



**HAL**  
open science

## Five-Precision GMRES-based iterative refinement

Patrick Amestoy, Alfredo Buttari, Nicholas Higham, Jean-Yves L'Excellent,  
Théo Mary, Bastien Vieublé

► **To cite this version:**

Patrick Amestoy, Alfredo Buttari, Nicholas Higham, Jean-Yves L'Excellent, Théo Mary, et al.. Five-Precision GMRES-based iterative refinement. *SIAM Journal on Matrix Analysis and Applications*, 2024, 45 (1), pp.529–552. 10.1137/23M1549079 . hal-03190686v1

**HAL Id: hal-03190686**

**<https://hal.science/hal-03190686v1>**

Submitted on 6 Apr 2021 (v1), last revised 1 Apr 2024 (v2)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# FIVE-PRECISION GMRES-BASED ITERATIVE REFINEMENT \*

PATRICK AMESTOY<sup>†</sup>, ALFREDO BUTTARI<sup>‡</sup>, NICHOLAS J. HIGHAM<sup>§</sup>,  
JEAN-YVES L'EXCELLENT<sup>†</sup>, THEO MARY<sup>¶</sup>, AND BASTIEN VIEUBLÉ<sup>||</sup>

**Abstract.** GMRES-based iterative refinement in three precisions (GMRES-IR3) uses a low precision LU factorization to accelerate the solution of a linear system without compromising numerical stability or robustness. GMRES-IR3 solves the update equation using GMRES preconditioned by the LU factors, where all operations within GMRES are carried out in the working precision  $u$ , except for the matrix-vector products and the application of the preconditioner, which require the use of extra precision  $u^2$ . The use of extra precision can be expensive, and is especially unattractive if it is not available in hardware; for this reason, existing implementations have not used extra precision, despite the absence of an error analysis for this approach. We relax the requirements on the precisions used within GMRES, allowing the use of arbitrary precisions  $u_p$  (for applying the preconditioner) and  $u_g$  (for the rest of the operations). We obtain the five-precision GMRES-based iterative refinement (GMRES-IR5) algorithm. We carry out a rounding error analysis that generalizes that of GMRES-IR3, obtaining conditions under which the forward and backward errors converge to their limiting values. Our analysis makes use of a new result on the backward stability of MGS-GMRES in two precisions. On hardware where up to five arithmetics are available, the number of possible combinations of precisions in GMRES-IR5 is extremely large, but our analysis identifies a small subset of relevant combinations. By choosing from within this subset one can achieve different levels of tradeoff between cost and robustness, which allows for a finer choice of precisions depending on the problem difficulty and the available hardware. Our numerical experiments on both random dense matrices and real-life sparse matrices from a wide range of applications show that the practical behavior of GMRES-IR5 is in good agreement with our theoretical analysis. GMRES-IR5 therefore has the potential to solve relatively badly conditioned problems in less time and memory than GMRES-IR3, thanks to the use of lower precision arithmetic in the GMRES iterations.

**Key words.** iterative refinement, GMRES, linear system, mixed precision, multiple precision, rounding error analysis, floating-point arithmetic, backward error, forward error, preconditioning

**AMS subject classifications.** 65G50, 65F05, 65F08, 65F50, 65F10

**1. Introduction.** Modern hardware increasingly supports low precision floating-point arithmetics that provide unprecedented speed, communication, and energy benefits. This has generated renewed interest in mixed precision algorithms that combine these lower precision arithmetics with higher precision ones to achieve both high performance and high accuracy. Among such mixed precision algorithms, iterative refinement for the solution of a linear system  $Ax = b$  is one of the oldest and most successful [1], [12, Chap. 12], [20], [24].

We recall, in Algorithm 1.1, iterative refinement in the general form recently proposed and analyzed by Carson and Higham [3]. Algorithm 1.1 includes and generalizes previous iterative refinement algorithms in two ways. First, each step of the algorithm may be carried out in a different precision: the initial LU factorization of  $A$  is performed in precision  $u_f$ , the residual  $r_i = b - Ax_i$  is computed in precision  $u_r$ , and the solution  $x_i$  is updated in the target, working precision  $u$ . Second, the correction term  $d_i$  is obtained by solving the system  $Ad_i = r_i$  with a generic solver

---

\*Version of April 3, 2021.

<sup>†</sup>Mumps Technologies, ENS Lyon, 46 Allée d'Italie, F-69007 Lyon, France (patrick.amestoy@mumps-tech.com; jean-yves.l.excellent@mumps-tech.com)

<sup>‡</sup>CNRS, IRIT, 2 Rue Charles Camichel, F-31071 Toulouse, France (alfredo.buttari@irit.fr)

<sup>§</sup>Department of Mathematics, The University of Manchester, Manchester, M13 9PL, UK (nick.higham@manchester.ac.uk)

<sup>¶</sup>Sorbonne Université, CNRS, LIP6, Paris, F-75005, France (theo.mary@lip6.fr)

<sup>||</sup>INPT, IRIT, 2 Rue Charles Camichel, F-31071 Toulouse, France (bastien.vieuble@irit.fr)

---

**Algorithm 1.1** Iterative refinement

---

**Input:** an  $n \times n$  matrix  $A$  and a right-hand side  $b$ .

**Output:** an approximate solution to  $Ax = b$ .

- 1: Compute the LU factorization  $A = LU$ .  $u_f$
  - 2: Initialize  $x_0$  (to, e.g.,  $U^{-1}L^{-1}b$ ).  $u_f$
  - 3: **while not** converged **do**
  - 4:   Compute  $r_i = b - Ax_i$ .  $u_r$
  - 5:   Solve  $Ad_i = r_i$ .  $u_s$
  - 6:   Compute  $x_{i+1} = x_i + d_i$ .  $u$
  - 7: **end while**
- 

assumed to compute  $d_i$  with relative accuracy of order  $u_s \geq u$ . The precision  $u_s$  is a measure of the quality of the solution the solver produces rather than an input to the solver.

The traditional choice of solver, which we call the LU solver, is to compute  $d_i = U^{-1}L^{-1}r_i$  by substitution in precision  $u_f$ , so that  $d_i$  is computed at accuracy  $u_s = O(u_f)$ . In this case, Algorithm 1.1 employs up to three precisions. We denote this variant by LU-IR3. While LU-IR3 can be very attractive for well-conditioned matrices, it is only guaranteed to converge when  $\kappa(A)u_f \ll 1$ , where  $\kappa(A) = \|A\|\|A^{-1}\|$  denotes the condition number of  $A$ . This condition can be quite restrictive, especially when low precision arithmetic is used for the factorization: for example, the condition becomes  $\kappa(A) \ll 2 \times 10^3$  with IEEE fp16 (half) precision, and  $\kappa(A) \ll 3 \times 10^2$  with bfloat16.

Carson and Higham [2] take  $u_f = u$  and solve  $Ad_i = r_i$  by GMRES preconditioned by the computed LU factors of  $A$ , with the preconditioner applied in higher precision  $u^2$  and the residual also computed at higher precision,  $u_r = u^2$ . Their motivation was to compute solutions with forward error of order  $u$  even when  $\kappa(A)$  is of order  $u^{-1}$ . Carson and Higham [3] subsequently generalized this algorithm in order to be able to exploit lower precision arithmetic for the factorization. This algorithm, called GMRES-based iterative refinement in three precisions (GMRES-IR3), is recalled in Algorithm 1.2. On line 5,  $\tilde{A}$  is not explicitly formed, but its action on a vector is obtained with a matrix multiplication followed by two triangular solves.

---

**Algorithm 1.2** GMRES based iterative refinement in three precisions (GMRES-IR3)

---

**Input:** an  $n \times n$  matrix  $A$  and a right-hand side  $b$ .

**Output:** Approximate solution  $\hat{x}$  to  $Ax = b$

- 1: Compute the LU factorization  $A = LU$ .  $u_f$
  - 2: Initialize  $x_0$  (e.g., to  $U^{-1}L^{-1}b$ ).  $u_f$
  - 3: **while not** converged **do**
  - 4:   Compute  $r_i = b - Ax_i$ .  $u_r$
  - 5:   Solve  $U^{-1}L^{-1}Ad_i = U^{-1}L^{-1}r_i$  by GMRES at precision  $u$  with matrix–vector products with  $\tilde{A} = U^{-1}L^{-1}A$  computed at precision  $u^2$ .
  - 6:   Compute  $x_{i+1} = x_i + d_i$ .  $u$
  - 7: **end while**
- 

GMRES-IR3 is shown in [3] to be guaranteed to converge as long as  $\kappa(A)^2u_f^2u \ll 1$ , and it is therefore able to handle much more ill-conditioned matrices than LU-IR3. However, the requirement that the preconditioner be applied in precision  $u^2$  is

a practical limitation, because this can be expensive; it is particularly inconvenient if the target accuracy  $u$  is double precision, as it requires applying the preconditioner in quadruple precision, an arithmetic that is not natively supported on most modern hardware and is an order of magnitude slower than double precision when implemented in software [13]. In fact, practical implementations of GMRES-IR3, as developed in [8], [9], [10] and implemented in the MAGMA library [19] and the NVIDIA cuSOLVER library [6], have relaxed this requirement by applying the preconditioner in double rather than quadruple precision—even though the error analysis of [3] does not cover this case. GMRES-IR variants for symmetric positive definite systems and least squares problems have also allowed just two precisions to be used [4], [17]. The question arises of whether we can use a lower precision to apply the preconditioner within GMRES and still obtain a GMRES-IR solver able to handle more ill-conditioned matrices than LU-IR.

These practical constraints and theoretical questions lead us to propose new variants of GMRES-based iterative refinement with relaxed requirements on the precisions used within the GMRES solver. We allow the preconditioner (the LU factors) to be applied in an arbitrary precision  $u_p$ , with  $u_p \geq u^2$ . We also allow the rest of the GMRES computations to be performed in an arbitrary precision  $u_g$ , with  $u_g \geq u$ . We obtain Algorithm 3.1, which has up to *five* different precisions in play and which we thus call GMRES-IR5. We generalize the rounding error analysis of Carson and Higham [2], [3] to this new GMRES-IR5 variant and obtain several novel combinations of precision parameters that achieve different levels of compromise between cost and robustness. GMRES-IR5 therefore fills the gap between LU-IR3 and GMRES-IR3, allowing for new variants that can be much faster than GMRES-IR3 but that remain more robust than LU-IR3. When five different arithmetics are available (such as bfloat16, fp16, fp32, fp64, and fp128), the number of possible combinations is extremely large, but our analysis identifies the small subset of these combinations that is relevant. These include combinations where  $u_p = u$  and, therefore, our results provide convergence guarantees for the existing implementations mentioned above. Moreover, some relevant combinations result in new variants of iterative refinement that employ, for the first time, four or even five different precisions.

The rest of this article is organized as follows. We begin, in section 2, by providing technical background on LU-IR3 and GMRES-IR3. In section 3, we present the generalized error analysis of GMRES-IR5 in five precisions and its different combinations. In section 4, we support our analysis with numerical experiments on both random dense matrices and on a wide range of real-life matrices from the SuiteSparse collection [7]. We provide our concluding remarks in section 5.

**2. Preliminaries on LU-IR3 and GMRES-IR3.** In this section we summarize our notation and assumptions for the analysis, and briefly recall the essential technical background on the LU-IR3 and GMRES-IR3 algorithms of Carson and Higham [2, 3].

