

# Balancing Inexactness in Mixed Precision Matrix Computations

Erin Claire Carson<sup>[0000-0001-9469-7467]</sup>

## 1 Introduction

The scale of available computational power is constantly growing. As of 2022, we have entered the “exascale era”, meaning that we have computers capable of performing  $10^{18}$  (a billion billion) double precision floating point operations per second. This level of computing power offers massive opportunity for scientific simulation and discovery, but we still have a significant challenge ahead, in that we need to actually design algorithms and software for real-world applications that can take advantage of this powerful underlying hardware.

The focus of this manuscript is on basic numerical linear algebra and matrix computations, which form the core of a number of critical large-scale applications, including computational fluid dynamics, climate and weather modeling, image analysis, AI, etc. A critical point is that in all these applications, the matrix computations that we perform are necessarily *inexact* matrix computations.

There are many different sources of inexactness that arise in the computational science process. We typically start with a given phenomenon to study and we write down a mathematical model that describes it. We of course can’t capture the full reality using a model, so we incur some *modeling error*. In order to actually solve our model equations and perform a simulation, we incur a *discretization error* and/or a *linearization error*. We eventually arrive at a linear algebra problem, and must pick an algorithm to solve it. This involves what we will call *algorithmic approximation errors*. For example, we might approximate our operator or use randomization to make the problem computationally feasible, or we might pick a certain stopping criterion for an iterative method based on the needs of the application. Finally, when we go to run this algorithm on real hardware, we necessarily incur *rounding errors* due to the use of finite precision.

---

Erin Claire Carson  
Faculty of Mathematics and Physics, Charles University, Czech Republic, e-mail:  
carson@karlin.mff.cuni.cz

These errors are intimately connected. For example, in order to determine an appropriate stopping criterion for an iterative method, one must have some understanding of its limits in finite precision, and the particular discretization used will result in different matrix properties, which can change the convergence behavior of an iterative solver. Unfortunately, these different sources of error are almost always studied in isolation. Further, the rounding error is frequently ignored completely.

Not accounting for rounding error may have made sense for some select applications in the past when all computations were run in a uniform double (64-bit) precision. On today's high-performance compute hardware, however, which features a variety of bit-width formats, ignoring the rounding error and how it interacts with other sources of inexactness is not only potentially hazardous, but is also a lost opportunity for improving the performance of scientific codes by exploiting the capabilities of today's hardware.

## 2 Floating Point Numbers and Mixed Precision Hardware

In computer hardware as we know it, we only have a finite number of bits to store and compute with numbers. The industry standard are the IEEE floating point numbers. A base-2 floating point number can be represented as

$$(-1)^{\text{sign}} \times 2^{(\text{exponent}-\text{offset})} \times 1.\text{significant}.$$

There is one sign bit, some number of exponent bits, and some number of bits for the significand (also called mantissa or fraction), depending on the particular format (which also defines the constant offset).

The number of bits assigned to the exponent will determine the range of representable numbers, and the number of bits assigned to the significand will affect the unit roundoff  $u$ , which is the precision with which numbers are stored. For example, for IEEE 754 double precision (64 bits), we have 11 exponent bits with an offset of 1023, and 52 explicitly stored significand bits (plus the implicit leading one). This means that the largest representable number is  $2^{(2^{\#\text{exponent bits}}-1)-1023} = 2^{1024} \approx 2 \cdot 10^{308}$ , the smallest (normal) representable number is  $2^{1-1023} = 2^{-1022} \approx 2 \cdot 10^{-308}$ , and the unit roundoff is  $2^{-\#\text{significand bits}} = 2^{-53} \approx 10^{-16}$ .

Whenever we perform a computation with two floating point numbers, the result is rounded to a floating point number, and this incurs a relative error that is bounded by the unit roundoff, i.e., for floating point numbers  $x$  and  $y$ ,

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

where  $\text{fl}()$  denotes the result of a computation performed in floating point arithmetic.

On modern high-performance computer hardware, we have many more precisions to work with than just double. For example, on the NVIDIA H100 GPU (featured in, for example, the JUPITER exascale computer), we have many different formats, not all of which are from the IEEE 754 standard. These formats along with their

properties are displayed in Table 1. From the table we can see that of course, with fewer bits, we generally have less to work with in both the exponent and significand fields, and thus we generally end up with a smaller range of representable numbers and a greater unit roundoff. From the final column, however, we can see that the potential performance gains from using lower precision are massive.

