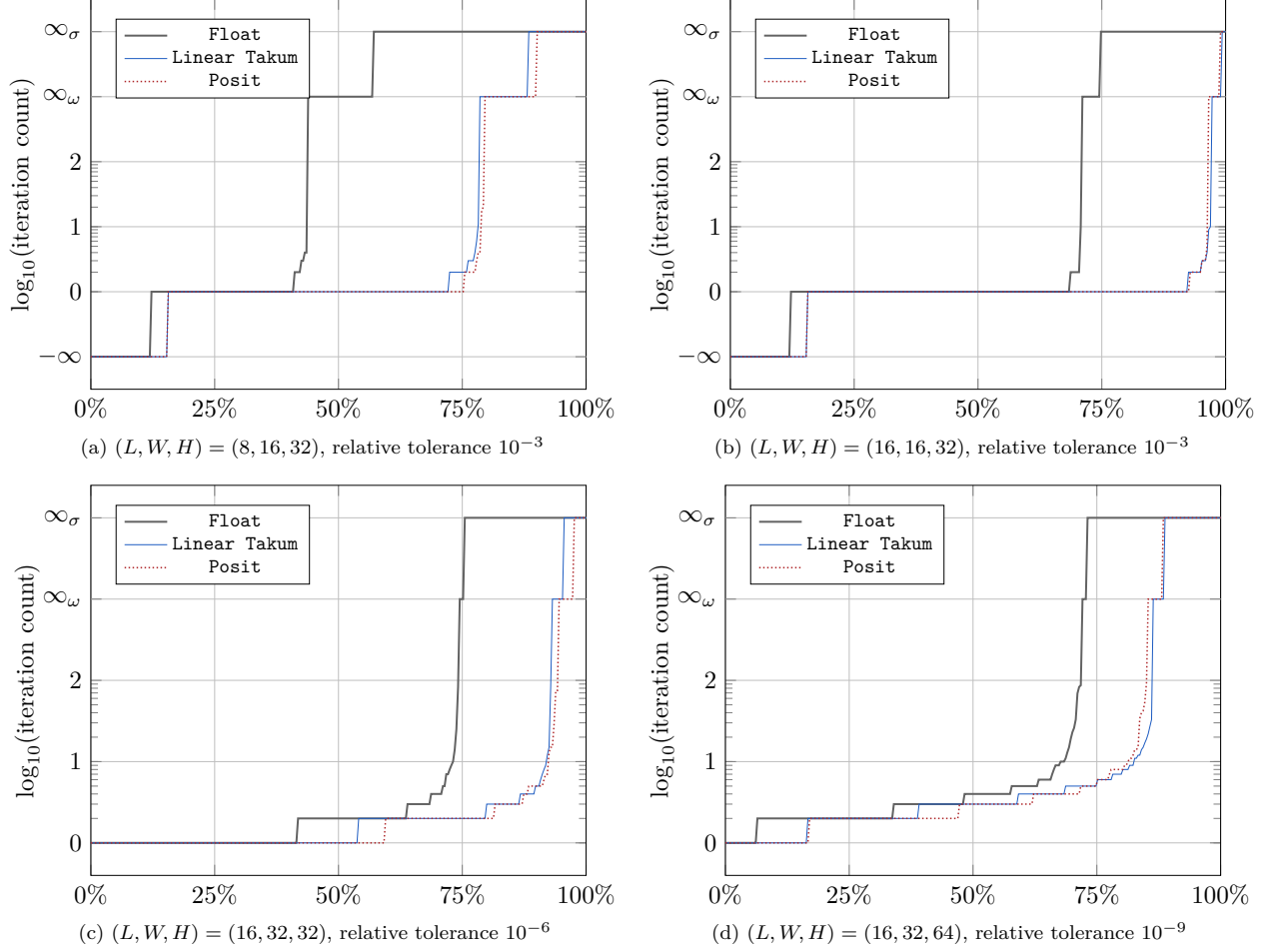


Figure 5: Cumulative distribution of the MPIR iteration counts using a range of machine number types. The symbol  $\infty_\sigma$  denotes where the initial low-precision LU decomposition yielded a singular system,  $\infty_\omega$  denotes where the maximum iteration count was reached without the residual going below the desired relative tolerance.