We use the standard model of floating-point arithmetic [12, sect. 2.2] and we use the notation  $\text{fl}(\cdot)$  to denote the computed value of a given expression. For any integer  $k$  we define

$$\gamma_k = \frac{ku}{1 - ku}.$$

A superscript on  $\gamma$  denotes that  $u$  carries that superscript as a subscript; thus  $\gamma_k^f = ku_f/(1 - ku_f)$ , for example. We also use the notation  $\tilde{\gamma}_k = \gamma_{ck}$  to hide modest constants  $c$ .

The error bounds obtained by our analysis depend on the problem dimension  $n$ , the number of iterations  $k$  of GMRES, and the growth factor  $\rho_n$  of the LU factorization. In parts of the analysis, we will gather these constants into a generic function  $f(n, k, \rho_n)$ . The constants depending on  $n$  are known to be pessimistic [5], [15], [16] and, with standard pivoting strategies,  $\rho_n$  is almost always small in practice and of order a constant [12, chap. 9] (for exceptions, see [11]). Therefore, for the sake of the readability, we do not always keep track of the precise value of  $f(n, k, \rho_n)$ . When we drop constants  $f(n, k, \rho_n)$  from an inequality we write the inequality using “ $\ll$ ”. A convergence condition expressed as “ $\kappa(A) \ll \theta$ ” can be read as “ $\kappa(A)$  is sufficiently less than  $\theta$ ”. Finally, we also use the notation  $\lesssim$  when dropping second order terms in the error bounds.

We consider linear systems  $Ax = b$ , where  $A \in \mathbb{R}^{n \times n}$  is nonsingular and  $b \in \mathbb{R}^n$ . The forward error of an approximation  $\hat{x}$  is  $\|x - \hat{x}\|/\|x\|$ , while the (normwise) backward error of  $\hat{x}$  is [12, sec. 7.1]

$$\min\{\epsilon : (A + \Delta A)\hat{x} = b + \Delta b, \|\Delta A\| \leq \epsilon\|A\|, \|\Delta b\| \leq \epsilon\|b\|\} = \frac{\|b - A\hat{x}\|}{\|A\| \|\hat{x}\| + \|b\|}.$$

We define the componentwise condition numbers

$$\text{cond}(A, x) = \frac{\| |A^{-1}| |A| |x| \|}{\|x\|}, \quad \text{cond}(A) = \| |A^{-1}| |A| \|,$$

where  $|A| = (|a_{ij}|)$ .

Our error analysis uses both the  $\infty$ -norm and the Frobenius norm, denoted by  $\|\cdot\|_\infty$  and  $\|\cdot\|_F$ , respectively, and we write  $\kappa_\infty(A)$  and  $\kappa_F(A)$  for the corresponding condition numbers of  $A$ . We will use unsubscripted norms or condition numbers when the constants depending on the problem dimensions have been dropped, since the norms are equivalent.

We denote by  $q$  the maximum number of nonzeros in any row of  $[Ab]$ ; thus  $q = n + 1$  for a dense matrix  $A$  and a vector  $b$ . As in [3], however, we do not make any assumptions on the sparsity of the LU factors  $L$  and  $U$ .

Carson and Higham perform the error analysis of Algorithm 1.1 with a general solver used at step 5. The solver is assumed to satisfy the conditions

$$\hat{d}_i = (I + u_s E_i) d_i, \quad u_s \|E_i\|_\infty < 1, \quad (2.1a)$$

$$\|\hat{r}_i - A\hat{d}_i\|_\infty \leq u_s (c_1 \|A\|_\infty \|\hat{d}_i\|_\infty + c_2 \|\hat{r}_i\|_\infty), \quad (2.1b)$$

where  $\hat{d}_i$  and  $\hat{r}_i$  denote the computed  $d_i$  and  $r_i$ , and where  $E_i$ ,  $c_1$ , and  $c_2$ , are functions of  $n$ ,  $A$ ,  $\hat{r}_i$ , and  $u_s$ . These conditions are needed for the normwise forward and backward error analyses, respectively.

In analyzing iterative refinement we aim to show that the forward error and backward error decrease until they reach a certain size called the limiting forward error or backward error. We will informally refer to attainment of this level as “convergence”, while recognizing that the error does not necessarily converge in the formal sense.

Corollary 3.3 of [3] shows that as long as

$$u_s \|E_i\|_\infty < 1 \quad (2.2)$$

for each  $i$ , the forward error converges to a limiting value of  $qu_r \text{cond}(A, x) + u$ . Here, we have simplified the convergence condition by ignoring a term that is shown in [3]

TABLE 2.1

Convergence conditions and limiting backward and forward errors for the LU-IR3 and GMRES-IR3 algorithms analyzed in [2], [3].

	Convergence condition		Limiting value
	LU-IR3	GMRES-IR3	
Backward error	$\kappa(A) \ll u_f^{-1}$	$\kappa(A) \ll u^{-1/2}u_f^{-1/2}$	$qu_r + u$
Forward error	$\kappa(A) \ll u_f^{-1}$	$\kappa(A) \ll u^{-1/2}u_f^{-1}$	$qu_r \text{cond}(A, x) + u$

to be dominated in practice by the left-hand side of (2.2). Corollary 4.2 of [3] shows that if

$$u_s(c_1\kappa_\infty(A) + c_2) \ll 1, \quad (2.3)$$

then the backward error converges to a limiting value of  $qu_r + u$ . The convergence conditions of Algorithm 1.1 therefore depend only on  $u_s$ , whereas the limiting accuracy and backward errors depend on  $u$  and  $u_r$ .

In the case of LU-IR3, we have  $u_s \equiv u_f$ ,  $\|E_i\|_\infty \equiv f(n, \rho_n)\kappa_\infty(A)$ , and  $c_1 = c_2 = f(n, \rho_n)$ , so the convergence condition is

$$\kappa(A) \ll u_f^{-1}, \quad (2.4)$$

for both the forward and backward errors.

In the case of GMRES-IR3, that is, Algorithm 1.2, we have instead  $u_s \equiv u$  and  $\|E_i\|_\infty \equiv f(n, k, \rho_n)\kappa_\infty(\tilde{A})$  [2, sec. 3], where  $\tilde{A} = \hat{U}^{-1}\hat{L}^{-1}A$  is the matrix preconditioned by the computed LU factors  $\hat{L}$  and  $\hat{U}$ . Using the bound

$$\kappa_\infty(\tilde{A}) \ll \kappa_\infty(A)^2 u_f^2 \quad (2.5)$$

from [3, Eq. (8.3)], which shows that  $\tilde{A}$  is usually better conditioned than  $A$ , (2.2) yields the condition

$$\kappa(A) \ll u^{-1/2}u_f^{-1}, \quad (2.6)$$

for the forward error to converge. This is a significantly less restrictive condition than (2.4) for LU-IR3. Similarly, for GMRES-IR3 we can show that  $c_1 = f(n, k, \rho_n)\|\tilde{A}\|_\infty$  and  $c_2 = f(n, k, \rho_n)\kappa_\infty(A)$ , so that after bounding  $\|\tilde{A}\|_\infty$  with [2, Eq. (3.2)] the convergence condition (2.3) for the backward error becomes

$$\kappa(A) \ll u^{-1/2}u_f^{-1/2}, \quad (2.7)$$

which is also less restrictive than (2.4) when  $u < u_f$ .

We summarize in Table 2.1 convergence conditions and limiting backward and forward errors for LU-IR3 and GMRES-IR3. Since the limiting backward error and forward error are solver independent (as long as the conditions (2.1) are satisfied), they are the same for LU-IR3 and GMRES-IR3.

**3. Error analysis of GMRES-IR5.** In this section we generalize the error analysis of GMRES-IR3 (Algorithm 1.2) [2], [3] to GMRES-IR5, which is defined in Algorithm 3.1. First, we extend in section 3.1 the analysis of Paige, Rozložník, and Strakoš [22] on the backward stability of MGS-GMRES to arbitrary matrix–vector products satisfying a generic error bound. Second, we use this generalized analysis in section 3.2 to show that GMRES-IR5 provides a backward stable solution for the problem  $\tilde{A}d_i = s_i$ , where  $\tilde{A} = \hat{U}^{-1}\hat{L}^{-1}A$  and  $s_i = \hat{U}^{-1}\hat{L}^{-1}r_i$ . We then use this new

analysis in section 3.3 to obtain conditions on  $\kappa(A)$  for GMRES-IR5 to converge (note that the limiting backward error and forward error are left unchanged because they only depend on  $u$  and  $u_r$ ). Finally, we use these generalized convergence conditions in section 3.4 to identify which combinations of precisions are meaningful.

---

**Algorithm 3.1** GMRES based iterative refinement in five precisions (GMRES-IR5)

---

**Input:** an  $n \times n$  matrix  $A$  and a right-hand side  $b$ .

**Output:** Approximate solution  $\hat{x}$  to  $Ax = b$

- 1: Compute the LU factorization  $A = LU$ .  $u_f$
  - 2: Initialize  $x_0$  (to, e.g.,  $U^{-1}L^{-1}b$ ).  $u_f$
  - 3: **while not converged do**
  - 4:   Compute  $r_i = b - Ax_i$ .  $u_r$
  - 5:   Solve  $U^{-1}L^{-1}Ad_i = U^{-1}L^{-1}r_i$  by GMRES at precision  $u_g$  with matrix–vector products with  $\tilde{A} = U^{-1}L^{-1}A$  computed at precision  $u_p$ .  $u_p$
  - 6:   Compute  $x_{i+1} = x_i + d_i$ .  $u$
  - 7: **end while**
- 

Our analysis makes use of the following two assumptions on the precisions:

- $u_g \geq u$ : since the solution computed by GMRES is stored in the working precision  $u$ , we do not expect running GMRES in precision  $u_g < u$  to give a significant benefit.
- $\kappa(A) \max(u, u_p) < 1$ : the assumption  $\kappa(A)u < 1$  is already present in the three-precision analysis [3]; here, we also require  $\kappa(A)u_p < 1$ .

**3.1. Error analysis of MGS-GMRES with arbitrary matrix–vector products.** We assume that the MGS-GMRES variant of GMRES is used in GMRES-IR5, so we need to bound the backward error of two-precision MGS-GMRES for the solution of the preconditioned system  $\tilde{A}d_i = s_i$ . The analysis of Paige et al. [22] is for fixed precision, unpreconditioned MGS-GMRES, and therefore is not directly applicable. Note that in the analysis of [2] the preconditioner is applied in precision  $u^2$ , and using the assumption  $\kappa(A)u < 1$  it is shown that the products with  $\tilde{A}$  in precision  $u^2$  are at least as accurate as the products with  $A$  in precision  $u$ , so that the backward stability result of [22] still holds. In our case, the same argument does not apply, and so we must generalize the backward stability result [22, Eq. (8.15)] to the case of arbitrary matrix–vector products satisfying a generic error bound. We state the conclusion of our analysis in the next theorem.