The performance advantages of low precision have been highlighted through the relatively recent addition of the HPL-MxP benchmark [19], a complement to the HPL benchmark in the TOP500 ranking [24]. As of writing, the current fastest computer in the world is El Capitan at Lawrence Livermore National Laboratory, which achieves 1.8 exaflops/s on the HPL benchmark, which is doing Gaussian elimination with partial pivoting to solve a dense linear system to double precision accuracy. The HPL-MxP benchmark also solves a dense linear system to double precision accuracy, but allows the use of *mixed precision* iterative refinement to do so. On this benchmark, El Capitan reaches 16.7 effective exaflops/s, demonstrating the power that comes from exploiting the available low precision hardware.

Indeed, this is a growing trend. The number of systems in the TOP500 ranking that feature mixed precision accelerators is only increasing over time, and further, on the top machines, the vast majority of floating point performance comes from accelerators; see, e.g., [9]. If we want to make effective use of these computers, we must thus redesign our algorithms so that they can make use of the mixed precision hardware available on GPUs and other specialized accelerators.

**Table 1** Floating point number formats on the NVIDIA H100 GPU [22].

	size(bits)	range	$u$	Tflops performance (NVIDIA H100 TC)
fp64	64	$10^{\pm 308}$	$1 \cdot 10^{-16}$	67
fp32	32	$10^{\pm 38}$	$6 \cdot 10^{-8}$	989
tf32	19	$10^{\pm 38}$	$5 \cdot 10^{-4}$	989
fp16	16	$10^{\pm 5}$	$5 \cdot 10^{-4}$	1979
bf16	16	$10^{\pm 38}$	$4 \cdot 10^{-3}$	1979
fp8-e5m2	8	$10^{\pm 5}$	$1 \cdot 10^{-1}$	3958
fp8-e4m3	8	$10^{\pm 2}$	$6 \cdot 10^{-2}$	3958

### 3 The Challenges of Low Precision

While it is clear we need to use low precision if we want to get anywhere close to peak performance on today’s machines, this is not something we can do blindly.

**Reduced Numerical Range** The first difficulty is that, as seen in Table 1, lower precisions have a smaller range of representable numbers. In floating point error analyses, it is often assumed that no overflow or underflow occurs during a computation. When using low precision, however, this assumption may not be valid. Overflow will usually cause a computation to fail entirely; underflow may be innocuous in some cases, but in others, it can cause us to lose important numerical properties, such as nonsingularity or positive definiteness.

The reduced numerical range of low precision brings challenges on many levels. From a software engineering perspective, one must carefully implement code to catch and resolve these errors; see, e.g., [23]. From an analysis and algorithm development perspective, we would like to identify properties of the problem and precision which are likely to result in underflow or overflow, and develop techniques to mitigate this. One example is the sophisticated scaling and shifting approach developed in [18] to “squeeze” a matrix into lower precision.

**Invalid Error Bounds** We also must consider whether existing analysis of stability and accuracy for our algorithms still apply when we use finite precision. In many cases, bounds on the backward or forward errors for a particular algorithm will usually hold only under some assumption like  $nu < 1$ , where  $n$  is the problem dimension. Using lower precision means that the unit roundoff  $u$  is larger, which means that error bounds only hold for much smaller problems. For example, in half precision where  $u \approx 10^{-4}$ , we have no guarantee of stability for problems with dimension  $10^4$  or larger (a relatively small problem in the world of exascale).

The question arises as to whether this is an actual limitation of the algorithm or the arithmetic system, or whether such conditions are overly restrictive as a result of the worst-case fashion of standard error analyses. For example, it has been shown that, using probabilistic approaches, one can obtain bounds in which  $n$  is replaced by  $\sqrt{n}$  (or even by 1 in some special cases) [16]. Further details and a discussion of other techniques are described in the blog post of Higham [15].

**Reduced Precision** Every time we store a number in lower precision, we lose accuracy, and every time we perform a floating point operation in lower precision, we lose accuracy. As mentioned in the introduction, in many instances in scientific computing, it is assumed that the model error, discretization error, etc., dominates the rounding error, and thus the rounding error is ignored in many analyses. When our rounding errors are on the order of  $10^{-4}$  instead of  $10^{-16}$ , this may no longer be a valid assumption. It is thus now more important than ever to carefully analyze these multiple sources of error *together*, and ultimately, to *balance* the errors. We will give three examples of this in Section 5.

