

Iterative Refinement of Schur decompositions

Zvonimir Bujanović* Daniel Kressner† Christian Schröder‡

March 22, 2022

Abstract

The Schur decomposition of a square matrix A is an important intermediate step of state-of-the-art numerical algorithms for addressing eigenvalue problems, matrix functions, and matrix equations. This work is concerned with the following task: Compute a (more) *accurate* Schur decomposition of A from a given *approximate* Schur decomposition. This task arises, for example, in the context of parameter-dependent eigenvalue problems and mixed precision computations. We have developed a Newton-like algorithm that requires the solution of a triangular matrix equation and an approximate orthogonalization step in every iteration. We prove local quadratic convergence for matrices with mutually distinct eigenvalues and observe fast convergence in practice. In a mixed low-high precision environment, our algorithm essentially reduces to only four high-precision matrix-matrix multiplications per iteration. When refining double to quadruple precision, it often needs only 3–4 iterations, which reduces the time of computing a quadruple precision Schur decomposition by up to a factor of 10–20.

1 Introduction

Given a matrix $A \in \mathbb{C}^{n \times n}$, a factorization of the form

$$T = Q^H A Q, \tag{1}$$

with $Q \in \mathbb{C}^{n \times n}$ unitary and $T \in \mathbb{C}^{n \times n}$ upper triangular is called *Schur decomposition* of A . This decomposition plays a central role in algorithms for solving eigenvalue problems, computing matrix functions, and solving matrix equations. Note that the diagonal of T contains the eigenvalues of A and they can be reordered to appear there in any desirable order [18].

*University of Zagreb, Faculty of Science, Department of Mathematics, Croatia. zbujanov@math.hr. Supported by the Croatian Science Foundation under grant IP-2019-04-6268.

†Institute of Mathematics, EPFL, Switzerland. daniel.kressner@epfl.ch

‡Freelancing numerical analyst, last academic position: Institut für Mathematik, TU Berlin, Germany. chris.schroeder@gmx.net

In this work, we consider the following question. Suppose that we have an approximate Schur decomposition at our disposal, that is, T is *nearly* upper triangular and Q is unitary or only *nearly* unitary. How do we refine Q, T to yield a more accurate Schur decomposition using an algorithm that is computationally and conceptually less demanding than computing the Schur decomposition of A from scratch by applying, e.g., the QR algorithm [11, 12]?

Partly motivated by the increased role of GPUs and TPUs in high-performance computing, there has been a revival of interest in exploiting the benefits of a mixed-precision environment in numerical computations; see [2] for a recent survey. Now, suppose that a Schur decomposition of A has been computed, inexpensively, in a certain (low) machine precision but the target application demands for higher accuracy. For the computed factor \hat{Q} , the entries of $\hat{Q}^H A \hat{Q}$ below the diagonal and the entries of $\hat{Q}^H \hat{Q} - I$ are roughly on the level of unit roundoff in low precision. The refinement procedure discussed in this work produces a Schur decomposition that is accurate in high machine precision essentially at the cost of a few matrix multiplications in high machine precision. Examples for this setting of current interest include mixing half or single with double precision on specialized processors as well as mixing double with quadruple or even higher precision on standard CPUs. In both cases, the operations performed in high precision are significantly more costly and should therefore be limited to the minimum. Other scenarios that could benefit from the refinement of Schur decompositions include the solution of parameter-dependent eigenvalue problems and continuation methods [8, 9].

The goal of this paper is to develop an efficient Newton-like method for refining a Schur decomposition. Closely related to the theme of this paper, the refinement of an individual eigenvector has been investigated quite thoroughly, going back to the works of Wielandt [41] and Wilkinson [42], and is usually addressed by some form of inverse iteration, see also [24]. When several eigenvalues and/or eigenvectors are of interest, it is sensible to refine the invariant subspace (belonging to the eigenvalues of interest) as a whole rather than individual eigenvectors. Various variants of the Newton algorithm have been proposed for this purpose, see, e.g., [5, 8, 9, 13, 15, 31]. Often, these algorithms require the solution of a Sylvester equation at each iteration. While the developments presented in this work bear similarities, we are not aware that any of these existing methods would allow for refining a Schur decomposition or, equivalently, for simultaneously refining the entire flag of invariant subspaces associated with the Schur decomposition. In principle, Jacobi algorithms could be used for this purpose; see [16, 19, 23, 38] for examples. While such algorithms converge locally and asymptotically quadratically under certain conditions [32], their efficient implementation requires significant attention to low-level implementation aspects while the critical parts of our algorithm are entirely based on matrix multiplications. Improving an earlier method by Davies and Modi [14], Ogita and Aishima [34, 35] recently proposed a Newton-type methods for refining spectral decompositions, addressing the task considered in this work when A is real symmetric/complex Hermitian matrix. We will discuss similarities and differences with our algorithm in Section 2.5. As a cu-

riosity, we also point out a discussion by Kahan [28] on refining eigendecompositions for nonsymmetric matrices; see also Jahn's 1948 work [25].

2 Algorithms

To motivate the approach pursued in this work, let us consider a unitary matrix \hat{Q} that nearly effects a Schur decomposition:

$$\hat{Q}^H A \hat{Q} = T + E, \quad \varepsilon := \|E\|_F \approx 0, \quad (2)$$

where T is upper triangular and E is strictly lower triangular. We will write $\text{stril}(\cdot)$ to denote the strictly lower triangular part of a matrix. A unitary matrix Q_L that transforms $T + E$ to Schur form satisfies the two matrix equations

$$\text{stril}(Q_L^H(T + E)Q_L) = 0, \quad Q_L^H Q_L = I. \quad (3)$$

In the following we linearize these equations, analogous to existing first-order perturbation analyses of Schur decompositions [29, 39]. Writing $Q_L = I + W$, the condition $Q_L^H Q_L = I$ becomes equivalent to

$$W + W^H + W^H W = 0 \quad (4)$$

Ignoring the second order term $W^H W$, this means that W is skew-Hermitian and can be written as

$$W = L + D - L^H, \quad \text{with } L = \text{stril}(W), \quad (5)$$

where D is diagonal with purely imaginary diagonal elements. Now, the first equation in (3) becomes

$$\text{stril}((I + L + D - L^H)^H(T + E)(I + L + D - L^H)) = 0.$$

Assuming $\|W\|_F = O(\varepsilon)$, where $\|\cdot\|_F$ denotes the Frobenius norm, and dropping all second order terms in ε , we arrive at the triangular matrix equation

$$\text{stril}(E - LT + TL) = 0. \quad (6)$$

If the eigenvalues of T are pairwise distinct, there is a unique strictly lower triangular matrix L satisfying (6), see [29, 39] and also Theorem 1 below. Note that the diagonal factor D has disappeared in (6); it can be chosen to contain arbitrary diagonal entries on the imaginary axis. In the following, we choose $D = 0$.

In order to attain a unitary factor Q_L , the matrix $I + W = I + L - L^H$ needs to be orthogonalized; we will discuss different orthogonalization strategies in Section 2.2.

Algorithm 1 summarizes the procedure outlined above. Note that we allow the input matrix \hat{Q} to be non-unitary, which is taken care of by an optional orthogonalization step in line 1.