**THEOREM 3.1.** *Consider the solution of a linear system*

$$Bx = d, \quad B \in \mathbb{R}^{n \times n}, \quad 0 \neq d \in \mathbb{R}^n \quad (3.1)$$

*with an MGS-GMRES solver carrying out its operations in precision  $u_g$ , except for the products with  $B$ , which satisfy instead*

$$\text{fl}(Bv) = Bv + f, \quad \|f\|_2 \lesssim \epsilon_p \|B\|_F \|v\|_2, \quad (3.2)$$

*where  $\epsilon_p > 0$  is a parameter quantifying the stability of the matrix–vector products. Provided that*

$$\sigma_{\min}(B) \gtrsim (k^{1/2} \epsilon_p + \tilde{\gamma}_{kn}^g) \|B\|_F, \quad (3.3)$$

there is a step  $k \leq n$  such that the algorithm produces a computed  $\hat{x}_k$  satisfying

$$(B + \Delta B)\hat{x}_k = d + \Delta d, \quad (3.4a)$$

$$\|\Delta B\|_F \lesssim (k^{1/2}\epsilon_p + \tilde{\gamma}_{kn}^g)\|B\|_F, \quad (3.4b)$$

$$\|\Delta d\|_2 \lesssim \tilde{\gamma}_{kn}^g\|d\|_2. \quad (3.4c)$$

*Proof.* The proof of Theorem 3.1 relies on the analysis of [22] and, more precisely, on [22, sec. 8] and [22, Eq. (4.3)] therein. For the sake of readability, we will not attempt to make this proof self-contained but, rather, we will highlight the differences with the analysis in [22] and refer the reader to that work for the full details. The original notation has been slightly adapted to be consistent with the notation of this article inherited from [2] and [3].

In our version of MGS-GMRES we consider a product with  $B$  satisfying (3.2). We now show that considering (3.2) with  $\epsilon_p \neq \gamma_n^g$  mainly changes [22, Eq. (4.3)]. Let us consider  $\hat{V}_k = [\hat{v}_1, \dots, \hat{v}_k]$ , the matrix of computed basis vectors, and  $\dot{V}_k = [\dot{v}_1, \dots, \dot{v}_k]$  the same matrix but with its columns correctly normalized; that is, for  $j \leq k$ ,

$$\begin{aligned} \hat{v}_j &= \dot{v}_j + \Delta v_j^{(1)}, & \|\Delta v_j^{(1)}\|_2 &\leq \tilde{\gamma}_n^g, \\ \hat{V}_k &= \dot{V}_k + \Delta V_k^{(1)}, & \Delta V_k^{(1)} &= [\Delta v_1^{(1)}, \dots, \Delta v_k^{(1)}], \end{aligned} \quad (3.5)$$

where  $\Delta v_j^{(1)}$  is the error for the normalization of  $\hat{v}_j$  and  $\Delta V_k^{(1)}$  is the accumulated error for the normalization of the basis at step  $k$ . By (3.2) and (3.5), we obtain

$$\begin{aligned} \text{fl}(B\hat{v}_j) &= B(\dot{v}_j + \Delta v_j^{(1)}) + f_j \\ &= B\dot{v}_j + \Delta v_j^{(2)}, \end{aligned}$$

where  $\Delta v_j^{(2)} = B\Delta v_j^{(1)} + f_j$  satisfies  $\|\Delta v_j^{(2)}\|_2 \lesssim (\epsilon_p + \tilde{\gamma}_n^g)\|B\|_F$  since  $\|\dot{v}_j\|_2 = 1$  and  $\|f_j\|_2 \lesssim \epsilon_p\|B\|_F\|\dot{v}_j + \Delta v_j^{(1)}\|_2$ . We therefore obtain

$$\text{fl}(B\hat{V}_k) \approx B\dot{V}_k + \Delta V_k^{(2)}, \quad \|\Delta V_k^{(2)}\|_F \lesssim k^{1/2}(\epsilon_p + \tilde{\gamma}_n^g)\|B\|_F, \quad (3.6)$$

where  $\Delta V_k^{(2)}$  contains the error for both the product and the normalization at the  $k$ th iteration. Equation (3.6) is our new version of [22, Eq. (4.3)]; adapting the remainder of [22, sec. 8] to take this change into account is straightforward. Consequently, we show that at the  $(\bar{m} - 1)$ st iteration, MGS-GMRES has computed a backward stable solution of the system, where  $\bar{m}$  satisfies [22, Eq. (6.1)]. From now on, we set  $k$  such that  $k \equiv \bar{m} - 1 \leq n$ , and rewrite [22, Eq. (8.2)] as

$$\begin{aligned} r_k(\hat{y}_k) &\equiv d_k - B_k\hat{y}_k, & d_k &\equiv d + \Delta d_k(\hat{y}_k), & B_k &\equiv B\dot{V}_k + \Delta V_k^{(3)}(\hat{y}_k), \\ \|\Delta d_k(\hat{y}_k)\|_2 &\leq \tilde{\gamma}_{kn}^g\|d\|_2, & \Delta V_k^{(3)}(y) &\equiv \Delta V_k^{(2)} + \Delta C_k(y), \\ \|\Delta V_k^{(3)}\|_F &\lesssim (k^{1/2}\epsilon_p + \tilde{\gamma}_{kn}^g)\|B\|_F, \end{aligned} \quad (3.7)$$

where  $\Delta C_k(y)$  and  $\Delta d_k(y)$  are the errors in the MGS least squares solution [22, sec. 7]. Using the scaling invariance of MGS to scale the right-hand side  $d_k$  by  $\phi'$  and making use of [22, Thm. 2.4] gives a bound on the residual [22, Eq. (8.9)]

$$\|r_k(\hat{y}_k)\|_2^2 \leq (\tilde{\gamma}_{kn}^g)^2(\|d_k\phi'\|_2^2 + \|B_k\|_F^2)2(\phi')^{-2}. \quad (3.8)$$



In addition to bounding  $\|B_k\|_F^2$  and  $\|d_k\phi'\|_2^2$ , we use the nonsingularity condition (3.3) in the same fashion as [22, Eq. (8.11)] to compute a bound for  $(\phi')^{-2}$ , which allows us to rewrite [22, Eq. (8.12)] as

$$\|r_k(\hat{y}_k)\|_2 \lesssim \tilde{\gamma}_{kn}^g (\|B\|_F \|\hat{y}_k\|_2 + \|d\|_2). \quad (3.9)$$

Since  $\epsilon_p$  appears in second order terms and higher, they have been dropped, making (3.9) equivalent to [22, Eq. (8.12)]. Considering now  $\hat{x}_k = \text{fl}(\hat{V}_k \hat{y}_k) = (\hat{V}_k + \Delta V_k^{(4)}) \hat{y}_k$  and using a standard matrix–vector product in precision  $u_g$  satisfying  $\|\Delta V_k^{(4)}\|_F \leq \tilde{\gamma}_k^g \|\hat{V}_k\|_F$  and  $\Delta B_k \equiv [\Delta V_k^{(3)}(\hat{y}_k) - B(\Delta V_k^{(4)} + \hat{V}_k - \dot{V}_k)] \hat{y}_k \frac{\hat{x}_k^T}{\|\hat{x}_k\|_2^2}$ , we can rewrite [22, Eq. (8.15)] as

$$\begin{aligned} r_k(\hat{y}_k) &= d + \Delta d_k(\hat{y}_k) - (B + \Delta B_k) \hat{x}_k, \\ \|r_k(\hat{y}_k)\|_2 &\lesssim \tilde{\gamma}_{kn}^g (\|B\|_F \|\hat{x}_k\|_2 + \|d\|_2), \\ \|\Delta d_k(\hat{y}_k)\|_2 &\leq \tilde{\gamma}_{kn}^g \|d\|_2, \\ \|\Delta B_k\|_F &\lesssim (k^{1/2} \epsilon_p + \tilde{\gamma}_{kn}^g) \|B\|_F. \end{aligned} \quad (3.10)$$

This leads to (3.4) and completes the proof of the theorem.  $\square$

In the original analysis,  $\text{fl}(Bv)$  is a standard matrix–vector product operation and  $u_p = u_g$ , so  $\text{fl}(Bv) = (B + \Delta B)v$  holds where  $\|\Delta B\|_F \leq \gamma_q^g \|B\|_F$  [12, sec. 3.5]. In this case we can set  $f = \Delta Bv$  and apply Theorem 3.1 with  $\epsilon_p = \gamma_q^g$ , recovering the result of [22] for an unpreconditioned MGS-GMRES in uniform precision, which produces a solution  $\hat{x}_k$  of the system  $Bx = d$  satisfying

$$(B + \Delta B) \hat{x}_k = d + \Delta d, \quad \|\Delta B\|_F \lesssim \tilde{\gamma}_{kn}^g \|B\|_F, \quad \|\Delta d\|_2 \lesssim \tilde{\gamma}_{kn}^g \|d\|_2.$$

### 3.2. Error analysis of GMRES-IR5 with general $u_g$ and $u_p$ precisions.

We proceed in three steps. First, we bound the error in  $\hat{s}_i = \text{fl}(\hat{U}^{-1} \text{fl}(\hat{L}^{-1} \hat{r}_i))$ . Second, we use our analysis of three-precision MGS-GMRES of the previous section to prove the backward stability of the solution to the system  $\tilde{A}d_i = \hat{s}_i$ . Third, we combine the previous two results to derive bounds of the type (2.1a)–(2.1b) for the solution of  $\tilde{A}d_i = s_i$ .

We begin by bounding the error introduced in forming the preconditioned right-hand side  $s_i = \hat{U}^{-1} \hat{L}^{-1} \hat{r}_i$  in precision  $u_p$ . The computed  $\hat{s}_i$  satisfies [12, Thm. 8.5]

$$(\hat{L} + \Delta L)(\hat{U} + \Delta U) \hat{s}_i = \hat{r}_i, \quad |\Delta L| \leq \gamma_n^p |\hat{L}|, \quad |\Delta U| \leq \gamma_n^p |\hat{U}|. \quad (3.11)$$

Also, considering [12, Thm. 9.3] the LU factors computed at precision  $u_f$  satisfy

$$\hat{L}\hat{U} = A + \Delta A^{(1)}, \quad |\Delta A^{(1)}| \leq \gamma_n^f |\hat{L}||\hat{U}|, \quad (3.12)$$

so

$$\begin{aligned} s_i - \hat{s}_i &= \hat{U}^{-1} \hat{L}^{-1} (\Delta L \hat{U} + \hat{L} \Delta U + \Delta L \Delta U) \hat{s}_i \\ &= (A + \Delta A^{(1)})^{-1} (\Delta L \hat{U} + \hat{L} \Delta U + \Delta L \Delta U) \hat{s}_i \\ &\approx (A^{-1} - A^{-1} \Delta A^{(1)} A^{-1}) (\Delta L \hat{U} + \hat{L} \Delta U + \Delta L \Delta U) \hat{s}_i \end{aligned}$$

and dropping second-order terms we obtain

$$\begin{aligned} \|s_i - \hat{s}_i\|_\infty &\lesssim \gamma_{2n}^p \|A^{-1} |\hat{L}||\hat{U}|\|_\infty \|\hat{s}_i\|_\infty \\ &\leq n^2 \rho_n \tilde{\gamma}_n^p \kappa_\infty(A) \|\hat{s}_i\|_\infty \lesssim n^2 \rho_n \tilde{\gamma}_n^p \kappa_\infty(A) \|s_i\|_\infty. \end{aligned} \quad (3.13)$$

where the second inequality comes from [12, Lem. 9.6].