## 4 A (Very) Brief History of Mixed Precision Iterative Refinement

Before jumping into our examples of balancing rounding error with other errors, we will give a (very) brief overview of some work on mixed precision iterative refinement. For a more complete history and further references, we direct the reader to Chapter 3 of the Ph.D. thesis of Vieublé [26].

Iterative refinement, shown in Algorithm 1, is an algorithm for iteratively improving the solution to linear systems  $Ax = b$ . The initial approximate solution in line 1 is typically obtained via an LU factorization, with the computed LU factors being reused for solving for the correction  $d_i$  in line 4.

**Algorithm 1** Iterative Refinement**Input:** nonsingular  $A \in \mathbb{R}^{n \times n}$ ,  $b \in \mathbb{R}^n$ , integer  $\text{maxit}$ **Output:** Approximate solution  $x_{i+1}$  to  $Ax = b$ 

- 
- 1: Solve  $Ax_0 = b$ .
  - 2: **for**  $i = 0:\text{maxit}$  **do**
  - 3:      $r_i = b - Ax_i$
  - 4:     Solve  $Ad_i = r_i$ .
  - 5:      $x_{i+1} = x_i + d_i$
  - 6: **end for**
- 

Wilkinson and his colleagues were using this approach as early as 1948 for solving linear systems on the Automatic Computing Engine [27], and Wilkinson wrote down a finite precision (fixed point) analysis of this algorithm in his 1963 book [28]. The traditional approach used by Wilkinson used two precisions: a working precision  $u$  for storage and all computations except the residual computation in line 3, which was done in precision  $u^2$  (double the working precision; here and in the remainder of the paper, we use “precision  $u$ ” as shorthand for “precision with unit roundoff  $u$ ”). This use of mixed precision was quite natural, as the hardware at the time computed an exact scalar product in precision  $u^2$  of two numbers in precision  $u$  at no extra cost. Wilkinson’s analysis showed that if  $3nu\kappa_\infty(A) < 1$ , where  $\kappa_\infty(A)$  is the infinity-norm condition number of  $A$ , then the limiting relative forward and backward errors, i.e.,  $\|x_i - x\|_\infty/\|x\|_\infty$  and  $\|b - Ax_i\|_\infty/(\|A\|_\infty\|x_i\|_\infty\|b\|_\infty)$ , are on the order  $O(u)$ .

In [5], it was shown that by changing the solver in line 4, more ill-conditioned systems can be solved. In particular, if we instead use preconditioned GMRES, where the preconditioner consists of the computed LU factors, the condition for convergence now contains a  $\kappa_\infty(U^{-1}L^{-1}A)$  instead of  $\kappa_\infty(A)$ , which can be substantially smaller, even if the LU factors are computed inexactly. This approach is called GMRES-based iterative refinement or simply GMRES-IR.

Motivated by emerging hardware that implemented half precision formats, [6] combined Wilkinson’s approach with low-precision factorization variants (see, e.g., [20]) in which the LU factorization (the most expensive part of the algorithm) is computed in half the working precision, and the rest of the computations are carried out in precision  $u$ . This resulted in bounds for a general *three-precision* variant of iterative refinement, which uses precision  $u_f$  for the factorization,  $u_r$  for the residual computation, and  $u$  for other computations, with  $u_r \leq u \leq u_f$ . The analysis also contains an “effective solve precision” which allows for a general solver in line 4.

Iterative refinement in general represents one case where it is natural to use mixed precision, that is, when the algorithm itself contains some type of “self-correction” mechanism, or some type of “inner-outer” solve scheme. This is a common occurrence in numerical algorithms.

## 5 Balancing Sources of Inexactness

Another case where it is natural to use low precision in parts of a computation is when there exist other significant sources of error in the algorithm, for example, the use of low-rank approximations, or coarse approximations of the domain. Again, while most existing analysis of these types of errors typically ignore finite precision errors, if we want to safely and effectively exploit low precision computations, we must look closely at what happens when we combine these types of analyses to determine how these errors interact. Our goal will ultimately be to balance these errors, and we can accomplish this by producing analyses which tell us how large the finite precision error can be relative to the other approximations we make so that it does not dominate.

Here we will present three examples of this balancing, dealing with sparsified matrices, randomized algorithms, and hierarchical matrix approximations. All numerical experiments were performed using MATLAB 2025a.