Algorithm 1 Template for refining a Schur decomposition

Input: Matrix $A \in \mathbb{C}^{n \times n}$ and (nearly) unitary matrix $\hat{Q} \in \mathbb{C}^{n \times n}$ such that $\hat{Q}^H A \hat{Q}$ is nearly upper triangular.

Output: Unitary matrix Q such that $Q^H A Q$ is upper triangular.

- 1: Compute (approximately) unitary matrix Q from orthogonalizing \hat{Q} .
- 2: **repeat**
- 3: Compute $\hat{T} \leftarrow Q^H A Q$.
- 4: Set $E \leftarrow \text{stril}(\hat{T})$, $T \leftarrow \hat{T} - E$.
- 5: Solve equation (6) for strictly lower triangular L .
- 6: Replace Q by (approximately) unitary matrix obtained from orthogonalizing $Q(I + L - L^H)$.
- 7: **until** convergence

In the following sections, we will provide details on Algorithm 1 and perform a convergence analysis.

2.1 Solution of triangular matrix equation

Let us first consider the triangular matrix equation from Step 5 of Algorithm 1:

$$\text{stril}(TL - LT) = -E, \quad (7)$$

where T is upper triangular and E as well as the desired solution L are strictly lower triangular. For $1 \leq j < i \leq n$, equating the (i, j) -entry of (7) yields

$$\sum_{k=i}^n t_{ik} \ell_{ki} - \sum_{k=1}^j \ell_{ik} t_{kj} = -e_{ij}, \quad (8)$$

which, assuming $t_{ii} \neq t_{jj}$, can be rewritten as

$$\ell_{ij} = -\frac{1}{t_{ii} - t_{jj}} \left(e_{ij} + \sum_{k=i+1}^n t_{ik} \ell_{ki} - \sum_{k=1}^{j-1} \ell_{i,k} t_{kj} \right). \quad (9)$$

Note that all entries of L appearing on the right-hand side are located either to the left or below the entry (i, j) in the matrix L . Thus, if $(i_1, j_1), (i_2, j_2), \dots, (i_N, j_N)$ with $N = n(n-1)/2$ describes an order in which the entries of L shall be determined, this order must satisfy

$$\nu < \mu \implies i_\nu > i_\mu \text{ or } j_\nu < j_\mu \quad (10)$$

for all $\nu, \mu \in \{1, \dots, N\}$. Interestingly, this corresponds to the elimination order of the northeast directed sweep sequences in the nonsymmetric Jacobi algorithm for which local quadratic convergence was proven in [32]. Each order satisfying (10) gives rise to a different

algorithm for solving (6). A possible choice is a bottom-to-top columnwise order, leading to Algorithm 2. Note that we use Matlab notation to refer to entries and submatrices of a matrix.

Algorithm 2 Successive substitution for solution of triangular matrix equation (7)

Input: Upper triangular matrix $T \in \mathbb{C}^{n \times n}$ with pairwise distinct diagonal entries. Strictly lower triangular matrix $E \in \mathbb{C}^{n \times n}$.

Output: Strictly lower triangular matrix $L \in \mathbb{C}^{n \times n}$ satisfying (7).

```

1: for  $j = 1 : n - 1$  do
2:   for  $i = n : -1 : j + 1$  do
3:      $E(i, j) \leftarrow E(i, j) + T(i, i + 1 : n) \cdot L(i + 1 : n, j) - L(i, 1 : j - 1) \cdot T(1 : j - 1, j)$ 
4:      $L(i, j) \leftarrow -E(i, j) / (T(i, i) - T(j, j))$ 
5:   end for
6: end for

```

Theorem 1 *There exists a unique strictly lower triangular matrix L satisfying (7) if and only if the diagonal entries of T are pairwise distinct.*

Proof. This result can be found (implicitly) in [29, Sec. 3], [39], and [40, Theorem 3.1]. For completeness, we provide an explicit proof.

By construction, Algorithm 2 shows that (7) has a solution for every E under stated condition on T . The unique solvability simply follows from the fact that (7) can be regarded as a linear system of $n(n - 1)/2$ equations in $n(n - 1)/2$ unknowns.

For the other direction, assume that the condition is violated, that is, T has two identical eigenvalues λ . Then there is a principal submatrix \tilde{T} of T taking the form

$$\tilde{T} = \begin{bmatrix} \lambda & T_{12} & T_{13} \\ 0 & T_{22} & T_{23} \\ 0 & 0 & \lambda \end{bmatrix},$$

where the eigenvalues of T_{22} are pairwise distinct and different from λ or T_{22} vanishes. We consider

$$\tilde{L} = \begin{bmatrix} 0 & 0 & 0 \\ L_{21} & L_{22} & 0 \\ 1 & L_{32} & 0 \end{bmatrix},$$

partitioned conformally with \tilde{T} . Then $\text{stril}(\tilde{T}\tilde{L} - \tilde{L}\tilde{T}) = 0$ if and only if

$$(T_{22} - \lambda I)L_{21} = -T_{23} \tag{11}$$

$$L_{32}(\lambda I - T_{22}) = T_{12} \tag{12}$$

$$\text{stril}(T_{22}L_{22} - L_{22}T_{22}) = \text{stril}(-T_{23}L_{32} + L_{21}T_{12}) \tag{13}$$

Since $\lambda \notin \Lambda(T_{22})$, equations (11) and (12) have unique solutions L_{21} and L_{32} . This determines the right-hand side of (13), which is an equation of the same type as (7).

Since T_{22} has pairwise distinct eigenvalues, the first part of this theorem implies that there is a strictly lower triangular solution L_{22} to (13). By embedding $L = \text{diag}(0, \tilde{L}, 0)$ we have thus found a nonzero $L \in \text{stril}(\mathbb{C}^{n \times n})$ such that $\text{stril}(TL - LT) = 0$. Because of linearity, this implies that the equation $\text{stril}(TL - LT) = -E$ is not uniquely solvable for any $E \in \text{stril}(\mathbb{C}^{n \times n})$ if the condition is violated. \square

The non-locality of the memory access pattern renders an actual implementation of Algorithm 2 slow for larger matrices. Other orderings satisfying (10), like parallel wave-front techniques [36], may lead to more efficient implementations. In the following we use a recursive formulation, as proposed for a variety of matrix equations in [26, 27], to achieve increased data locality in a convenient manner. For this purpose, let us partition

$$T = \begin{bmatrix} T_{11} & T_{12} \\ 0 & T_{22} \end{bmatrix}, \quad E = \begin{bmatrix} E_{11} & 0 \\ E_{21} & E_{22} \end{bmatrix}, \quad L = \begin{bmatrix} L_{11} & 0 \\ L_{21} & L_{22} \end{bmatrix},$$

with $T_{11}, E_{11}, L_{11} \in \mathbb{C}^{n_1 \times n_1}$. Inserted into (7), this gives

$$\begin{aligned} & \begin{bmatrix} \text{stril}(T_{11}L_{11} + T_{12}L_{21} - L_{11}T_{11}) & 0 \\ T_{22}L_{21} - L_{21}T_{11} & \text{stril}(T_{22}L_{22} - L_{21}T_{12} - L_{22}T_{22}) \end{bmatrix} \\ = & - \begin{bmatrix} E_{11} & 0 \\ E_{21} & E_{22} \end{bmatrix}. \end{aligned} \quad (14)$$

The (2, 1) entry is a triangular Sylvester equation,

$$T_{22}L_{21} - L_{21}T_{11} = -E_{21}, \quad (15)$$

which has a unique solution L_{21} , provided that T_{11} and T_{22} have no eigenvalue in common, and can be addressed using, e.g., the software package RECSY [26, 27]. Once L_{21} is determined, L_{11} and L_{22} can be obtained as solutions of

$$\begin{aligned} \text{stril}(T_{11}L_{11} - L_{11}T_{11}) &= -(E_{11} + \text{stril}(T_{12}L_{21})), \\ \text{stril}(T_{22}L_{22} - L_{22}T_{22}) &= -(E_{22} - \text{stril}(L_{21}T_{12})). \end{aligned}$$

Note that both equations are of the same type as (6), but of size n_1 and $n - n_1$, respectively. Their solutions can be obtained by again subdividing into 2×2 blocks or, for smaller problems, by using Algorithm 2. The described strategy leads to the following recursive algorithm for solving (6).