Next, we show that this new version of GMRES (with general  $u_p$  and  $u_g$  precisions) provides a backward stable solution to the system  $\tilde{A}d_i = s_i$ , where  $\tilde{A} = \hat{U}^{-1}\hat{L}^{-1}A$  and  $s_i = \hat{U}^{-1}\hat{L}^{-1}r_i$ . We rely on Theorem 3.1, which provides backward error bounds for MGS-GMRES with general matrix–vector products. Our aim is therefore to prove that (3.2) holds for some  $\epsilon_p$  when the matrix–vector products are computed with matrix  $\tilde{A} = \hat{U}^{-1}\hat{L}^{-1}A$  and in precision  $u_p$ . Let  $z_j = \tilde{A}\hat{v}_j$  be computed in precision  $u_p$  by a matrix product followed by two triangular solves. Then

$$\begin{aligned}(A + \Delta A^{(2)})\hat{v}_j &= \hat{w}_j, & |\Delta A^{(2)}| &\leq \gamma_q^p |A|, \\ (\hat{L} + \Delta L)\hat{g}_j &= \hat{w}_j, & |\Delta L| &\leq \gamma_n^p \hat{L}, \\ (\hat{U} + \Delta U)\hat{z}_j &= \hat{g}_j, & |\Delta U| &\leq \gamma_n^p \hat{U}.\end{aligned}$$

The computed vector  $\hat{z}_j$  can therefore be written as

$$\begin{aligned}\hat{z}_j &= (\hat{U} + \Delta U)^{-1}(\hat{L} + \Delta L)^{-1}(A + \Delta A^{(2)})\hat{v}_j \\ &\approx (\hat{U}^{-1} - \hat{U}^{-1}\Delta U\hat{U}^{-1})(\hat{L}^{-1} - \hat{L}^{-1}\Delta L\hat{L}^{-1})(A + \Delta A^{(2)})\hat{v}_j \\ &= \tilde{A}\hat{v}_j + f_j,\end{aligned}$$

where

$$\begin{aligned}f_j &\approx (\hat{U}^{-1}\hat{L}^{-1}\Delta A^{(2)} - \hat{U}^{-1}\hat{L}^{-1}\Delta L\hat{L}^{-1}A - \hat{U}^{-1}\Delta U\hat{U}^{-1}\hat{L}^{-1}A)\hat{v}_j \\ &= (\tilde{A}A^{-1}\Delta A^{(2)} - \hat{U}^{-1}\hat{L}^{-1}\Delta L\hat{U}\tilde{A} - \hat{U}^{-1}\Delta U\tilde{A})\hat{v}_j,\end{aligned}$$

and so

$$\|f_j\|_2 \lesssim \gamma_n^p(\kappa_F(A) + \kappa_F(\hat{U})\kappa_F(\hat{L}) + \kappa_F(\hat{U}))\|\tilde{A}\|_F\|\hat{v}_j\|_2.$$

Since we can expect  $\kappa_F(\hat{L})$  to be of modest size (because  $\hat{L}$  is a unit triangular matrix with off-diagonal elements bounded by 1 if, for example, partial pivoting is used in the factorization), and since  $\kappa_F(\hat{U}) \lesssim \kappa_F(A)\kappa_F(\hat{L})$ , we obtain

$$\|f_j\|_2 \lesssim \tilde{\gamma}_n^p \kappa_F(A)\|\tilde{A}\|_F\|\hat{v}_j\|_2 \leq n\tilde{\gamma}_n^p \kappa_\infty(A)\|\tilde{A}\|_F\|\hat{v}_j\|_2. \quad (3.14)$$

Condition (3.2) is thus satisfied for  $\epsilon_p = n\tilde{\gamma}_n^p \kappa_\infty(A)$ , and Theorem 3.1 is therefore applicable. From (3.4) we obtain

$$(\tilde{A} + \Delta\tilde{A})\hat{d}_i = \hat{s}_i + \Delta\hat{s}_i, \quad (3.15a)$$

$$\|\Delta\tilde{A}\|_F \lesssim (\tilde{\gamma}_{kn}^g + n\tilde{\gamma}_{k^{1/2}n}^p \kappa_\infty(A))\|\tilde{A}\|_F, \quad (3.15b)$$

$$\|\Delta\hat{s}_i\|_2 \lesssim \tilde{\gamma}_{kn}^g \|\hat{s}_i\|_2 \lesssim n^{1/2}\tilde{\gamma}_{kn}^g \|s_i\|_\infty. \quad (3.15c)$$

Rewriting (3.15a) as

$$s_i - \tilde{A}\hat{d}_i = \Delta\tilde{A}\hat{d}_i - (\hat{s}_i - s_i) - \Delta\hat{s}_i \quad (3.16)$$

and using (3.15b), (3.13), and (3.15c) to bound the three terms on the right-hand side, we obtain

$$\begin{aligned}\|s_i - \tilde{A}\hat{d}_i\|_\infty &\leq \|\Delta\tilde{A}\|_\infty \|\hat{d}_i\|_\infty + \|\hat{s}_i - s_i\|_\infty + \|\Delta\hat{s}_i\|_\infty \\ &\lesssim n(\tilde{\gamma}_{kn}^g + n\tilde{\gamma}_{k^{1/2}n}^p \kappa_\infty(A))\|\tilde{A}\|_\infty \|\hat{d}_i\|_\infty \\ &\quad + n^2 \rho_n \tilde{\gamma}_n^p \kappa_\infty(A) \|s_i\|_\infty + n^{1/2} \tilde{\gamma}_{kn}^g \|s_i\|_\infty \\ &\lesssim n^3 \max(k^{1/2}, \rho_n) (\tilde{\gamma}_1^g + \tilde{\gamma}_1^p \kappa_\infty(A)) (\|\tilde{A}\|_\infty \|\hat{d}_i\|_\infty + \|s_i\|_\infty).\end{aligned}$$

In conclusion, the normwise relative backward error of the system  $\tilde{A}\hat{d}_i = s_i$  is bounded by

$$\frac{\|s_i - \tilde{A}\hat{d}_i\|_\infty}{\|\tilde{A}\|_\infty \|\hat{d}_i\|_\infty + \|s_i\|_\infty} \lesssim f(n, k, \rho_n) (\tilde{\gamma}_1^g + \tilde{\gamma}_1^p \kappa_\infty(A)) \quad (3.17)$$

and the relative error of the computed  $\hat{d}_i$  therefore satisfies

$$\frac{\|d_i - \hat{d}_i\|_\infty}{\|d_i\|_\infty} \lesssim f(n, k, \rho_n) (\tilde{\gamma}_1^g + \tilde{\gamma}_1^p \kappa_\infty(A)) \kappa_\infty(\tilde{A}), \quad (3.18)$$

where  $f(n, k, \rho_n) = n^3 \max(k^{1/2}, \rho_n)$ .

**3.3. Convergence conditions on  $\kappa(A)$ .** We can now use this analysis together with (2.2) and (2.3) to determine sufficient conditions for the convergence of the forward and backward errors for general  $u_p$  and  $u_g$  parameters.

Beginning with the forward error, (3.18) shows that (2.1a) holds with

$$u_s \|E_i\|_\infty \equiv f(n, k, \rho_n) (\tilde{\gamma}_1^g + \tilde{\gamma}_1^p \kappa_\infty(A)) \kappa_\infty(\tilde{A}),$$

and so, dropping constants, (2.2) becomes

$$(u_g + u_p \kappa(A)) \kappa(\tilde{A}) \ll 1. \quad (3.19)$$

By using the bound (2.5) on  $\kappa(\tilde{A})$ , we obtain the condition

$$(u_g + u_p \kappa(A)) \kappa(A)^2 u_f^2 \ll 1 \quad (3.20)$$

for the forward error to converge to its limiting value of order  $qu_r \text{cond}(A, x) + u$ .

Next we determine a condition for the backward error to converge. First we need to bound the backward error of the original correction equation  $Ad_i = r_i$ . By observing that the computed residual  $\hat{r}_i$  satisfies  $\hat{r}_i - Ad_i = \hat{L}\hat{U}(s_i - \tilde{A}\hat{d}_i)$ , since  $s_i = \hat{U}^{-1}\hat{L}^{-1}\hat{r}_i$ , and using the bound (3.17) and the relation (3.12) we obtain

$$\begin{aligned} \|\hat{r}_i - Ad_i\|_\infty &\lesssim f(n, k, \rho_n) (\tilde{\gamma}_1^g + \tilde{\gamma}_1^p \kappa_\infty(A)) \|\hat{L}\hat{U}\|_\infty (\|\tilde{A}\|_\infty \|\hat{d}_i\|_\infty + \|s_i\|_\infty) \\ &\lesssim f(n, k, \rho_n) (\tilde{\gamma}_1^g + \tilde{\gamma}_1^p \kappa_\infty(A)) (\|\tilde{A}\|_\infty \|A\|_\infty \|\hat{d}_i\|_\infty + \kappa_\infty(A) \|\hat{r}_i\|_\infty). \end{aligned}$$

We have thus shown that (2.1b) holds with  $u_s = u_g + u_p \kappa_\infty(A)$ ,  $c_1 = f(n, k, \rho_n) \|\tilde{A}\|_\infty$ , and  $c_2 = f(n, k, \rho_n) \kappa_\infty(A)$ . Dropping constants, condition (2.3) then becomes

$$(u_g + u_p \kappa(A)) (\|\tilde{A}\| + 1) \kappa(A) \ll 1.$$

Using (3.12) we have

$$\tilde{A} = \hat{U}^{-1}\hat{L}^{-1}A = (A + \Delta A^{(1)})^{-1}A \approx I - A^{-1}\Delta A^{(1)}$$

and using [12, Lem. 9.6] we obtain the bound

$$\|\tilde{A}\|_\infty \leq 1 + \|A^{-1}\|_\infty \|\Delta A^{(1)}\|_\infty \ll 1 + u_f \kappa_\infty(A).$$

We finally obtain the condition

$$(u_g + u_p \kappa(A)) (1 + \kappa(A) u_f) \kappa(A) \ll 1 \quad (3.21)$$

TABLE 3.1

Parameters for floating-point arithmetics: symbol used in this paper, number of bits in significand, number of bits in exponent, unit roundoff, and range.

Arithmetic	Symbol	Significand	Exponent	Unit roundoff	Range
bfloat16	B	8	8	$3.91 \times 10^{-3}$	$10^{\pm 38}$
fp16	H	11	5	$4.88 \times 10^{-4}$	$10^{\pm 5}$
fp32	s	24	8	$5.96 \times 10^{-8}$	$10^{\pm 38}$
fp64	D	53	11	$1.11 \times 10^{-16}$	$10^{\pm 308}$
fp128	Q	113	15	$9.63 \times 10^{-35}$	$10^{\pm 4932}$

for the backward error to converge to its limiting value of order  $qu_r + u$ .

Curiously, condition (3.21) is stricter than the condition for the forward error, since (3.20) has an extra  $u_f$  term. However, note that since the backward error is bounded by the forward error (see, for example, [12, Prob. 7.2]), (3.20) is also an obvious condition for the backward error to converge to the limiting value of the forward error, of order  $q \text{cond}(A, x)u_r + u$ . In particular, if  $u_r = u^2$ , condition (3.21) is not useful because (3.20) also guarantees a backward error of order  $u$  with a less restrictive condition on  $\kappa(A)$ , since by assumption  $\kappa(A)u < 1$ .

As a check, we compare our results with [14], which analyzes the case  $u_p = u_g = u_r = u$ , that is, GMRES-IR in two precisions, and uses a different argument to that here. Our condition (3.20) for convergence of the forward error is  $\kappa(A)^3 uu_f^2 \ll 1$  and our condition (3.21) for convergence of the backward error is  $\kappa(A)^2 u(1 + \kappa(A)u_f) \ll 1$ . These conditions agree with [14, Eqs (3.8), (3.6)].