### 5.1 Example 1: Mixed Precision Sparse Approximate Inverse Preconditioners

Sparse approximate inverses are commonly-used algebraic preconditioners for Krylov subspace methods. For a thorough overview, see [2]. The goal is to construct a sparse matrix  $M$  that is an approximation of  $A^{-1}$ . In the classical approach, due to Grote and Huckle [13], shown in Algorithm 2, each column of  $M$  can be constructed independently. While this is in theory highly parallelizable, construction and memory requirements can still be costly for large scale problems; see, e.g., [14]. We thus aim to improve the performance and memory requirements by using low precision to construct this sparse approximate inverse if possible.

---

#### Algorithm 2 Sparse Approximate Inverse Construction

---

**Input:** nonsingular  $A \in \mathbb{R}^{n \times n}$ , initial sparsity structure  $\mathcal{J}$ , tolerance  $\tau$

**Output:** Sparse approximate inverse  $M \approx A^{-1}$

- 1: **for** each column  $k$  **do**
  - 2:     Compute QR factorization of submatrix of  $A$  defined by  $\mathcal{J}$ .
  - 3:     Use QR to solve  $\min_{m_k} \|e_k - Am_k\|_2$  where  $e_k$  is the  $k$ th column of the identity.
  - 4:     **if**  $\|r_k\|_2 = \|e_k - Am_k\|_2 \leq \tau$  **then break;**
  - 5:     **else** Add select nonzeros to  $\mathcal{J}$ , repeat from line 2.
  - 6: **end for**
- 

In [7], it is shown that we can execute the entire computation in some precision  $u_s$  and still reach a solution with  $\|\hat{r}_k\|_2 \leq \tau$  (where  $\hat{\cdot}$  denotes a computed quantity) if

$$n^3 u_s \| |e_k| + |A| |\hat{m}_k| \|_2 \leq \tau.$$

In other words, the problem must not be so ill-conditioned with respect to  $u_s$  that we incur an error greater than  $\tau$  just computing the residual if we want the stopping criterion to be satisfiable.

We can turn this into a looser but perhaps more useful a priori bound, and say that the sparse approximate inverse can be constructed in precision  $u_s$  as long as

$$u_s \text{cond}_2(A) \lesssim \tau,$$

where  $\text{cond}_2(A) = \|A^{-1}\|_2 \|A\|_2$ .<sup>1</sup> This tells us how our errors, coming from the stopping criterion in sparse approximate inverse construction, and the rounding error from the algorithm, must be *balanced*. For a given matrix  $A$  and a desired value of  $\tau$ , we have a limit on how large we can make  $u_s$ .

The results in [7] also show that this sparse approximate inverse can be used within GMRES-based iterative refinement (in place of the usual LU factorization) under certain constraints; here, we also have a constraint on how the working precision in iterative refinement (Algorithm 1) must be chosen relative to  $\tau$ .

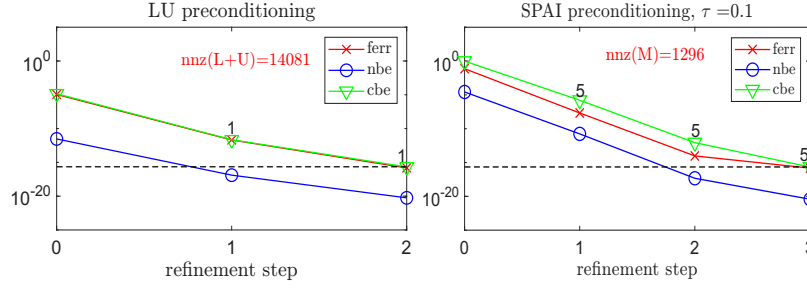
In Figure 1 we give one example of the use of low precision sparse approximate inverses within GMRES-based IR. Here we use the matrix `steam1` from SuiteSparse [10]. We compare GMRES-based IR with LU preconditioning (left) and sparse approximate inverse preconditioning (right). The preconditioners (both LU and the sparse approximate inverse) were computed in single precision, and for the refinement process we use  $u = u_r = \text{double}$ . In all cases, the GMRES convergence tolerance was set to  $10^{-6}$ . For SPAI, we set  $\tau = 0.1$ . The plots in Figure 1 show the forward error ('ferr'), normwise backward error ('nbe'), and componentwise backward error ('cbe'); the numbers above the markers indicate how many GMRES iterations were performed in each refinement step (i.e., in the solve in line 4 of Algorithm 1). In the LU and sparse approximate inverse plots, we list the size of the preconditioner in terms of number of nonzeros. We can see here that low-precision sparse approximate inverses are a viable alternative to LU-based preconditioners; while they require slightly more total GMRES iterations to converge, they result in a much less expensive preconditioner to apply and store. For further details and examples, see [7].