Algorithm 3 Recursive block algorithm for triangular matrix equation (7)

Input: Upper triangular matrix $T \in \mathbb{C}^{n \times n}$ with pairwise distinct diagonal entries. Strictly lower triangular matrix $E \in \mathbb{C}^{n \times n}$. Integer $n_{\min} \geq 2$.

Output: Strictly lower triangular matrix $L \in \mathbb{C}^{n \times n}$ satisfying (7).

- 1: **if** $n \leq n_{\min}$ **then**
- 2: Apply Algorithm 2.

3: **else**

4: Choose n_1 with $1 < n_1 < n$ (e.g., $n_1 = \lfloor n/2 \rfloor$) and define $\text{ind}_1 = 1 : n_1$, $\text{ind}_2 = n_1 + 1 : n$.

5: Solve Sylvester equation

$$T(\text{ind}_2, \text{ind}_2)X - XT(\text{ind}_1, \text{ind}_1) = -E(\text{ind}_2, \text{ind}_1).$$

6: Set $L(\text{ind}_2, \text{ind}_1) \leftarrow X$.

7: Update $E(\text{ind}_1, \text{ind}_1) \leftarrow E(\text{ind}_1, \text{ind}_1) + \text{stril}(T(\text{ind}_1, \text{ind}_2)L(\text{ind}_2, \text{ind}_1))$.

8: Update $E(\text{ind}_2, \text{ind}_2) \leftarrow E(\text{ind}_2, \text{ind}_2) - \text{stril}(L(\text{ind}_2, \text{ind}_1)T(\text{ind}_1, \text{ind}_2))$.

9: Use Algorithm 3 with input $T(\text{ind}_1, \text{ind}_1)$, $E(\text{ind}_1, \text{ind}_1)$ to compute $L(\text{ind}_1, \text{ind}_1)$.

10: Use Algorithm 3 with input $T(\text{ind}_2, \text{ind}_2)$, $E(\text{ind}_2, \text{ind}_2)$ to compute $L(\text{ind}_2, \text{ind}_2)$.

11: **end if**

The main cost of Algorithm 3 is in the two matrix multiplications in Lines 7 and 8, which allows to leverage the efficiency of level 3 BLAS operations.

2.2 Orthogonalization procedures

In the following, we discuss algorithms for carrying out the (approximate) orthogonalization in lines 1 and 6 of Algorithm 1. A suitable orthogonalization procedure needs to satisfy two goals. On the one hand, it should improve orthogonality. On the other hand, it should not modify the matrix more than needed. With these two goals in mind, we require the following: For $\varepsilon > 0$, consider any matrices $Q, W \in \mathbb{C}^{n \times n}$ such that W is skew-Hermitian and

$$\|Q^H Q - I\|_F \leq c_Q \varepsilon^2, \quad \|W\|_F \leq c_W \varepsilon, \quad (16)$$

with constants c_Q, c_W independent of ε . Then the matrix Q_{new} returned by the orthogonalization procedure applied to $Q(I + W)$ satisfies

$$\|Q_{\text{new}}^H Q_{\text{new}} - I\|_F = O(\varepsilon^4), \quad \|Q_{\text{new}} - Q(I + W)\|_F = O(\varepsilon^2). \quad (17)$$

2.2.1 Orthogonalization by QR decomposition

When Q is unitary (that is, (16) is satisfied with $c_Q = 0$) then $Q(I + W)$ is an element of the tangent space of the manifold of unitary matrices at Q . Retractions, a concept popularized in the context of Riemannian optimization [4, 10], map an element from the tangent space back to the manifold. Thus, the first relation in (17) is trivially satisfied by a retraction. The second relation in (17) follows from the fact that retractions approximate the exponential map in first order. Suitable retractions and, hence, suitable orthogonalization procedures are obtained from extracting the orthogonal factors of a QR or polar decomposition [10, Sec. 7.3]. More specifically, if one computes a QR decomposition

$$I + W = Q_L R,$$

such that Q_L is unitary and the upper triangular matrix R has real and positive entries on its diagonal then $Q_{\text{new}} = QQ_L$ satisfies (17).

2.2.2 Orthogonalization by Newton-Schulz iteration

The Newton-Schulz iteration $Q_{k+1} = \frac{1}{2}Q_k(3I - Q_k^H Q_k)$ converges quadratically towards the unitary factor of the polar decomposition of Q_0 if Q_0 is invertible and $\|Q_0\|_2 < \sqrt{3}$ [21, Ch. 8]. The following lemma shows that one step of the Newton-Schulz iteration yields a suitable orthogonalization procedure.

Lemma 2 *Consider $\hat{Q} = Q(I + W)$ satisfying (16). Then*

$$Q_{\text{new}} = \frac{1}{2}\hat{Q}(3I - \hat{Q}^H \hat{Q})$$

satisfies (17).

Proof. The first relation of (17) follows from $Q_{\text{new}}^H Q_{\text{new}} - I = -\frac{3}{4}(\hat{Q}^H \hat{Q} - I)^2 + \frac{1}{4}(\hat{Q}^H \hat{Q} - I)^3$; see [21, Problem 8.20]. To show the second relation, we set $\Delta := Q^H Q - I$ and compute

$$\hat{Q}^H \hat{Q} - I = (I - W)(I + \Delta)(I + W) - I = \Delta - W^2 - W\Delta + \Delta W - W\Delta W.$$

Combined with $Q_{\text{new}} = \hat{Q} - \frac{1}{2}\hat{Q}(\hat{Q}^H \hat{Q} - I)$, this shows

$$\|Q_{\text{new}} - \hat{Q}\|_F \leq \frac{c_Q + c_W^2}{2}\varepsilon^2 + O(\varepsilon^3).$$

□

2.3 Local convergence analysis

We now perform a local convergence analysis of Algorithm 1. This analysis makes use of the quantity

$$\phi(T) := \min \{ \|\text{stril}(TL - LT)\|_F : L \in \text{stril}(\mathbb{C}^{n \times n}), \|L\|_F = 1 \} \quad (18)$$

for an upper triangular matrix T , where $\text{stril}(\mathbb{C}^{n \times n})$ denotes the set of all strictly lower triangular $n \times n$ matrices. Note that $1/\phi(T)$ governs the first-order sensitivity of a Schur decomposition [29, 39]. By Theorem 1, $\phi(T) > 0$ if and only if the diagonal entries of T are pairwise distinct. The following lemma shows that ϕ is Lipschitz continuous with constant $\sqrt{2}$; see [40, Theorem 3.2] for a related result.