- documents/ocp-8-bit-floating-point-specification-ofp8-revision-1-0-2023-06-20-pdf.
- [6] N. Buon cristiani *et al.*, “Evaluating the numerical stability of posit arithmetic,” in *2020 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, 2020, pp. 612–621. DOI: 10.1109/IPDPS47924.2020.00069.
  - [7] D. Mallasén *et al.*, “Big-PERCIVAL: Exploring the native use of 64-bit posit arithmetic in scientific computing,” *IEEE Transactions on Computers*, vol. 73, no. 6, pp. 1472–1485, 2024. DOI: 10.1109/TC.2024.3377890.
  - [8] J. Quinlan and E. T. L. Omtzigt, “Iterative refinement with low-precision posit arithmetic,” in *Conference on Next Generation Arithmetic*, Springer, 2024, pp. 74–90.
  - [9] F. de Dinechin *et al.*, “Posits: The Good, the Bad and the Ugly,” ser. CoNGA’19, Singapore, Singapore: Association for Computing Machinery, 2019. DOI: 10.1145/3316279.3316285.
  - [10] L. Hunhold and J. Quinlan, *MuFoLAB - Multi-Format Linear Algebra Benchmarks*, version v1.1.0, Mar. 2025. DOI: 10.5281/zenodo.14984597.
  - [11] T. A. Davis and Y. Hu, “The University of Florida Sparse Matrix Collection,” *ACM Transactions on Mathematical Software (TOMS)*, vol. 38, no. 1, pp. 1–25, 2011.
  - [12] E. Carson and N. J. Higham, “A new analysis of iterative refinement and its application to accurate solution of ill-conditioned sparse linear systems,” *SIAM Journal on Scientific Computing*, vol. 39, no. 6, A2834–A2856, 2017.