## 5.2 Example 2: Mixed Precision Randomized Nyström Approximation

Our next example involves computing a rank  $k$  approximation of an  $n \times n$  symmetric positive semidefinite matrix  $A$  using a randomized Nyström approach. The randomized Nyström approximation has the form  $A_N = (A\Omega)(\Omega^T A\Omega)^\dagger (A\Omega)^T$ , where  $\Omega$  is an  $n \times k$  sampling matrix, and  $\dagger$  indicates the Moore-Penrose pseudoinverse. The Nyström approximation arises in many applications, such as approximating kernel matrices and constructing spectral limited memory preconditioners. In some applications, the matrix is prohibitively large, and in some applications, its entries can

---

<sup>1</sup> Here we use  $\lesssim$  to indicate that this is more of a heuristic than a rigorous bound, since we have dropped dimensional constants (which are usually overestimates in practice anyway).



**Fig. 1** Comparison of GMRES-IR with LU preconditioning (left) and SPAI preconditioning (right) for the matrix `steam1` from SuiteSparse [10]. Numbers above markers give the number of GMRES iterations performed in each refinement step; in all cases the GMRES convergence tolerance was set to  $10^{-6}$ . For the top left plot, AMD ordering was applied before computing the LU factors. In all cases, we use GMRES-IR with  $u_f = \text{single}$ ,  $u = \text{double}$ ,  $u_r = \text{double}$ .

only be accessed once. This motivated the single-pass variant of the Nyström method [25], summarized in Algorithm 3.

---

### Algorithm 3 Single-Pass Nyström Approximation [25]

---

**Input:** Symmetric positive semidefinite  $A \in \mathbb{R}^{n \times n}$ , target rank  $k$

**Output:**  $U \in \mathbb{R}^{n \times k}$ , whose columns give approximate eigenvectors,  $\Theta \in \mathbb{R}^{k \times k}$ , whose diagonal entries give approximate eigenvalues

- 1:  $G = \text{randn}(n, k)$ ,  $[\Omega, \sim] = \text{qr}(G, 0)$
  - 2:  $Y = A\Omega$
  - 3: Compute shift  $\nu$  and compute  $Y_\nu = Y + \nu\Omega$ ,  $B = \Omega^T Y_\nu$
  - 4:  $C = \text{chol}(B + B^T)/2$ , Solve  $F = Y_\nu/C$
  - 5:  $[U, \Sigma, \sim] = \text{svd}(F, 0)$ ,  $\Theta = \max(0, \Sigma^2 - \nu I)$
- 

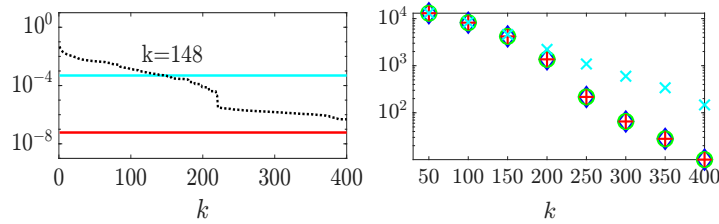
The matrix-matrix product in line 2 is overwhelmingly the dominant cost in this regime. Our goal will thus be to use a working precision  $u$  for all computations except this line, which is performed in precision  $u_p \geq u$ . Note that since this is the only time  $A$  is accessed,  $A$  should thus also be stored in precision  $u_p$ . Letting  $\hat{A}_N$  denote the Nyström approximation computed in mixed precision, we can write  $\|A - \hat{A}_N\|_2 \leq \|A - A_N\|_2 + \|A_N - \hat{A}_N\|_2$ ; in other words, we bound the overall error as the sum of the exact approximation error and the finite precision error. Bounds for the first term have been previously derived; see, e.g., [12], [11]. In [4], we prove the first bound on the finite precision error. The result essentially says that as long as we choose precisions  $u_p$  and  $u$  such that  $\kappa_2(A_k)\tilde{\kappa}(\Omega)^2 \ll u_p^{-2}$  and  $\kappa_2(A_k)\tilde{\kappa}(\Omega)^2 \ll u^{-1}$ , then

$$\|A_N - \hat{A}_N\|_F \lesssim \|A - A_N\|_F + k^{1/2} n u_p \kappa_2(A_k) \tilde{\kappa}(\Omega)^2 \|A\|_F,$$