Lemma 3 *Consider upper triangular matrices T, \hat{T} , each having pairwise distinct diagonal entries. Then $|\phi(T) - \phi(\hat{T})| \leq 2\|T - \hat{T}\|_2$.*

Proof. By the definition (18), the quantity $\phi(T)$ is the smallest singular value of the linear operator $\mathcal{L}_T : \text{stril}(\mathbb{C}^{n \times n}) \rightarrow \text{stril}(\mathbb{C}^{n \times n})$, $\mathcal{L}_T(L) : L \mapsto \text{stril}(TL - LT)$, with the norm $\|\mathcal{L}_T\| = \sup\{\|\mathcal{L}_T(L)\|_F : \|L\|_F = 1, L \in \text{stril}(\mathbb{C}^{n \times n})\}$. From

$$\begin{aligned} \|(\mathcal{L}_T - \mathcal{L}_{\hat{T}})(L)\|_F &= \|\text{stril}((T - \hat{T})L - L(T - \hat{T}))\|_F \\ &\leq \|(T - \hat{T})L - L(T - \hat{T})\|_F \leq 2\|T - \hat{T}\|_2 \|L\|_F \end{aligned}$$

it follows that $\|\mathcal{L}_T - \mathcal{L}_{\hat{T}}\| \leq 2\|T - \hat{T}\|_2$. Since Weyl's inequalities [22, Corollary 7.3.5] imply that singular values are Lipschitz continuous with constant 1, this concludes the proof. \square

Suppose that A has mutually distinct eigenvalues. Then its Schur decomposition $A = QTQ^H$ is unique up to the order of the eigenvalues on the diagonal of T and unitary diagonal transformations. While $\phi(T)$ is invariant under the latter, it does depend on the eigenvalue order. To circumvent this difficulty, we introduce

$$\psi(A) := \min\{\phi(T) : A = QTQ^H \text{ is a Schur decomposition}\}. \quad (19)$$

Because there are only finitely many different eigenvalue orderings, the minimum is assumed and positive. The following lemma shows that $\psi(\cdot)$ remains positive in a neighborhood of A .

Lemma 4 *Let $A \in \mathbb{C}^{n \times n}$ have mutually distinct eigenvalues and consider $\varepsilon > 0$ sufficiently small. Then for all $\tilde{A} \in \mathbb{C}^{n \times n}$ with $\|\tilde{A} - A\|_F \leq \varepsilon$, it holds that $\psi(\tilde{A}) \geq \psi(A) - 2(1 + 4\|A\|_2/\psi(A))\varepsilon$.*

Proof. For fixed $\Delta \in \mathbb{C}^{n \times n}$ with $\|\Delta\|_F = 1$ and $\varepsilon > 0$, set $A(\varepsilon) := A + \varepsilon\Delta$ and consider a Schur decomposition $A(\varepsilon) = Q(\varepsilon)T(\varepsilon)Q(\varepsilon)^H$. Let $\lambda_1(\varepsilon), \dots, \lambda_n(\varepsilon)$ denote the diagonal elements of $T(\varepsilon)$. Because of continuity of eigenvalues, we may assume w.l.o.g. that each $\lambda_i(\cdot)$ is a continuous function.

For sufficiently small $\varepsilon > 0$, the perturbation results from [29, 39, 40] imply that there is a Schur decomposition $A(\varepsilon) = \tilde{Q}(\varepsilon)\tilde{T}(\varepsilon)\tilde{Q}(\varepsilon)^H$ with $\tilde{T}(\varepsilon) \approx T(0)$. Specifically, the computation of the condition number of T on Page 391 of [29] shows that

$$\|\tilde{T}(\varepsilon) - T(0)\|_F \leq \left(1 + 2\sqrt{2} \frac{\|T(0)\|_2}{\phi(T(0))}\right)\varepsilon + O(\varepsilon^2) \leq \left(1 + 4 \frac{\|T(0)\|_2}{\phi(T(0))}\right)\varepsilon.$$

By Lemma 3, it follows that

$$\phi(\tilde{T}(\varepsilon)) \geq \phi(T(0)) - 2\left(1 + 4 \frac{\|T(0)\|_2}{\phi(T(0))}\right)\varepsilon \geq \psi(A) - 2\left(1 + 4 \frac{\|A\|_2}{\psi(A)}\right)\varepsilon. \quad (20)$$

The perturbation bound on $\tilde{T}(\varepsilon)$ also implies that the eigenvalues $\lambda_1(\varepsilon), \dots, \lambda_n(\varepsilon)$ need to appear in this order on the diagonal of $\tilde{T}(\varepsilon)$ for sufficiently small $\varepsilon > 0$. Thus, $\tilde{T}(\varepsilon)$ and $T(\varepsilon)$ only differ by a transformation with a unitary diagonal matrix and therefore $\phi(\tilde{T}(\varepsilon)) = \phi(T(\varepsilon))$. Combined with (20), this completes the proof. \square

Theorem 5 *Assume that the eigenvalues of $A \in \mathbb{C}^{n \times n}$ are pairwise distinct. Given $\varepsilon > 0$ consider any matrix $Q \in \mathbb{C}^{n \times n}$ satisfying*

1. $\|\text{stril}(Q^H A Q)\|_F \leq \varepsilon$;
2. $\|Q^H Q - I\|_F \leq \varepsilon^2$.

Let Q_{new} denote the output of one iteration of Algorithm 1 applied to A, Q using an orthogonalization procedure that satisfies (17). Then, for $\varepsilon > 0$ sufficiently small, it holds that

$$\|\text{stril}(Q_{\text{new}}^H A Q_{\text{new}})\|_F = O(\varepsilon^2), \quad \|Q_{\text{new}}^H Q_{\text{new}} - I\|_F = O(\varepsilon^4). \quad (21)$$

Proof. Let $E = \text{stril}(Q^H A Q)$ and $T = Q^H A Q - E$. By the definition of $\phi(T)$, it follows that the solution L of (6) satisfies

$$\|L\|_F \leq \frac{\|E\|_F}{\phi(T)} \leq \frac{\varepsilon}{\phi(T)}.$$

In order to proceed from here, we need to show that $\phi(T)$, which depends on Q , admits a uniform lower bound. By the polar decomposition there is a unitary matrix \tilde{Q} such that $\|\tilde{Q} - Q\|_F \leq \|Q^H Q - I\|_F \leq \varepsilon^2$; see, e.g., [20, P. 380]. In turn, we obtain the perturbed Schur decomposition

$$\begin{aligned} \tilde{Q}^H A \tilde{Q} &= Q^H A Q + \tilde{Q}^H A (\tilde{Q} - Q) + (\tilde{Q} - Q)^H A Q \\ &= T + E + \Delta, \quad \|\Delta\|_F \leq 3\|A\|_2 \varepsilon^2, \end{aligned}$$