We also compare with the analysis of [3] for GMRES-IR3, the conclusions from which we summarized in section 2. For  $u_f = u_g = u$  and  $u_p = u^2$ , our condition (3.20) for convergence of the forward error is  $(u + \kappa(A)u^2)\kappa(A)^2 u_f^2 \ll 1$ , that is,  $\kappa(A)^2 uu_f^2 \ll 1$ , since we are assuming  $\kappa(A)u < 1$ . This agrees with (2.6). Our condition (3.21) for convergence of the backward error is, similarly,  $\kappa(A)u(1 + \kappa(A)u_f) \ll 1$ , which is weaker than the condition  $\kappa(A)u \ll 1$  of (2.7). The reason for this discrepancy is that in [3] the assumption is made that  $c_1$  in (2.3) is of order  $f(n, k, \rho_n)$ , which is too optimistic in general.

**3.4. Identifying meaningful combinations of precisions.** Table 3.1 shows the five most widely used floating-point arithmetics. If all these arithmetics are available then GMRES-IR5 has over three thousand different combinations of its five precision parameters. It is therefore important to identify the subset of meaningful combinations. We say a combination is meaningful if none of the precisions it employs can be reduced without degrading the numerical properties (convergence condition or limiting accuracy). Every meaningful combination attains a tradeoff between performance, robustness (ability to converge for ill-conditioned matrices), and accuracy (ability to converge to small errors).

As an example, the meaningful combinations for LU-IR3 must satisfy  $u^2 \leq u_r \leq u \leq u_f$ . Indeed, the limiting backward and forward errors (recalled in Table 2.1) show that we should have  $u_r \leq u$  and that setting  $u_r < u^2$  is not useful since  $u_r = u^2$  is already enough to ensure a forward error of order  $u$  (since by assumption  $\kappa(A)u < 1$ ).

Meaningful combinations for GMRES-IR5 also satisfy  $u^2 \leq u_r \leq u \leq u_f$ , so our aim now is to discuss the choice of the two new precision parameters  $u_g \geq u$  and  $u_p \geq u^2$ . To compute the bounds on  $\kappa(A)$  given by the conditions (3.20) and (3.21), we solve the equalities  $(u_g + u_p x)u_f^2 x^2 = 1$  and  $(u_g + u_p x)(1 + u_f x)x = 1$ , respectively. We now state some observations that can be deduced from our analysis, in particular

from condition (3.20).

- $u_p \leq u_g$ . This first observation comes from the term  $u_g + u_p \kappa(A)$  that appears in the convergence conditions. We would like the components of this term to be balanced, so that  $u_p \approx u_g / \kappa(A) \leq u_g$ . So  $u_p < u_g$  may be required, but  $u_p > u_g$  is not meaningful. We also note that there is no advantage to taking  $u_p < u_g / \kappa(A)$ .
- $u_p < u_f$ . This second observation comes from the fact that if  $u_p = u_f$ , condition (3.19) requires  $\kappa(A)u_f \kappa(\tilde{A}) \ll 1$ , which is worse than the condition  $\kappa(A)u_f \ll 1$  for LU-IR3.
- $u_p < u$ ,  $u_p = u$ , and  $u_p > u$  are all meaningful. This is one of the main conclusions of our analysis. We know from [3] (in which  $u_p = u^2$  and  $u_g = u$ ) that setting  $u_p = u^2$  provides the least restrictive convergence conditions (2.6) and (2.7), but the precise role of  $u_p$  in the convergence was not analyzed. With the new conditions (3.20) and (3.21) obtained from our generalized analysis, we can now understand what the conditions become if  $u_p$  is taken larger than  $u^2$ . Crucially, setting  $u_p = u$  and even  $u_p > u$  can potentially yield conditions that remain less restrictive than the LU-IR3 condition  $\kappa(A)u_f \ll 1$ , and therefore represent meaningful combinations. Let us illustrate this observation with a practical example. Assume  $u_f$  is set to fp16 and  $u = u_g$ , with  $u_g$  set to fp64. Then the condition for the forward error to converge with LU-IR3 is  $\kappa(A) \ll 2 \times 10^3$ . Instead, with GMRES-IR5:
  - If  $u_p$  is set to fp128 (as in [3]), condition (3.20) becomes  $\kappa(A) \ll 2 \times 10^{11}$ . We recover the same condition as in [3].
  - If  $u_p$  is set to fp64 (and thus  $u_p = u$ ), condition (3.20) becomes  $\kappa(A) \ll 3 \times 10^7$ , which is still much better than the LU-IR3 condition. This version of GMRES-IR was analyzed in [14] and has been successfully used in practice, for example in [9], [10].
  - If  $u_p$  is set to fp32 (and thus  $u_p > u$ ), condition (3.20) becomes  $\kappa(A) \ll 4 \times 10^4$ , which is still over an order of magnitude better than the LU-IR3 condition. Note that this variant uses up to four different precisions if  $u_r$  is set to fp128.

These examples illustrate how the  $u_p$  precision can be tuned to achieve different levels of tradeoff between robustness (ability to handle ill-conditioned matrices) and performance (cost of the application of LU factors within GMRES).

- $u_g = u$  and  $u_g > u$  are both meaningful. This is also an important conclusion of our analysis. Whereas in [3],  $u_g$  is set to  $u$  to obtain the least restrictive convergence condition, our analysis reveals that setting  $u_g > u$  can also be meaningful. In fact, as long as  $u_g < 1$ , condition (3.20) is better than the LU-IR3 condition  $\kappa(A)u_f \ll 1$ , so we have much flexibility in choosing  $u_g$ . Let us again illustrate this with a practical example. Assume  $u_f$  is set to fp16 and  $u = u_p$  are set to fp64. With  $u_g = u$ , as previously stated the condition (3.20) is  $\kappa(A) \ll 3 \times 10^7$ . However, if  $u_g$  is set to fp32, the condition is  $\kappa(A) \ll 8 \times 10^6$ , which is only slightly worse and remains much better than the LU-IR3 condition. Note that with  $u_g$  set to fp32, setting  $u_p$  to fp128 instead of fp64 does not improve the condition  $\kappa(A) \ll 8 \times 10^6$ : this shows that the meaningful values of  $u_p$  can be influenced by the choice of  $u_g$ , and vice versa. In conclusion the  $u_g$  precision can also be tuned to achieve different levels of tradeoff between robustness and performance (the cost of GMRES

and, in particular, the memory footprint of the Krylov basis).

- $u_g < u_f$ ,  $u_g = u_f$ , and  $u_g > u_f$  are all meaningful. This final observation is that  $u_g$  and  $u_f$  can be independent. We have already illustrated in previous examples where  $u_f$  is fp16 that  $u_g < u_f$  and  $u_g = u_f$  are both meaningful. This final observation states that even  $u_g > u_f$  can be meaningful. To see why, let us take another example with  $u_f$  set to fp32 and  $u = u_p$  set to fp64. Then setting  $u_g$  to fp16 yields the condition  $\kappa(A) \ll 7 \times 10^8$ , which is better than any combination with the same  $u$  and  $u_p$  but with  $u_f$  set to fp16 (the best possible condition being  $\kappa(A) \ll 3 \times 10^7$  for  $u_g$  at least in fp64). Usually, one would expect an fp16 LU factorization and fp32 GMRES to be faster than the converse (fp32 factorization and fp16 GMRES), so in practice the combinations with  $u_g > u_f$  are only relevant for a narrow range of  $\kappa(A)$  (in this example, for  $\kappa(A)$  such that  $\kappa(A) \ll 3 \times 10^7$  is not satisfied but  $\kappa(A) \ll 7 \times 10^8$  is).

In Table 3.2, we summarize all the possible combinations of  $u_p$  and  $u_g$  when  $u_f$  is set to bfloat16, fp16, or fp32 (these three precisions being the ones that determine the convergence conditions). The underlined, red terms correspond to combinations of these three precisions that are not meaningful. Interestingly, when the working precision  $u$  is fp64 and the residual precision  $u_r$  is fp128, Table 3.2 contains thirteen meaningful variants that use at least four different arithmetics, including two that use all five (BHS and HBS).

All the observations made in this part reduce the number of meaningful combinations of the five precision parameters to a subset of 112 combinations, corresponding to 3.5% of all the possibilities. We summarize the numerical properties (convergence condition and limiting accuracy, for both the forward and backward errors) of a selected subset of 42 GMRES-IR5 variants in Table 3.3. Importantly, the table includes variants used in existing implementations, for which our analysis provides new theoretical guarantees. It also includes new combinations of precisions not proposed previously that achieve new, finer tradeoffs between the convergence conditions and the precisions used.

**4. Numerical experiments.** We now perform numerical experiments to assess the validity of the convergence conditions of GMRES-IR5 derived in the previous section. Throughout our experiments, we focus on the forward error convergence (condition (3.20)) and we fix  $u = \text{D}$  and  $u_r = \text{Q}$  and analyze the role of the factorization precision  $u_f$  and that of the newly introduced precisions  $u_p$  (preconditioner precision) and  $u_g$  (GMRES precision).

We have written a Julia code that implements GMRES-IR5 and LU-IR3, where half precision arithmetics (fp16 and bfloat16) are simulated. We have made this code publicly available<sup>1</sup>.

**4.1. Random dense matrices.** We first use random dense matrices with specified condition 2-norm number  $\kappa(A)$ , which are generated in Julia with the command `matrixdepot('randsvd', n, kappa, mode)` where `mode = 2`. Note that this class of matrices leads to unusually large growth factors  $\rho_n$  of order  $n$  [11]. However, we use only small matrices ( $n = 50$ ), for which  $\rho_n$  does not exceed 20.

We take  $\kappa(A) = 10^c$ , for  $c = 0:17$  and, for each value of  $\kappa(A)$ , we generate 100 random  $50 \times 50$  matrices of corresponding condition number. Then, we run LU-IR3 and GMRES-IR5 on each matrix and compute their success rate, that is, the

<sup>1</sup><https://github.com/bvieuble/Itref.jl>

TABLE 3.2

Bound on  $\kappa(A)$ , rounded to one significant figure, given by conditions (3.20) and (3.21) for the forward and backward errors to converge with GMRES-IR5, depending on the precisions  $u_f$ ,  $u_g$ , and  $u_p$ , and assuming the working precision  $u$  is double. We recall that the forward error convergence condition for LU-IR3 is  $\kappa(A) \ll 3 \times 10^2$  for  $u_f = \text{B}$ ,  $\kappa(A) \ll 2 \times 10^3$  for  $u_f = \text{H}$ , and  $\kappa(A) \ll 2 \times 10^7$  for  $u_f = \text{S}$ . The red, underlined terms denote combinations of precisions that are not meaningful. Each variant is presented in the form of a triplet  $(u_f, u_g, u_p)$ , hence HDQ means  $u_f = \text{H}$ ,  $u_g = \text{D}$ , and  $u_p = \text{Q}$ . We recall that  $u_r$  does not play a role in the convergence bounds.