where  $A_k$  is the best rank- $k$  approximation of  $A$  and  $\tilde{\kappa}(\Omega) = \|\Omega\|_F \|(W_1^T \Omega)^\dagger\|_2$ , where  $W_1$  are the eigenvectors for the leading  $k$  eigenvalues of  $A$ . Under some assumptions we can take this bound and derive a heuristic which says that it is likely that the finite precision error is less than the approximation error when

$$u_p \leq n^{-1/2} \lambda_k / \lambda_1,$$

where  $\lambda_1$  and  $\lambda_k$  are the largest and  $k$ th largest eigenvalues of  $A$ , respectively. This tells us again how the errors should be balanced: the greater the approximation error (the closer  $\lambda_k$  is to  $\lambda_1$ ), the lower the precision we can safely use. In Figure 2, we give one example of this for the matrix `bcsstm07` from SuiteSparse [10]. From the left plot, the heuristic above indicates that for  $k \leq 148$ , we can safely use half precision, and we can always safely use single precision. Indeed, this seems to be confirmed in the plot of the total error on the right; we see no difference with the exact case for single precision, and only start to see the finite precision error dominate for half precision around  $k > 150$ . For more examples, an alternative (potentially more rigorous) heuristic, and insight on choosing  $\Omega$  and  $\nu$ , see [4].



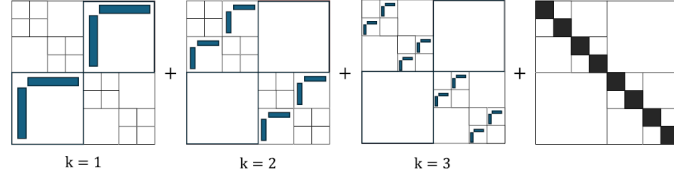
**Fig. 2** Experiments for matrix `bcsstm07`. The left plot shows the scaled spectrum  $\lambda_k/\lambda_1$  (black dotted), and the values of  $\sqrt{n}u_p$  for half precision (cyan) and single precision (red). The heuristic indicates that the finite precision error should not dominate for  $k$  less than the point where the black dotted line intersects the colored lines. The right plot shows the Frobenius norm of the mean total error,  $\|A - \hat{A}_N\|_F$ , taken over five runs, for exact arithmetic (blue diamonds),  $u_p, u = \text{double}$  (green circles),  $u_p = \text{single}$  and  $u = \text{double}$  (red +s), and  $u_p = \text{half}$  and  $u = \text{double}$  (cyan x).

### 5.3 Example 3: Mixed Precision Hierarchical Matrix Approximations

The final example involves hierarchical matrix approximations, and in particular, HODLR (hierarchical off-diagonal low-rank) matrices. HODLR matrices have a fixed hierarchical block structure; given a number of levels  $\ell$ , we store the off-diagonal blocks in each level as rank- $p$  approximations, and the diagonal blocks in the final level as dense matrices. See the depiction with  $\ell = 3$  in Figure 3. Such representations reduce the cost of computations, e.g.,  $O(n^2)$  computations become  $O(pn \log n)$ .

The idea is that since we are already approximating the off-diagonal blocks, we can likely store them in lower precision as well without losing a significant amount of additional information. This is analyzed in [3].

Let  $\varepsilon$  be the maximum relative error in the approximation of the off-diagonal blocks of a HODLR matrix  $H$  (coming, for example, from using a truncated SVD). Let  $\hat{H}$  be the mixed precision representation of  $H$  where the off-diagonal blocks in level  $k$  are stored in precision  $u_k$  and the diagonal blocks in the final level are stored in a working precision  $u < u_k$ . Then we can write a bound for the global



**Fig. 3** A depiction of a HODLR matrix with  $\ell = 3$  levels

approximation error as

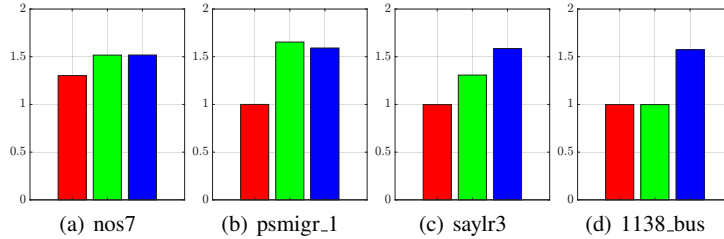
$$\frac{\|H - \hat{H}\|_F}{\|H\|_F} \lesssim 2\sqrt{2} \left( \sum_{k=1}^{\ell} 2^k \xi_k^2 u_k^2 \right)^{1/2} + \varepsilon,$$

where  $\xi_k = \max_{|i-j|=1} \|H_{ij}^{(k)}\|_F / \|H\|_F$  and  $H_{ij}^{(k)}$ ,  $|i-j| = 1$  denotes any off-diagonal block from the  $k$ th level.