where we used that $\|Q\|_2 = \sqrt{\|Q^H Q\|_2} \leq \sqrt{1 + \varepsilon^2} \leq 2$ for ε sufficiently small. Equivalently, this can be viewed as a Schur decomposition of a perturbed matrix: $A - \tilde{Q}(E + \Delta)\tilde{Q}^H = \tilde{Q}T\tilde{Q}^H$. From Lemma 4 it follows that $\phi(T) \geq \psi(A)/2$ holds for ε sufficiently small. Hence, $W = L - L^H$ satisfies $\|W\|_F \leq c_W \varepsilon$ with $c_W = 4/\psi(A)$. This means that the matrix $\hat{Q} = Q(I + W)$ satisfies the conditions (16) and therefore the matrix Q_{new} returned by the orthogonalization procedure in line 6 of Algorithm 1 satisfies (17). This already establishes the second relation in (21). To establish the first relation, we note that

$$\begin{aligned} \hat{Q}^H A \hat{Q} &= (I + W)^H (T + E) (I + W) \\ &= \underbrace{E - LT + TL}_{\text{upper triangular}} + \underbrace{T + L^H T - TL^H}_{\text{upper triangular}} + W^H (T + E) W + W^H E + EW, \end{aligned}$$

and, hence,

$$\begin{aligned} \|\text{stril}(\hat{Q}^H A \hat{Q})\|_F &\leq \|T + E\|_F \|W\|_F^2 + 2\|W\|_F \|E\|_F \\ &\leq \|Q\|_2^2 \|A\|_F c_W^2 \varepsilon^2 + 2c_w \varepsilon^2 \leq (4\|A\|_F c_W^2 + 2c_w) \varepsilon^2. \end{aligned}$$

Combined with $\|Q_{\text{new}} - \hat{Q}\| = O(\varepsilon^2)$, this proves the first relation in (21). \square

If Algorithm 1 is started with a unitary matrix \hat{Q} sufficiently close to a unitary matrix Q that transforms A to Schur form then the relations (17) imply that the matrix Q obtained from the orthogonalization procedure applied in line 1 satisfies the conditions of Theorem 5. Hence, Theorem 5 establishes local quadratic convergence of Algorithm 1.

Remark 6 *The idea of imposing constraints (such as unitary matrix structure) only asymptotically upon convergence of an iterative procedure is not new. Similar ideas have been proposed in the context of differential-algebraic equations (Baumgarte's method [6]), constrained optimization (interior point method [33]), and – more recently – Riemannian optimization ([17], landing algorithm [3]). However, Algorithm 1 does not seem to fit into any of these existing developments.*

2.4 Complete algorithm for mixed precision

In this section we specialize the template Algorithm 1 to computing a Schur decomposition of a given matrix A in a mixed precision environment. The decomposition is first computed in a lower precision (lp in the following), and then refined to a higher precision (hp in the following). A typical scenario uses double precision as lp, and quadruple or 100-digits precision as hp. The initial lp decomposition produces the matrices \hat{Q} and T_{lp} such that $A \approx \hat{Q}T_{lp}\hat{Q}^H$, and the matrices A and \hat{Q} are the input to the refinement algorithm. In order to ensure convergence (see Theorem 5 and the discussion in the previous section), orthogonality of the matrix \hat{Q} is first improved by applying one step of the Newton-Schulz iteration:

$$Q = \frac{1}{2}\hat{Q}(3I - \hat{Q}^H\hat{Q}). \quad (22)$$

This computation is done entirely in hp; the matrix \hat{Q} is first converted to hp. The most time consuming parts of (22) are the two matrix-matrix multiplications in hp.

We then enter the loop in Algorithm 1: computation of $\hat{T} = Q^H A Q$ in Line 3 is also done in hp, requiring two additional matrix-matrix multiplications. Equation (6) in Line 5 can be solved entirely in lp. Note that numerical instabilities in computing L are expected when there are clustered eigenvalues in T . In hope that these instabilities do not spread through the entire matrix L , the initial lp Schur decomposition is reordered so that clustered eigenvalues appear on neighboring diagonal elements of T_{lp} . This is done using `ordschur` in Matlab; the target permutation of eigenvalues is determined by the order in which they appear after being projected on a random line in the complex plane.

Finally, we need to compute the matrix Q_{new} by improving the orthogonality of $Q(I + L - L^H)$. In line with Lemma 2, this is done by applying one step of the Newton-Schulz iteration to $Q(I + L - L^H)$. To lower computational complexity and avoid unnecessary matrix-matrix multiplications in hp, we proceed in the following way: let $W = L - L^H$,

and $Y = Q^H Q - I$. Then (22) with $\hat{Q} = Q(I + W)$ becomes

$$\begin{aligned} Q_{\text{new}} &= \frac{1}{2}Q(I + W)(3I - (I + W)^H Q^H Q(I + W)) \\ &= \frac{1}{2}Q(2I + 2W - Y - YW + W^2 + W^3 + W^2Y + W^2YW), \end{aligned} \quad (23)$$

where we used $W^H = -W$. We observe that the products between the matrices Y and W can all be computed in lp while the summation of the terms needs to be carried out in hp. In any case, only 2 hp matrix-matrix multiplications are needed to compute the update: one to compute Y , and one to compute the product of Q with $(2I + 2W - Y - \dots)$. The cost can be reduced slightly by observing that high-order terms tend to become very small and may even fall below the unit roundoff in lp. In our implementation, we actually use the approximation $Q_{\text{new}} \approx \frac{1}{2}Q(2I + 2W - Y - YW + W^2 + W^3)$.

The complete procedure is summarized in Algorithm 4. As the numerical experiments will demonstrate, the computational complexity is almost entirely concentrated in hp matrix multiplication. The algorithm does 2 such multiplications before the loop, and then 4 in each loop iteration. The iteration is stopped in line 5 when E becomes small enough. Note that lines 9–11 can be skipped if the norms of Y and W indicate that $Q(I + W)$ is already unitary in hp. For this purpose, the norm of $\|Y\|_F \sim \mathbf{u}_{\text{hp}}$, where \mathbf{u}_{hp} denotes unit roundoff in hp, while $\|W\|_F \sim \sqrt{\mathbf{u}_{\text{hp}}}$ due to the fact that W is skew-Hermitian (see (4)). This potentially saves 1 hp matrix multiplication in the penultimate iteration.

Algorithm 4 Mixed-precision computation of Schur decomposition

Input: Matrix $A \in \mathbb{C}^{n \times n}$ in hp.

Output: Matrix $Q \in \mathbb{C}^{n \times n}$ such that Q is unitary in hp and $Q^H A Q$ is upper triangular in hp.

- 1: Compute (ordered) Schur decomposition $A \approx \hat{Q} T_{\text{lp}} \hat{Q}^H$ in lp.
- 2: Convert \hat{Q} to hp and update $Q \leftarrow \frac{1}{2} \hat{Q} (3I - \hat{Q}^H \hat{Q})$ in hp (one step of Newton-Schulz).

3: **repeat**

4: Compute $\hat{T} \leftarrow Q^H A Q$ in hp.

5: Set $E \leftarrow \text{stril}(\hat{T})$, $T \leftarrow \hat{T} - E$ and convert to lp.

6: Solve equation (6) for strictly lower triangular L in lp, using Algorithm 3.

7: Convert L to hp.

8: Set $W = L - L^H$, compute $Y = Q^H Q - I$ in hp.