	Variants $(u_f, u_g, u_p)$	Forward	Backward	
$u_f = \text{B}$	BBH <u>BHH BSH BDH</u>	$5 \times 10^2$	$4 \times 10^1$	§
	BBS <u>BBD BBQ</u>	$4 \times 10^3$	$2 \times 10^2$	
	BHS	$8 \times 10^3$	$6 \times 10^2$	
	BHD <u>BHQ</u>	$1 \times 10^4$	$6 \times 10^2$	
	BSS <u>BDS</u>	$1 \times 10^4$	$2 \times 10^3$	
	BSD <u>BDQ</u>	$1 \times 10^6$	$7 \times 10^4$	
	BDD	$8 \times 10^6$	$1 \times 10^6$	
	BDQ	$2 \times 10^{10}$	$2 \times 10^9$	
$u_f = \text{H}$	HBS <u>HBD HBQ</u>	$3 \times 10^4$	$2 \times 10^2$	
	HHS	$4 \times 10^4$	$1 \times 10^3$	
	HSS <u>HDS</u>	$4 \times 10^4$	$3 \times 10^3$	
	HHD <u>HHQ</u>	$9 \times 10^4$	$1 \times 10^3$	
	HSD <u>HSQ</u>	$8 \times 10^6$	$2 \times 10^5$	
	HDD	$3 \times 10^7$	$3 \times 10^6$	
	HDQ	$2 \times 10^{11}$	$4 \times 10^9$	
$u_f = \text{S}$	SBD <u>SBQ</u>	$3 \times 10^8$	$3 \times 10^2$	
	SHD <u>SHQ</u>	$7 \times 10^8$	$2 \times 10^3$	
	SSD	$1 \times 10^{10}$	$1 \times 10^7$	
	SDD	$1 \times 10^{10}$	$5 \times 10^7$	
	SSQ	$7 \times 10^{10}$	$1 \times 10^7$	
	SDQ	$2 \times 10^{15}$	$4 \times 10^{11}$	

§ BHH is slightly better than BBH on Forward and Backward before rounding.

percentage of matrices for which a forward error  $\|x - \hat{x}\|_2 / \|x\|_2 \leq 4.44 \times 10^{-16}$  is achieved, since with  $u = \text{D}$  and  $u_r = \text{Q}$  the forward error should reach full double precision with no dependency on  $\text{cond}(A, x)$ .

Figure 4.1 reports the success rate of LU-IR3 and nine variants of GMRES-IR5, corresponding to all possible combinations of the  $u_p$  and  $u_g$  parameters over the values  $u_p = \text{S}, \text{D}, \text{Q}$  and  $u_g = \text{B}, \text{S}, \text{D}$ , with  $u_f = \text{B}$  fixed for all variants. This experiment allows us to obtain an empirical bound on the value of  $\kappa(A)$  at which each variant stops being able to converge. For example, the success rate of LU-IR3 is 100% as long as  $\kappa(A) \leq 10^2$ , but starts decreasing for larger  $\kappa(A)$ , and quickly becomes 0%. This experimentally confirms the theoretical condition  $\kappa(A) \ll 2 \times 10^3$  given by (2.4).

Let us now analyze the GMRES-IR5 variants, starting with the role of the preconditioner precision  $u_p$  with fixed  $u_g = \text{D}$  (top graph of Figure 4.1). We can observe that convergence is achieved with 100% success rate as long as  $\kappa(A)$  is smaller than  $10^7$  for  $u_p = \text{S}$  and  $10^{15}$  for  $u_p = \text{D}$  or  $\text{Q}$ . The relative robustness of each method is therefore consistent with the theoretical bounds of Table 3.2. However, the variants with  $u_p = \text{S}$  and  $u_p = \text{D}$  both perform much better than expected. This is not entirely surprising since the analysis can be pessimistic, especially in the bound (2.5) on  $\kappa(\tilde{A})$ , which in practice has been observed to often be of order  $\kappa(A)u_f$  rather than the worst-case bound  $(\kappa(A)u_f)^2$  [2], [21], [23].

Next we analyze the role of the GMRES precision  $u_g$  by comparing the top graph of Figure 4.1 with the middle and bottom ones. When  $u_p = \text{S}$ , switching from  $u_g = \text{D}$

TABLE 3.3

Bounds on  $\kappa(A)$ , rounded to one significant figure, such that GMRES-IR5 with the precisions in the first five columns is guaranteed to converge to the indicated limiting error values.

$u$	$u_r$	$u_f$	$u_g$	$u_p$	Forward error		Backward error	
					$\kappa(A)$	Limit	$\kappa(A)$	Limit
S	S	B	B	S	$4 \times 10^3$	$\text{cond}(A, x) \times S$	$2 \times 10^2$	S
S	S	B	H	S	$8 \times 10^3$	$\text{cond}(A, x) \times S$	$6 \times 10^2$	S
S	S	B	S	S	$1 \times 10^4$	$\text{cond}(A, x) \times S$	$2 \times 10^3$	S
S	S	H	H	S	$4 \times 10^4$	$\text{cond}(A, x) \times S$	$1 \times 10^3$	S
S	S	H	S	S	$4 \times 10^4$	$\text{cond}(A, x) \times S$	$3 \times 10^3$	S
S	S	B	S	D	$1 \times 10^6$	$\text{cond}(A, x) \times S$	$7 \times 10^4$	S
S	S	H	S	D	$8 \times 10^6$	$\text{cond}(A, x) \times S$	$2 \times 10^5$	S
S	D	B	B	S	$4 \times 10^3$	S	$2 \times 10^2$	S
S	D	B	H	S	$8 \times 10^3$	S	$6 \times 10^2$	S
S	D	B	S	S	$1 \times 10^4$	S	$2 \times 10^3$	S
S	D	H	H	S	$4 \times 10^4$	S	$1 \times 10^3$	S
S	D	H	S	S	$4 \times 10^4$	S	$3 \times 10^3$	S
S	D	B	S	D	$1 \times 10^6$	S	$7 \times 10^4$	S
S	D	H	S	D	$8 \times 10^6$	S	$2 \times 10^5$	S
D	D	B	B	S	$4 \times 10^3$	$\text{cond}(A, x) \times D$	$2 \times 10^2$	D
D	D	B	H	S	$8 \times 10^3$	$\text{cond}(A, x) \times D$	$6 \times 10^2$	D
D	D	B	S	S	$1 \times 10^4$	$\text{cond}(A, x) \times D$	$2 \times 10^3$	D
D	D	H	H	S	$4 \times 10^4$	$\text{cond}(A, x) \times D$	$1 \times 10^3$	D
D	D	H	S	S	$4 \times 10^4$	$\text{cond}(A, x) \times D$	$3 \times 10^3$	D
D	D	B	S	D	$1 \times 10^6$	$\text{cond}(A, x) \times D$	$7 \times 10^4$	D
D	D	H	S	D	$8 \times 10^6$	$\text{cond}(A, x) \times D$	$2 \times 10^5$	D
D	D	B	D	D	$8 \times 10^6$	$\text{cond}(A, x) \times D$	$1 \times 10^6$	D
D	D	H	D	D	$3 \times 10^7$	$\text{cond}(A, x) \times D$	$3 \times 10^6$	D
D	D	S	H	D	$7 \times 10^8$	$\text{cond}(A, x) \times D$	$2 \times 10^3$	D
D	D	S	D	D	$1 \times 10^{10}$	$\text{cond}(A, x) \times D$	$5 \times 10^7$	D
D	D	B	D	Q	$2 \times 10^{10}$	$\text{cond}(A, x) \times D$	$2 \times 10^9$	D
D	D	H	D	Q	$2 \times 10^{11}$	$\text{cond}(A, x) \times D$	$4 \times 10^9$	D
D	D	S	D	Q	$2 \times 10^{15}$	$\text{cond}(A, x) \times D$	$4 \times 10^{11}$	D
D	Q	B	B	S	$4 \times 10^3$	D	$2 \times 10^2$	D
D	Q	B	H	S	$8 \times 10^3$	D	$6 \times 10^2$	D
D	Q	B	S	S	$1 \times 10^4$	D	$2 \times 10^3$	D
D	Q	H	H	S	$4 \times 10^4$	D	$1 \times 10^3$	D
D	Q	H	S	S	$4 \times 10^4$	D	$3 \times 10^3$	D
D	Q	B	S	D	$1 \times 10^6$	D	$7 \times 10^4$	D
D	Q	B	D	D	$8 \times 10^6$	D	$1 \times 10^6$	D
D	Q	H	S	D	$8 \times 10^6$	D	$2 \times 10^5$	D
D	Q	H	D	D	$3 \times 10^7$	D	$3 \times 10^6$	D
D	Q	S	H	D	$7 \times 10^8$	D	$2 \times 10^3$	D
D	Q	S	D	D	$1 \times 10^{10}$	D	$5 \times 10^7$	D
D	Q	B	D	Q	$2 \times 10^{10}$	D	$2 \times 10^9$	D
D	Q	H	D	Q	$2 \times 10^{11}$	D	$4 \times 10^9$	D
D	Q	S	D	Q	$2 \times 10^{15}$	D	$4 \times 10^{11}$	D

(top) to  $u_g = S$  (middle) has no impact on the success rate, which equals 100% as long as  $\kappa(A) \leq 10^7$ ; reducing the precision even further by setting  $u_g = B$  (bottom) only has a slight impact: the success rate remains at 100% for  $\kappa(A) \leq 10^5$ . When  $u_p = D$  or Q, reducing the GMRES precision has a much more visible impact: a 100% success rate is achieved only when  $\kappa(A) \leq 10^9$  ( $u_g = S$ , middle) or  $\kappa(A) \leq 10^5$  ( $u_g = B$ , bottom).

These experiments also show that the success rate is independent of  $u_p$  when



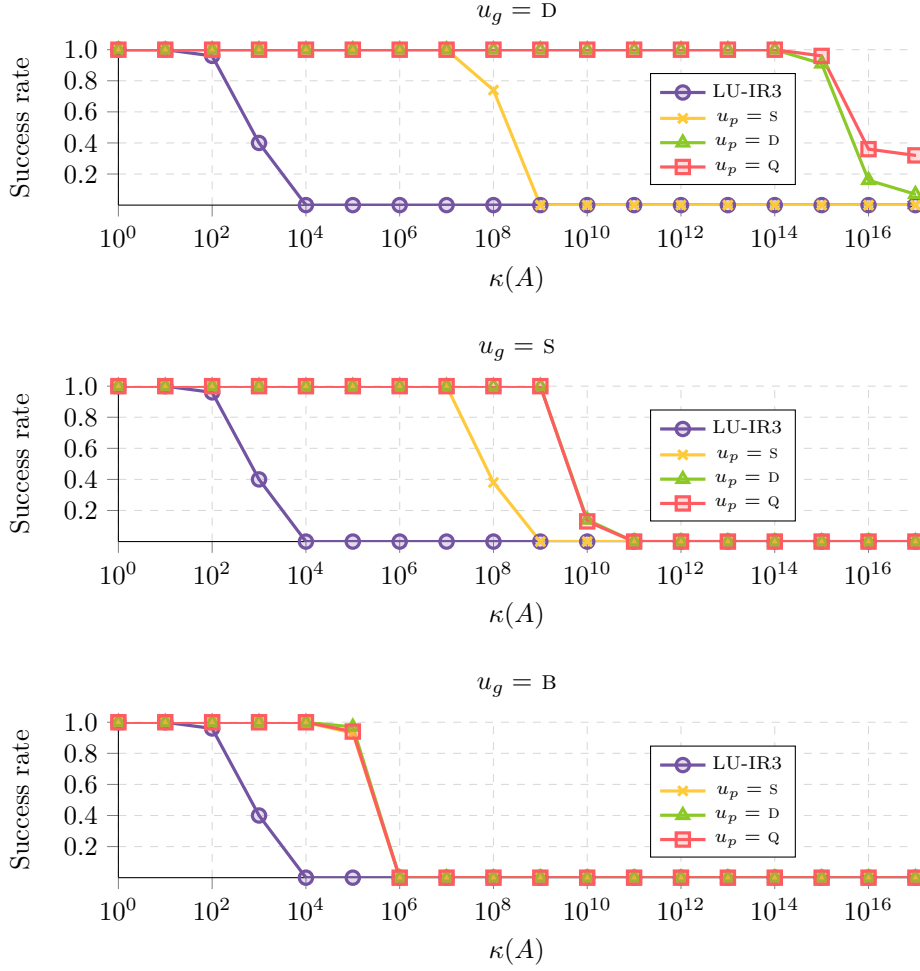


FIG. 4.1. Proportion of matrices for which LU-IR3 and GMRES-IR5 with varying  $u_p$  and  $u_g$  converged to double precision forward error as a function of  $\kappa(A)$ . For all variants,  $u_f = B$ ,  $u = D$ , and  $u_r = Q$ .