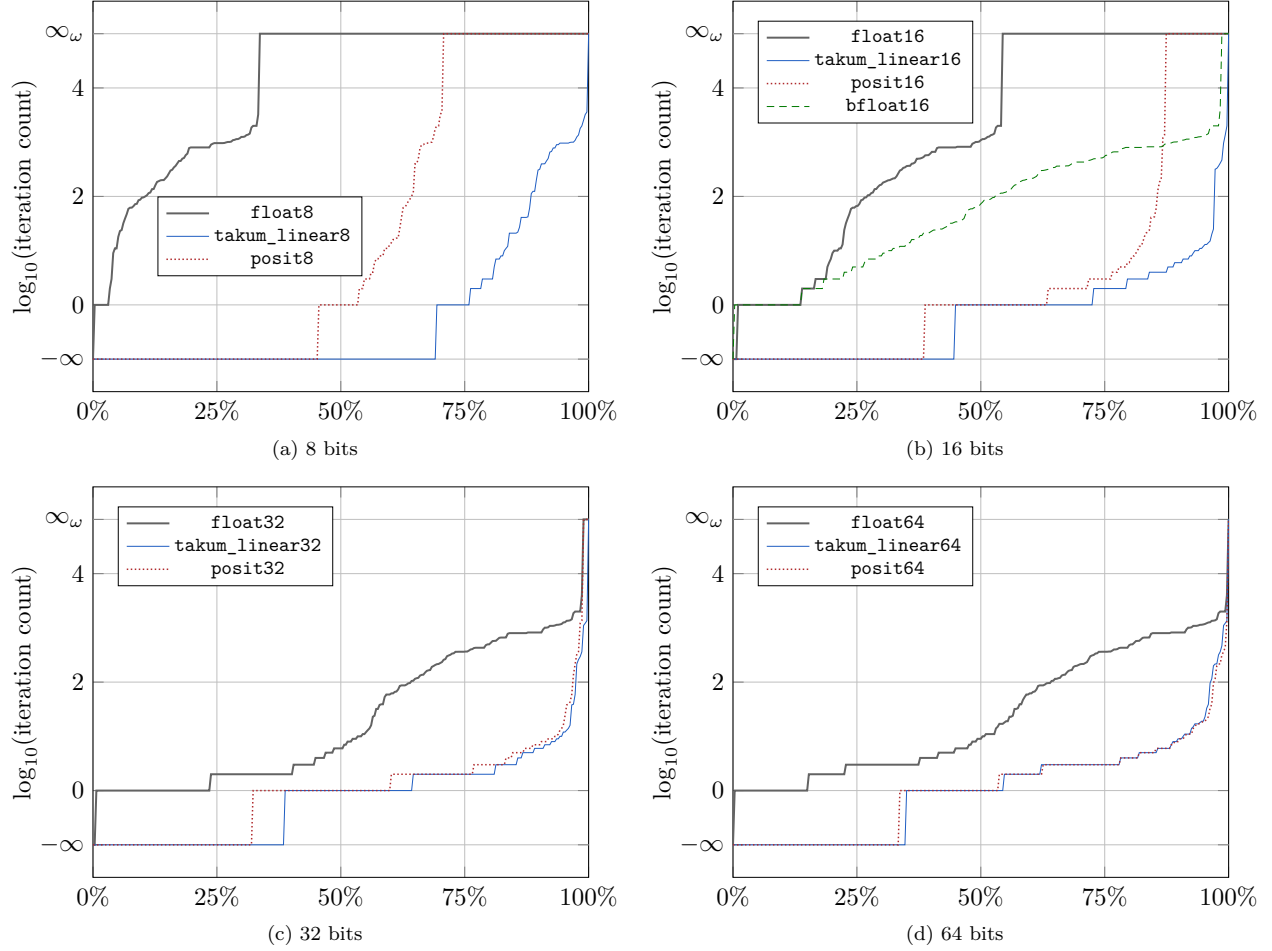


Figure 6: Cumulative distribution of the GMRES iteration counts using a range of machine number types. The symbol  $\infty$  denotes where the maximum iteration count was exceeded without reaching the desired tolerance.

- [13] N. J. Higham *et al.*, “Squeezing a Matrix into Half Precision, with an Application to Solving Linear Systems,” *SIAM Journal on Scientific Computing*, vol. 41, no. 4, A2536–A2551, 2019. doi: 10.1137/18M1229511.
- [14] T. A. Davis, “UMFPACK User Guide,” Jun. 2024. [Online]. Available: [https://web.archive.org/web/20241005011222/https://fossies.org/linux/SuiteSparse/UMFPACK/Doc/UMFPACK\\_UserGuide.pdf](https://web.archive.org/web/20241005011222/https://fossies.org/linux/SuiteSparse/UMFPACK/Doc/UMFPACK_UserGuide.pdf).
- [15] T. A. Davis *et al.*, “User’s Guide for SuiteSparseQR, a multifrontal multithreaded sparse QR factorization package (with optional GPU acceleration),” Jun. 2024. [Online]. Available: [https://web.archive.org/web/20241221222735/https://fossies.org/linux/SuiteSparse/SPQR/Doc/spqr\\_user\\_guide.pdf](https://web.archive.org/web/20241221222735/https://fossies.org/linux/SuiteSparse/SPQR/Doc/spqr_user_guide.pdf).
- [16] J. H. Wilkinson, *Progress report on the automatic computing engine*. Mathematics Division, Department of Scientific & Industrial Research, National Physical Laboratory, 1948.
- [17] N. J. Higham, *Accuracy and stability of numerical algorithms*. SIAM, 2002.
- [18] E. Carson and N. J. Higham, “Accelerating the solution of linear systems by iterative refinement in three precisions,” *SIAM Journal on Scientific Computing*, vol. 40, no. 2, A817–A847, 2018.
- [19] P. Amestoy *et al.*, “Five-Precision GMRES-based iterative refinement,” *SIAM Journal on Matrix Analysis and Applications*, vol. 45, no. 1, pp. 529–552, 2024.