9: Compute products YW , W^2 , W^3 , W^2Y , W^2YW in lp.

10: Compute $\Sigma = 2I + 2W - Y - YW + W^2 + W^3$ in hp.

11: Update $Q \leftarrow \frac{1}{2} Q \Sigma$ in hp (one step of Newton-Schulz).

12: **until** convergence

2.5 Symmetric A

When A is real and symmetric, its Schur form becomes diagonal and our algorithms can be simplified significantly. Instead of Steps 4 of Algorithms 1 and 4, \hat{T} is now decomposed in its diagonal part T and its off-diagonal part E . In turn, the solution L of the linear system (6) becomes nearly trivial; its entries are given by $l_{ij} = e_{ij}/(t_{ii} - t_{jj})$ for $i > j$. The resulting simplified algorithms bear close similarity with the method RefSyEv proposed by Ogita and Aishima [34]. However, in contrast to our approach, RefSyEv performs the improvement of diagonality and orthogonality in reverse order and integrates them in a single update. More concretely, using our notation, one iteration of RefSyEv updates $Q \leftarrow Q(I + W)$ in the following manner. It first computes $Y = Q^T Q - I$ and $\hat{T} = Q^T A Q$. Orthogonality is improved by setting the symmetric part of W to $-Y/2$, which is equivalent to one step of Newton-Schulz as pointed out in [34, Appendix A]. Then the diagonal of \hat{T} is updated and the skew-symmetric part of W is chosen to improve diagonality by solving an equation of the form (6), with the right-hand side updated to reflect improvement of orthogonality. RefSyEv actually does not treat the symmetric and skew-symmetric parts separately but derives a simple and elegant formula to update all entries of W in a single step. The most significant advantage of RefSyEv is that it avoids the initial orthogonalization in Algorithms 1 and 4. This saves 2 hp matrix-matrix multiplications in the initial phase; note, however, that both RefSyEv and Algorithm 4 need the same number (4) of hp matrix-matrix multiplications in each iteration. In summary, there seems to be little that speaks in favor of our approach for symmetric eigenvalue problems, especially when taking into account that RefSyEv and its further development described in [35] take precautions for (nearly) multiple eigenvalues in order to still attain fast convergence in such a critical case. One (small) advantage of our approach is that the separation of the orthogonalization step allows for the use of other orthogonalization procedures, which may be of interest for future developments.

3 Numerical experiments

In this section, we report the results of a number of numerical experiments in order to demonstrate the correctness and effectiveness of the proposed algorithm. The experiments were executed on a notebook computer with Intel Core i5-1135G7 CPU and 24GB RAM, running Ubuntu 21.10. Algorithm 4 was implemented in Matlab R2021b. We have also implemented a straightforward extension of this algorithm to real Schur decompositions [18]. The main difference to the complex case is that some attention needs to be paid to 2×2 diagonal blocks containing complex pairs of eigenvalues. In Algorithm 3, the value of n_{\min} was set to 4 in the complex case, and to 1 in the real case. To solve Sylvester equations of the form (15) we used the internal Matlab solver `matlab.internal.math.sylvester_tri`.

Our test procedure starts by generating an input matrix A in high precision, for which we discuss two scenarios: quadruple precision (i.e., 34 decimal digits of precision), and 100

decimal digits of precision. The matrix is given as input to Algorithm 4. There we use the standard double precision as low precision, and the builtin Matlab command `schur` to compute the initial lp Schur decomposition (either real or complex). The refined factors Q, T returned by Algorithm 4 are verified for correctness: we check if T has proper (quasi) triangular form, analyze orthogonality of Q by computing $\|I - Q^H Q\|_F$, and compute the residual $\|T - Q^H A Q\|_F$.

The performance of our algorithm heavily relies on the efficiency of matrix-matrix multiplication in the target precision. For that purpose we use the Matlab toolbox `acc` based on the Ozaki scheme [37], which was also used in [34]. The toolbox uses a configurable number of standard doubles to represent a single high-precision number; we use 2 doubles to represent a quadruple precision number, and 6 doubles to represent a 100-digit number.

The time needed by the whole refinement procedure (including the initial Schur decomposition in double precision) is compared with the time needed for computing the Schur decomposition immediately in quadruple/100 digits precision. For the latter we use Advanpix [1], a multiprecision computing toolbox for Matlab.¹

Example 7 *This first example aims at verifying the accuracy of our algorithm by applying it to the companion matrix for the Wilkinson polynomial*

$$p_{20}(x) = \prod_{i=1}^{20} (x - i).$$

The coefficients of the polynomial and, hence, entries of the companion matrix vary wildly in magnitude and cannot be stored exactly in double precision. As a consequence, computing the eigenvalues via the Matlab command `roots(poly(1:20))` incurs a large error; see Table 1.

Storing the companion matrix in quadruple precision allows for more accurate eigenvalues. Table 1 compares the accuracy of the results obtained from applying Advanpix's `schur` (with 34 digits) with the ones from Algorithm 4 using mixed double-quadruple precision. In both cases the obtained accuracy is on a similar level and significantly improves upon the double precision computation. However, it can also be seen that the use of the `acc` toolbox (which uses 2 doubles representing each number) results in a slight loss of accuracy compared to using Advanpix (34 digits) within Algorithm 4. A similar effect is seen when using 100 digits in Advanpix versus 6 doubles per high-precision number in `acc`. As the `acc` toolbox is significantly faster, we will use it in all subsequent experiments. Note that ongoing and future modifications of the Ozaki scheme, such as the ones presented in [30], will likely yield further improvements of the accuracy and efficiency of this approach.

¹Note that Matlab's Symbolic Toolbox `vpa` (variable precision arithmetic) supports eigenvalue computations but it does not support the computation of Schur decompositions. Executing `eig` in `vpa` takes 33.70s for a 100×100 matrix in quadruple precision, compared to only 0.30s needed by Advanpix for the complete Schur decomposition.