This again tells us how to balance the errors: if we choose

$$u_k \leq \varepsilon / (2^{k/2} \xi_k),$$

then the right-hand side of the above bound simplifies to  $(2\sqrt{2}\ell + 1)\varepsilon$ ; in other words, we do not see the effect of storing the off-diagonal blocks in lower precision. In [3], an adaptive precision approach to HODLR matrix construction is developed based on this idea. We show an example of the resulting storage savings for various  $\varepsilon$  in Figure 4 for a few problems from SuiteSparse [10].



**Fig. 4** Storage savings of adaptive-precision HODLR matrices relative to uniform (double) precision HODLR matrices. The depth  $\ell = 8$ ; red bars correspond to  $\varepsilon = 10^{-7}$ , green bars indicate  $\varepsilon = 10^{-4}$ , and blue bars indicate  $\varepsilon = 10^{-1}$ .

In [3] we also determine how to set the working precision  $u$  in subsequent computations with mixed precision HODLR matrices so that the backward error of the computation does not exceed the error resulting from inexact representation of the matrix. The theorems say that for both matrix-vector products and LU factorizations with HODLR matrices stored in the mixed precision format, as long as  $u \leq \varepsilon/n$ ,

the backward error in these computations is dominated by the matrix representation error. Again, this gives us a theoretical guideline by which we can safely balance these errors. See [3] for details and examples.

## 6 Other Examples and Outlook

There are many other examples of balancing finite precision error with other errors, e.g., discretization and sampling errors [21], and hardware failures [8], that we did not have space to include here. For readers interested in mixed precision numerical algorithms, we refer to the surveys [1] and [17].

Looking forward, it is clear that we are moving more and more towards the era of extreme heterogeneity in large-scale systems, meaning that it is more important than ever to take a holistic approach to algorithm design, and truly consider the whole computational science process, including the hardware we run our codes on.

There is a great amount of work going on in developing new non-IEEE number formats and arithmetic systems, developing approximate hardware, and the emergence of new computing paradigms like quantum computing, neuromorphic computing and their hybrids. The task of analyzing errors from multiple sources and determining how to balance errors in an optimal way in these new paradigms will surely provide a source of interesting, challenging problems for a long time to come.

We end with a challenge to the reader to consider their own application area, the sources of error involved, and where it might make sense to use low or mixed precision computations.

**Acknowledgements** The author is supported by the European Union (ERC, inEXASCALE, 101075632). Views and opinions expressed are those of the authors only and do not necessarily reflect those of the European Union or the European Research Council. Neither the European Union nor the granting authority can be held responsible for them. The author additionally acknowledges support from the Charles University Research Centre program No. UNCE/24/SCI/005.

This manuscript is dedicated to the memory of Nicholas J. Higham, a mentor and friend, who first piqued my interest in mixed precision algorithms.

**Competing Interests** The author has no conflicts of interest to declare that are relevant to the content of this chapter.

## References

1. Abdelfattah, A., Anzt, H., Boman, E.G., Carson, E., Cojean, T., Dongarra, J., Fox, A., Gates, M., Higham, N.J., Li, X.S., et al.: A survey of numerical linear algebra methods utilizing mixed-precision arithmetic. *Int. J. High Perf. Comput.Appl.* **35**(4), 344–369 (2021)
2. Benzi, M.: Preconditioning techniques for large linear systems: a survey. *J. Comp. Phys.* **182**(2), 418–477 (2002)