$u_g = B$ , and almost independent of  $u_g$  when  $u_p = S$ . This illustrates that the meaningful choices of  $u_p$  depend on the choice of  $u_g$ , and vice versa. Overall, the experiments of Figure 4.1 are therefore in good agreement with the theoretical bounds. Importantly, they confirm that even with relaxed requirements on the precisions  $u_p$  and  $u_g$ , GMRES-IR5 can still handle matrices that are much more ill-conditioned than LU-IR3.

Finally, we evaluate in Figure 4.2 the two meaningful variants of Table 3.2 that use five different precisions: BHS ( $u_f = B$ ,  $u_g = H$ ,  $u_p = S$ ) and HBS ( $u_f = H$ ,  $u_g = B$ ,  $u_p = S$ ). We compare these two variants with LU-IR3 with  $u_f = B$  or  $H$ , and with the BBS and HHS GMRES-IR5 variants ( $u_f = B$ ,  $u_g = B$ ,  $u_p = S$  and  $u_f = H$ ,  $u_g = H$ ,  $u_p = S$ ), which are the four-precision variants right below and above in terms of convergence condition. The figure experimentally confirms that both of these five-precision variants are in between BBS and HHS and are able to handle more

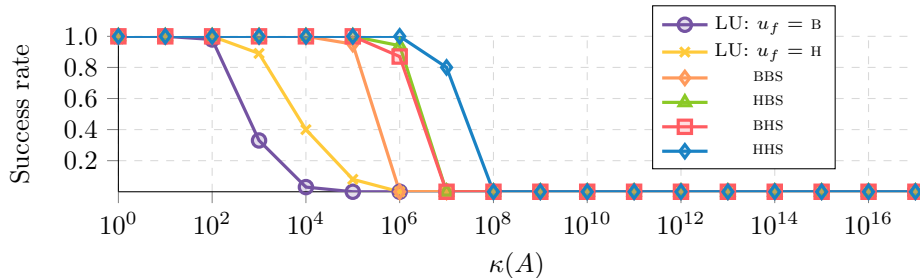


FIG. 4.2. Proportion of matrices for which LU-IR and GMRES-IR variants (including five meaningful precisions combinations, with  $u_f$ ,  $u_g$ ,  $u_p$  as specified in the legend) converged to double precision forward error as a function of  $\kappa(A)$ . For all variants  $u = D$  and  $u_r = Q$ .

ill-conditioned matrices than LU-IR3, showing that they are meaningful not only theoretically but also in practice.

**4.2. Real-life matrices from SuiteSparse.** It is also important to check whether reducing the preconditioner and/or GMRES precisions impacts the number of iterations required to converge. To answer this question, we use a set of 230 matrices from the SuiteSparse collection. These matrices are all real, square, and of dimension between 500 and 3000. We compare LU-IR3 with the meaningful GMRES-IR5 variants over the values  $u_f = B, H, S$ . Rather than counting the number of iterations, we count the number of calls to LU triangular solves: this is because in the case of GMRES-IR5, an extra LU solve is required at the start of each IR step (external iterations).

For the experiments that use an LU factorization in fp16 arithmetic, we must pay attention to the narrow range of this arithmetic (see Table 3.1). Many matrices in our set have entries outside of this range, and so it is essential to address this issue. We use the diagonal scaling method of [18], which first normalizes every row and column by its maximum value (preventing overflow), and then scales the matrix by a quantity  $\lambda$  close to the maximum representable value (to minimize underflow).

One difficulty is that the number of iterations is very sensitive to the choice of the GMRES stopping criterion  $\tau$  and the scaling factor  $\lambda$ , and that the optimal choice of  $\tau$  and  $\lambda$  is different for each variant and matrix. For the comparison to be as fair as possible, for each variant and each matrix, we have tested eight different values of  $\tau$  ( $10^{-10}$ ,  $10^{-8}$ ,  $10^{-6}$ ,  $10^{-4}$ ,  $10^{-3}$ ,  $10^{-2}$ ,  $10^{-1}$ , and  $5 \times 10^{-1}$ ) and five different values of  $\lambda$  ( $10^4$ ,  $10^3$ ,  $10^2$ ,  $10^1$ , and  $10^0$ ), and taken the values that leads to the lowest number of iterations. We note in passing that the parameter  $\tau$  could be taken into account in the error analysis. Indeed,  $\tau$  provides a bound on the backward error (3.15) of the solution of the system  $\tilde{A}d_i = s_i$  by GMRES, which readily affects the convergence conditions.

We present in Figure 4.3 performance profiles for  $u_f = B, H, S$ , which plot the percentage  $\phi$  of matrices for which a given variant converges in less than  $\alpha$  times the number of LU solves required by the best variant. We also provide detailed results on a representative sample of these matrices for the case  $u_f = B$  in Table 4.1.

These results are in agreement with our theoretical study and with the previous experiments of section 4.1, and provide a confirmation of the relative robustness of each variant on a dataset of real-life applications matrices. For example, with  $u_f = B$ ,

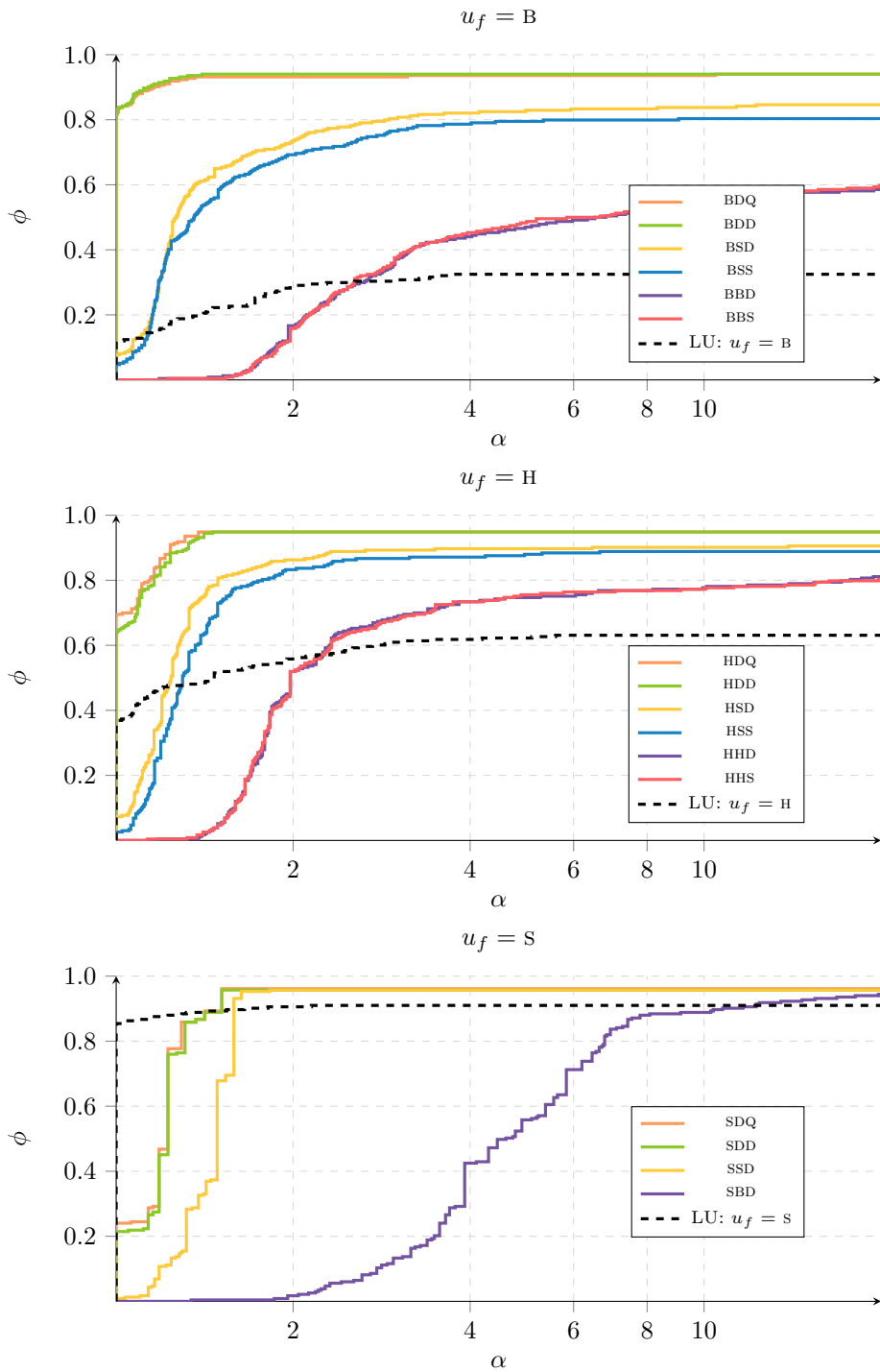


FIG. 4.3. Performance profile of LU-IR3 and GMRES-IR5 variants (with  $u_f$ ,  $u_g$ ,  $u_p$  as specified in the legend) on 230 SuiteSparse matrices.  $\phi$  (y-axis) indicates the percentage of matrices for which a given variant requires less than  $\alpha$  (x-axis) times the number of LU solves required by the best variant.

TABLE 4.1

Number of LU solves on a sample of SuiteSparse matrices for LU-IR3 and GMRES-IR5 variants with  $u_f = B$ .  $\kappa(A)$  is the original condition number,  $\kappa(A_s)$  is the condition number after scaling. The matrices are sorted by increasing  $\kappa(A_s)$ . A “—” indicates failure to converge.