exact	roots(poly(1:20))	Schur in high-precision, Advanpix with 34 digits	Algorithm 4 (double → quad, using Advanpix with 34 digits)	Algorithm 4 (double → quad, using acc with 2 cells)
20	$1.25 \cdot 10^{-4}$	$-2.55 \cdot 10^{-23}$	$-2.19 \cdot 10^{-23}$	$-1.02 \cdot 10^{-21}$
19	$-1.29 \cdot 10^{-3}$	$8.41 \cdot 10^{-22}$	$-2.70 \cdot 10^{-22}$	$9.17 \cdot 10^{-21}$
18	$6.32 \cdot 10^{-3}$	$-6.20 \cdot 10^{-21}$	$2.25 \cdot 10^{-21}$	$-3.67 \cdot 10^{-20}$
17	$-1.85 \cdot 10^{-2}$	$2.27 \cdot 10^{-20}$	$-7.09 \cdot 10^{-21}$	$9.03 \cdot 10^{-20}$
16	$4.02 \cdot 10^{-2}$	$-5.08 \cdot 10^{-20}$	$1.31 \cdot 10^{-20}$	$-1.60 \cdot 10^{-19}$
15	$-5.93 \cdot 10^{-2}$	$7.55 \cdot 10^{-20}$	$-1.67 \cdot 10^{-20}$	$2.27 \cdot 10^{-19}$
14	$6.98 \cdot 10^{-2}$	$-7.72 \cdot 10^{-20}$	$1.58 \cdot 10^{-20}$	$-2.66 \cdot 10^{-19}$
13	$-6.26 \cdot 10^{-2}$	$5.40 \cdot 10^{-20}$	$-1.15 \cdot 10^{-20}$	$2.52 \cdot 10^{-19}$
12	$4.11 \cdot 10^{-2}$	$-2.40 \cdot 10^{-20}$	$6.47 \cdot 10^{-21}$	$-1.86 \cdot 10^{-19}$
11	$-2.24 \cdot 10^{-2}$	$4.73 \cdot 10^{-21}$	$-2.71 \cdot 10^{-21}$	$1.03 \cdot 10^{-19}$
10	$8.80 \cdot 10^{-3}$	$1.68 \cdot 10^{-21}$	$8.29 \cdot 10^{-22}$	$-4.24 \cdot 10^{-20}$
9	$-2.71 \cdot 10^{-3}$	$-1.70 \cdot 10^{-21}$	$-1.70 \cdot 10^{-22}$	$1.25 \cdot 10^{-20}$
8	$6.05 \cdot 10^{-4}$	$6.72 \cdot 10^{-22}$	$2.59 \cdot 10^{-23}$	$-2.62 \cdot 10^{-21}$
7	$-9.69 \cdot 10^{-5}$	$-1.58 \cdot 10^{-22}$	$-2.48 \cdot 10^{-24}$	$3.76 \cdot 10^{-22}$
6	$1.04 \cdot 10^{-5}$	$2.33 \cdot 10^{-23}$	$1.60 \cdot 10^{-25}$	$-3.46 \cdot 10^{-23}$
5	$-7.05 \cdot 10^{-7}$	$-2.07 \cdot 10^{-24}$	$-1.18 \cdot 10^{-26}$	$1.70 \cdot 10^{-24}$
4	$2.61 \cdot 10^{-8}$	$1.02 \cdot 10^{-25}$	$9.03 \cdot 10^{-28}$	$-2.17 \cdot 10^{-26}$
3	$-4.44 \cdot 10^{-10}$	$-2.44 \cdot 10^{-27}$	$-6.45 \cdot 10^{-29}$	$-1.14 \cdot 10^{-27}$
2	$1.61 \cdot 10^{-12}$	$2.85 \cdot 10^{-29}$	$-1.30 \cdot 10^{-30}$	$-8.12 \cdot 10^{-28}$
1	$5.11 \cdot 10^{-14}$	$-1.95 \cdot 10^{-31}$	$-2.64 \cdot 10^{-29}$	$-8.16 \cdot 10^{-28}$

Table 1: First column: Eigenvalues of the companion matrix for the Wilkinson polynomial. Other columns: Absolute errors when eigenvalues are computed in different ways.

Example 8 *This experiment focuses on the performance of Algorithm 4. We generate a series of random matrices of increasing sizes with random entries from the standard normal distribution. Both the real and the complex Schur forms are computed in high precision. As can be seen from Figure 1, Algorithm 4 is much faster than Advanpix’s `schur`. For this example, our algorithm always converges within 3 iterations to quadruple precision, while 7–8 iterations are needed to attain 100-digit precision. The computed factors Q and T satisfy*

$$\|I - Q^*Q\|_F \leq \begin{cases} 9 \cdot 10^{-32}, & \text{in quadruple precision;} \\ 3 \cdot 10^{-97}, & \text{in 100-digits precision;} \end{cases}$$

$$\|\text{stril}(Q^*AQ)\|_F / \|A\|_F \leq \begin{cases} 3 \cdot 10^{-33}, & \text{in quadruple precision;} \\ 2 \cdot 10^{-98}, & \text{in 100-digits precision,} \end{cases}$$

for all input matrices A .

For more insight, Table 2 shows detailed timings for the largest matrix of size 1000. We note that Algorithm 4 spends essentially all of its time on high-precision matrix-matrix multiplications.

Example 9 *This example aims to provide insight into how Algorithm 4 fails to converge in exceptional situations, when eigenvalues are poorly separated or, equivalently, the triangular matrix equation (6) is very ill conditioned. To illustrate how such a situation affects Algorithm 4, we generate the matrix $A = XDX^{-1}$ of size 150. Here, X is a random matrix of condition 10^5 . The matrix D is a diagonal matrix with diagonal elements chosen uniformly at random between -10 and 10 , with the exception of two clusters of size 10. In*

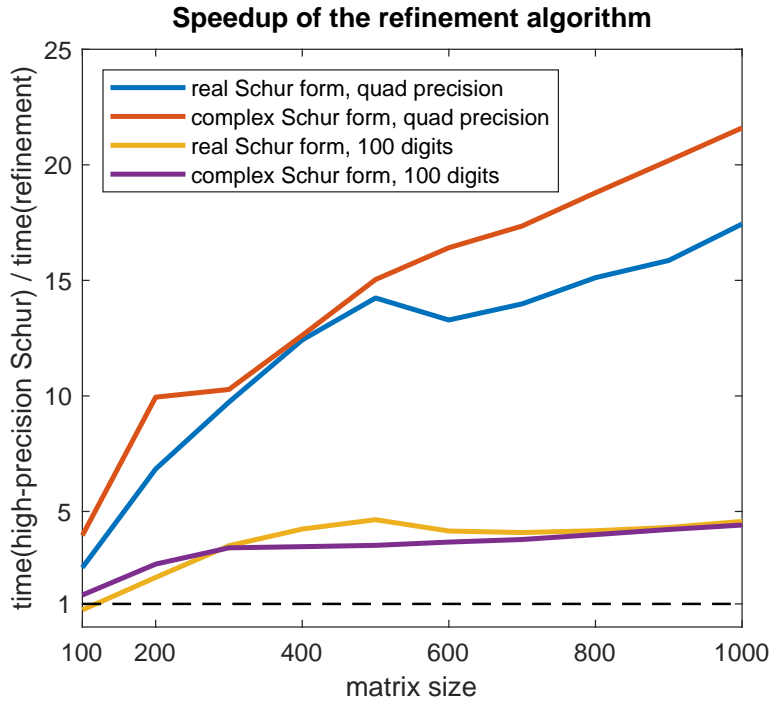


Figure 1: Speedup attained by Algorithm 4 for Example 8. Time needed by Advanpix for computing high-precision Schur decomposition divided by time needed by Algorithm 4 (including the time needed to compute Schur decomposition in double precision).

	quad, real	quad, complex	100 digits, real	100 digits, complex
high-precision Schur	165.14s	614.61s	713.66s	2178.66s
refinement algorithm: total time	9.47s	28.45s	156.18s	493.72s
matrix multiplication	8.59s	27.41s	150.68s	490.72s
other operations	0.88s	1.04s	5.50s	3.00s

Table 2: Time in seconds needed by Advanpix (line 2) and Algorithm 4 (lines 3–5).