3. Carson, E., Chen, X., Liu, X.: Mixed precision HODLR matrices. *SIAM J. Sci. Comput.* **47**(3), A1408–A1435 (2025)
4. Carson, E., Daužickaitė, I.: Single-pass nyström approximation in mixed precision. *SIAM J. Matrix Anal. Appl.* **45**(3), 1361–1391 (2024)
5. Carson, E., Higham, N.J.: A new analysis of iterative refinement and its application to accurate solution of ill-conditioned sparse linear systems. *SIAM J. Sci. Comput.* **39**(6), A2834–A2856 (2017)
6. Carson, E., Higham, N.J.: Accelerating the solution of linear systems by iterative refinement in three precisions. *SIAM J. Sci. Comput.* **40**(2), A817–A847 (2018)
7. Carson, E., Khan, N.: Mixed precision iterative refinement with sparse approximate inverse preconditioning. *SIAM J. Sci. Comput.* **45**(3), C131–C153 (2023)
8. Carson, E.C., Hercík, J.: The detection and correction of silent errors in pipelined Krylov subspace methods. *Numer. Algs.* pp. 1–36 (2025)
9. Chalmers, N., Kurzak, J., McDougall, D., Bauman, P.: Optimizing high-performance Linpack for exascale accelerated architectures. In: *Proc. Int. Conf. High Perf. Comput., Netw., Stor. Anal.*, pp. 1–12 (2023)
10. Davis, T.A., Hu, Y.: The University of Florida sparse matrix collection. *ACM Trans. Math. Soft.* **38**(1), 1–25 (2011)
11. Frangella, Z., Tropp, J.A., Udell, M.: Randomized Nyström preconditioning. *SIAM J. Matrix Anal. Appl.* **44**(2), 718–752 (2023)
12. Gittens, A., Mahoney, M.W.: Revisiting the Nyström method for improved large-scale machine learning. *J. Mach. Learn. Res.* **17**(1), 3977–4041 (2016)
13. Grote, M.J., Huckle, T.: Parallel preconditioning with sparse approximate inverses. *SIAM J. Sci. Comput.* **18**(3), 838–853 (1997)
14. He, G., Yin, R., Gao, J.: An efficient sparse approximate inverse preconditioning algorithm on GPU. *Conc. and Comp.: Prac. and Exp.* **32**(7), e5598 (2020)
15. Higham, N.J.: Can we solve linear algebra problems at extreme scale and low precisions? <https://nhigham.com/2021/09/14/can-we-solve-linear-algebra-problems-at-extreme-scale-and-low-precisions/> (2021)
16. Higham, N.J., Mary, T.: A new approach to probabilistic rounding error analysis. *SIAM J. Sci. Comput.* **41**(5), A2815–A2835 (2019)
17. Higham, N.J., Mary, T.: Mixed precision algorithms in numerical linear algebra. *Acta Numer.* **31**, 347–414 (2022)
18. Higham, N.J., Pranesh, S., Zounon, M.: Squeezing a matrix into half precision, with an application to solving linear systems. *SIAM J. Sci. Comput.* **41**(4), A2536–A2551 (2019)
19. HPL-MxP mixed-precision benchmark. <https://hpl-mxp.org> (2024)
20. Langou, J., Langou, J., Luszczek, P., Kurzak, J., Buttari, A., Dongarra, J.: Exploiting the performance of 32 bit floating point arithmetic in obtaining 64 bit accuracy (revisiting iterative refinement for linear systems). In: *Proc. 2006 ACM/IEEE Conf. Supercomput.* (2006). DOI 10.1109/SC.2006.30. URL <https://doi.org/10.1109/SC.2006.30>
21. Martínek, J., Carson, E., Scheichl, R.: Exploiting inexact computations in multilevel sampling methods. *arXiv preprint arXiv:2503.05533* (2025)
22. NVIDIA: NVIDIA H100 Tensor Core GPU Architecture. <https://resources.nvidia.com/en-us-hopper-architecture/nvidia-h100-tensor-c> (2022)
23. Scott, J., Tůma, M.: Developing robust incomplete Cholesky factorizations in half precision arithmetic. *Numer. Algs.* pp. 1–22 (2025)
24. TOP500: The List. <https://top500.org/> (1993). Accessed January 2026
25. Tropp, J.A., Yurtsever, A., Udell, M., Cevher, V.: Fixed-rank approximation of a positive-semidefinite matrix from streaming data. *Adv. Neural Info. Proc. Sys.* **30** (2017)
26. Vieublé, B.: Mixed precision iterative refinement for the solution of large sparse linear systems. Phd thesis, University of Toulouse, Toulouse, France (2022)
27. Wilkinson, J.H.: Progress report on the automatic computing engine. Tech. Rep. MA/17/1024, Mathematics Division, Department of Scientific and Industrial Research, National Physical Laboratory, Teddington, UK (1948). 127 pp.
28. Wilkinson, J.H.: Rounding Errors in Algebraic Processes, *Notes on Applied Science*, vol. 32. Her Majesty’s Stationery Office, London (1963)