Name	$n$	$\kappa(A)$	$\kappa(A_s)$	LU-IR3	BBS	BBD	BSS	BSD	BDD	BDQ
dw256B	512	3.7E+00	1.8E+00	8	24	24	10	10	10	10
gre_512	512	1.6E+02	1.7E+02	16	28	28	17	17	14	14
mahindas	1258	2.1E+13	9.8E+02	12	34	36	14	12	12	12
dw1024	2048	2.1E+03	1.9E+03	18	36	36	20	20	19	19
fs_541.4	541	1.2E+10	6.9E+03	—	48	48	33	35	28	28
rajat12	1879	6.9E+05	8.1E+03	—	94	69	29	29	24	24
sherman1	1000	1.6E+04	9.3E+03	—	78	78	30	26	26	26
watt_2	1856	1.4E+11	3.0E+04	—	522	497	32	32	26	26
bp_600	822	1.5E+06	6.6E+04	—	56	54	23	25	19	19
bwm2000	2000	2.4E+05	2.0E+05	—	—	—	237	197	148	148
meg1	2904	1.4E+12	3.4E+05	13	28	26	17	17	14	14
linsp_511	511	3.3E+15	5.8E+05	—	302	268	39	30	23	23
hangGlider_2	1647	1.4E+09	1.2E+06	—	52	56	32	33	31	31
tub1000	1000	1.3E+06	1.3E+06	—	—	—	354	254	114	114
1138.bus	1138	8.6E+06	2.4E+06	—	—	—	238	163	83	83
gre_1107	1107	3.2E+07	2.3E+07	—	—	—	244	185	75	75
rajat19	1157	1.1E+10	1.4E+08	—	—	—	—	—	70	69
spaceStation_7	1134	3.9E+11	1.7E+08	—	—	—	—	220	127	127
bcsstk19	817	1.3E+11	1.2E+09	—	—	—	—	—	307	209
reorientation_3	2513	1.5E+21	8.9E+11	—	—	—	—	—	144	188
fs_760.3	760	9.8E+19	2.6E+12	—	—	—	—	—	—	—
nnc1374	1374	3.7E+14	5.2E+12	—	—	—	—	—	85	92
lung1	1650	5.1E+19	1.1E+13	—	—	—	—	—	27	27

LU-IR3 only converges for about 35% of the matrices, whereas with GMRES-IR3 (BDQ variant), we are able to process about 95% of the matrices (the 5% remaining being highly ill-conditioned matrices). Crucially, the new GMRES-IR5 variants (with relaxed  $u_p$  and/or  $u_g$ ) are all more robust than LU-IR3, converging for a much larger percentage of matrices. In particular, on this set of matrices, the BDD variant is as robust as BDQ and requires the same number of LU solves. Therefore, these experiments show that, in practice, we can switch from  $u_p = Q$  to  $u_p = D$  with no impact on the convergence of GMRES-IR.

The other two performance profiles with  $u_f = H$  and  $S$  show similar trends. Unsurprisingly, as we increase the precision of the factorization, LU-IR3 is able to converge on a wider range of matrices and so the range of matrices where GMRES-IR is relevant is narrower. This is especially the case for an fp32 factorization, where LU-IR3 is able to converge on almost 90% of the matrices. Note that this observation heavily depends on the dataset: for these matrices, the distribution of the  $\kappa(A)$  is not uniform and is mainly concentrated between  $10^3$  and  $10^6$ , which explains why LU-IR3 performs well with an fp32 factorization.

A performance comparison of the actual runtime of the variants is outside our scope, but we can nevertheless extrapolate some performance trends based on Figure 4.3 and on the assumptions that (1) the LU solves dominate the overall runtime of the iterative phase of the solver; and (2) a bfloat16 LU solve is twice faster than an fp32 one, which is itself twice faster than an fp64 one. With this performance model and considering for example  $u_f = B$ , we can expect LU-IR3 (which uses bfloat16 LU solves) to be the fastest method as long as it does not require more than twice the number of LU solves of a variant with  $u_p = S$ , that is, for about 30% of the matrices.

ces. Similarly, the BSS variant should outperform BDD for over 70% of the matrices, making it the best variant for about  $70 - 30 = 40\%$  of the matrices. Finally, on our set of matrices, the BDD variant never requires more LU solves than the BDQ one and therefore it should be the best variant for the remaining 30% of the matrices. An high performance implementation of GMRES-IR5 is necessary to confirm these predictions, and will be the object of future work.

**5. Conclusions.** We have addressed the solution of linear systems of equations by means of iterative refinement in mixed precision. The renewed interest in this method stems from the trend for computer hardware to support fast low precision floating-point arithmetic. Our baseline is the work by Carson and Higham [3], who propose a method, GMRES-IR3, that employs up to three precisions and in which a preconditioned mixed precision GMRES method is used for solving the correction equation in order to converge on badly conditioned problems. In its original form, this method requires the use of a precision higher than the working precision, which can be expensive, especially when the working precision is double precision. By relaxing the assumptions on the precision of the operations within the GMRES solver we have proposed a GMRES-based iterative refinement method with up to five different precisions in play, which we call GMRES-IR5. We extended the rounding error analysis of Carson and Higham to cover GMRES-IR5. As a key component of this analysis, we extended the work of Paige et al. [22] to prove the backward stability of a mixed precision GMRES method. Based on this result we derived conditions on  $\kappa(A)$  that guarantee the convergence of GMRES-IR5.

Our results show that GMRES-IR5 can accurately and reliably solve relatively badly conditioned problems in potentially lower time and memory than GMRES-IR3 thanks to the use of lower precision arithmetic in the GMRES iterations. Although the combined use of five different precisions results in thousands of different variants, we provided a list of rules that we used to identify a much smaller set of practical interest; these include variants for which only an experimental evaluation was available prior to this work. Finally, we presented a rich experimental analysis of the GMRES-IR5 method on randomly generated matrices as well as matrices from real life applications. The experimental results are in good agreement with our theoretical analysis and show that GMRES-IR5 is a robust and versatile method for solving linear systems of equations.

#### REFERENCES

- [1] A. BUTTARI, J. DONGARRA, J. LANGOU, J. LANGOU, P. LUSZCZEK, AND J. KURZAK, *Mixed precision iterative refinement techniques for the solution of dense linear systems*, International Journal of High Performance Computing Applications, 21 (2007), <https://doi.org/10.1177/1094342007084026>.
- [2] 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 J. Sci. Comput., 39 (2017), pp. A2834–A2856, <https://doi.org/10.1137/17M1122918>.
- [3] E. CARSON AND N. J. HIGHAM, *Accelerating the solution of linear systems by iterative refinement in three precisions*, SIAM J. Sci. Comput., 40 (2018), pp. A817–A847, <https://doi.org/10.1137/17M1140819>.
- [4] E. CARSON, N. J. HIGHAM, AND S. PRANESH, *Three-precision GMRES-based iterative refinement for least squares problems*, SIAM J. Sci. Comput., 42 (2020), pp. A4063–A4083, <https://doi.org/10.1137/20m1316822>.
- [5] M. P. CONNOLLY, N. J. HIGHAM, AND T. MARY, *Stochastic rounding and its probabilistic backward error analysis*, SIAM J. Sci. Comput., 43 (2021), pp. A566–A585, <https://doi.org/10.1137/20m1334796>.

- [6] *cuSOLVER: Direct Linear Solvers on NVIDIA GPUs*. <https://developer.nvidia.com/cusolver>.
- [7] T. A. DAVIS AND Y. HU, *The University of Florida Sparse Matrix Collection*, ACM Trans. Math. Software, 38 (2011), pp. 1:1–1:25, <https://doi.org/10.1145/2049662.2049663>.
- [8] A. HAIDAR, A. ABDEFATTAH, M. ZOUNON, P. WU, S. PRANESH, S. TOMOV, AND J. DONGARRA, *The design of fast and energy-efficient linear solvers: On the potential of half-precision arithmetic and iterative refinement techniques*, in Computational Science—ICCS 2018, Y. Shi, H. Fu, Y. Tian, V. V. Krzhizhanovskaya, M. H. Lees, J. Dongarra, and P. M. A. Sloot, eds., Springer, Cham, Switzerland, 2018, pp. 586–600, [https://doi.org/10.1007/978-3-319-93698-7\\_45](https://doi.org/10.1007/978-3-319-93698-7_45).
- [9] A. HAIDAR, H. BAYRAKTAR, S. TOMOV, J. DONGARRA, AND N. J. HIGHAM, *Mixed-precision iterative refinement using tensor cores on GPUs to accelerate solution of linear systems*, Proc. Roy. Soc. London Ser. A, 476 (2020), p. 20200110, <https://doi.org/10.1098/rspa.2020.0110>.
- [10] A. HAIDAR, S. TOMOV, J. DONGARRA, AND N. J. HIGHAM, *Harnessing GPU tensor cores for fast FP16 arithmetic to speed up mixed-precision iterative refinement solvers*, in Proceedings of the International Conference for High Performance Computing, Networking, Storage, and Analysis, SC18 (Dallas, TX), Piscataway, NJ, USA, 2018, pp. 47:1–47:11, <https://doi.org/10.1109/SC.2018.00050>.
- [11] D. J. HIGHAM, N. J. HIGHAM, AND S. PRANESH, *Random matrices generating large growth in LU factorization with pivoting*, SIAM J. Matrix Anal. Appl., 42 (2021), p. 185–201, <https://doi.org/10.1137/20M1338149>.
- [12] N. J. HIGHAM, *Accuracy and Stability of Numerical Algorithms*, Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, second ed., 2002, <https://doi.org/10.1137/1.9780898718027>.
- [13] N. J. HIGHAM, *How fast is quadruple precision arithmetic?* <https://nhigham.com/2017/08/31/how-fast-is-quadruple-precision-arithmetic/>, June 2016.
- [14] N. J. HIGHAM, *Error analysis for standard and GMRES-based iterative refinement in two and three-precisions*, MIMS EPrint 2019.19, Manchester Institute for Mathematical Sciences, The University of Manchester, UK, Nov. 2019, <http://eprints.maths.manchester.ac.uk/2735/>.
- [15] N. J. HIGHAM AND T. MARY, *A new approach to probabilistic rounding error analysis*, SIAM J. Sci. Comput., 41 (2019), pp. A2815–A2835, <https://doi.org/10.1137/18M1226312>.
- [16] N. J. HIGHAM AND T. MARY, *Sharper probabilistic backward error analysis for basic linear algebra kernels with random data*, SIAM J. Sci. Comput., 42 (2020), pp. A3427–A3446, <https://doi.org/10.1137/20M1314355>.
- [17] N. J. HIGHAM AND S. PRANESH, *Exploiting lower precision arithmetic in solving symmetric positive definite linear systems and least squares problems*, SIAM J. Sci. Comput., 43 (2021), pp. A258–A277, <https://doi.org/10.1137/19M1298263>.
- [18] N. J. HIGHAM, S. PRANESH, AND M. ZOUNON, *Squeezing a matrix into half precision, with an application to solving linear systems*, SIAM J. Sci. Comput., 41 (2019), pp. A2536–A2551, <https://doi.org/10.1137/18M1229511>.
- [19] *Matrix algebra on GPU and multicore architectures (MAGMA)*. <http://icl.cs.utk.edu/magma/>.
- [20] C. B. MOLER, *Iterative refinement in floating point*, J. ACM, 14 (1967), p. 316–321, <https://doi.org/10.1145/321386.321394>, <https://doi.org/10.1145/321386.321394>.
- [21] T. OGITA, *Accurate matrix factorization: Inverse LU and inverse QR factorizations*, SIAM J. Matrix Anal. Appl., 31 (2010), pp. 2477–2497, <https://doi.org/10.1137/090754376>.
- [22] C. C. PAIGE, M. ROZLOŽNÍK, AND Z. STRAKOŠ, *Modified Gram-Schmidt (MGS), least squares, and backward stability of MGS-GMRES*, SIAM J. Matrix Anal. Appl., 28 (2006), pp. 264–284, <https://doi.org/10.1137/050630416>.
- [23] S. M. RUMP, *Inversion of extremely ill-conditioned matrices in floating-point*, Japan Journal of Industrial and Applied Mathematics, 26 (2009), pp. 249–277.
- [24] J. H. WILKINSON, *Rounding Errors in Algebraic Processes*, Notes on Applied Science No. 32, Her Majesty’s Stationery Office, London, 1963. Also published by Prentice-Hall, Englewood Cliffs, NJ, USA. Reprinted by Dover, New York, 1994.