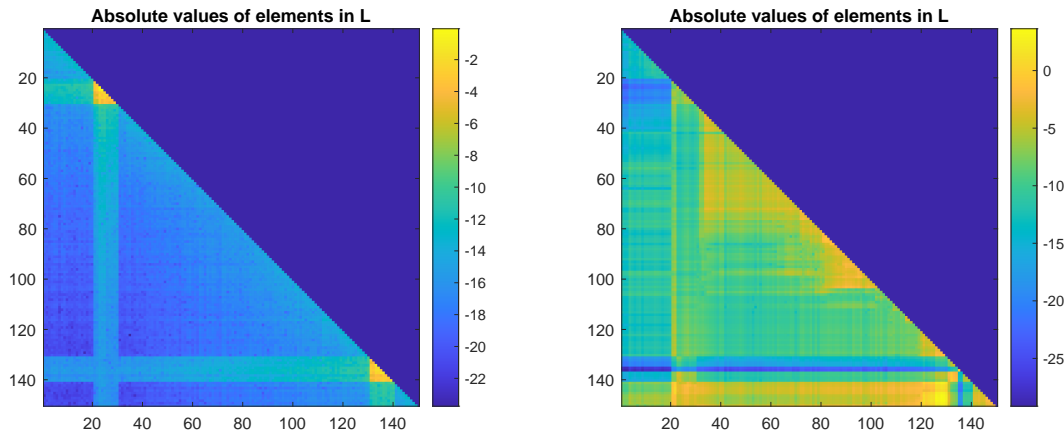


Figure 2: Absolute values of entries of the matrix L obtained from solving the triangular matrix equation (6) within Algorithm 4 applied to the matrix from Example 9. Left: Second iteration of Algorithm 4. Right: Sixth iteration. Color indicates base-10 logarithm of each matrix element’s absolute value.

each cluster, the cluster center is chosen randomly, and the cluster elements are perturbed randomly by at most 10^{-5} from the center.

We execute Algorithm 4 for computing the complex Schur form using double-quadruple precision. Figure 2 shows the absolute values of the computed matrix L during the second and sixth iteration of the algorithm. We note that, initially, all the entries of L are quite small with the exception of those that correspond to the two clusters. In later iterations, these elements have polluted the entire lower triangular part of L . In the eighth iteration, the algorithm fails as the matrix L contains NaNs.

Note that if this example is slightly softened (e.g., lowering the condition of X to 10^4 or the increasing the cluster radius to 10^{-4}) then Algorithm 4 converges in 6 iterations and yields a Schur decomposition in quadruple precision.

Example 10 Finally, we report on the performance of Algorithm 4 for a number of well-known eigenvalue benchmark examples from MatrixMarket [7]. Table 3 details the timing, the errors in the computed factors, and the number of iterations needed for each of these test cases. Note that the matrix size is indicated in the matrix name.

Algorithm 4 fails to converge for matrices denoted by asterisks, due to the reasons explained in Example 9. However, there is a simple trick to avoid the appearance of large numbers in the computed factor L : during its computation, each value larger than, e.g., 10^{-5} in absolute value is immediately set to zero. Currently, there is little theoretical support that such a modification will lead to convergence of Algorithm 4—in fact, it still fails for the matrix from Example 9. However, for all matrices reported in Table 3 this

matrix	time (quad Schur)	#matmuls	time outside of matmuls	$\ I - Q^*Q\ _F$	$\frac{\ \text{stril}(Q^*AQ)\ _F}{\ A\ _F}$	#iterations
olm100	0.17s (0.12s)	12	0.04s	2.73e-32	1.16e-33	3
tols90	0.09s (0.09s)	12	0.02s	2.75e-32	2.72e-34	3
tub100	0.11s (0.15s)	12	0.03s	2.69e-32	2.92e-33	3
bwm200	0.21s (0.96s)	12	0.06s	4.26e-32	3.22e-33	3
olm500	1.25s (7.00s)	12	0.17s	6.48e-32	1.56e-33	3
olm1000	13.61s (50.56s)	16	0.83s	8.07e-32	8.02e-34	4
tub1000	8.48s (90.58s)	12	0.87s	7.51e-32	2.64e-33	3
bwm2000	60.15s (737.59s)	12	4.24s	1.28e-31	3.77e-33	3
dw2048	98.50s (982.43s)	12	4.07s	1.27e-31	3.48e-33	3
rdb200*			did not converge			
rdb2001*			did not converge			
tols1090*			did not converge			
rdb1250*			did not converge			

Table 3: Real Schur decomposition of various matrices from MatrixMarket; double precision is refined to quadruple precision. The second column shows the total time in seconds for Algorithm 4; the time needed by Advanpix’s `schur` is shown in parentheses. The third column shows the total number of quadruple-precision matrix-matrix multiplications needed by Algorithm 4, and the fourth column shows the time the algorithm spent outside these operations. The next two columns show the errors in the computed refined factors Q and T . The last column shows the number of refinement iterations. Algorithm 4 fails to converge for the matrices denoted by asterisk.

matrix	time (quad Schur)	#matmuls	time outside of matmuls	$\ I - Q^*Q\ _F$	$\frac{\ \text{stril}(Q^*AQ)\ _F}{\ A\ _F}$	#iterations
rdb200	0.91s (1.86s)	16	0.07	3.99e-32	3.36e-33	4
rdb2001	0.88s (2.27s)	16	0.07	3.99e-32	2.71e-33	4
tols1090	135.45s (67.77s)	20	1.95	8.93e-32	2.22e-35	5
rdb1250	143.98s (528.12s)	16	1.90	9.79e-32	3.58e-33	4

Table 4: The trick described in Example 10 resolves the convergence failures reported in Table 3 when refining the complex Schur decomposition.

trick successfully resolves convergence problems. We implemented this trick only for the complex Schur form and the obtained results are shown in Table 4.

4 Conclusions

In this work, we have developed iterative algorithms for refining approximate Schur decompositions that exhibit rapid convergence, in theory and in practice. Using the Newton-Schulz iteration for orthogonalization yields an algorithm that carries out most operations in terms of matrix-matrix multiplications, allowing for a simple and efficient implementation. In particular, when refining a double precision decomposition to high precision this allows to leverage the Ozaki scheme and attain significant speedup over existing implementations of algorithms for computing high-precision Schur decompositions.

A number of points deserve further investigation. It is not unlikely that a suitable extension of the modifications described for the symmetric case in [35] will address the

convergence failures observed in Examples 9 and 10. Also, it would be valuable to further study possibilities for merging the improvement of orthogonality and triangularity in a single step, as in [34], and avoiding the need for the initial orthogonalization in Algorithm 4. Finally, to attain very high precision it would certainly be beneficial to study the effective use of more than two levels of precision.

Acknowledgments. The authors thank Takeshi Ogita for providing the Matlab toolbox `acc` based on [37]. They are also grateful to Nicolas Boumal and Christian Lubich for discussions related to Remark 6.

References

- [1] Advanpix: Multiprecision computing toolbox for Matlab. <http://www.advanpix.com/>. Accessed: 2021-12-15.
- [2] Ahmad Abdelfattah, Hartwig Anzt, Erik G. Boman, Erin Carson, Terry Cojean, Jack Dongarra, Alyson Fox, Mark Gates, Nicholas J. Higham, Xiaoye S. Li, Jennifer Loe, Piotr Luszczek, Srikara Pranesh, Siva Rajamanickam, Tobias Ribizel, Barry F. Smith, Kasia Swirydowicz, Stephen Thomas, Stanimire Tomov, Yaohung M. Tsai, and Ulrike Meier Yang. A survey of numerical linear algebra methods utilizing mixed-precision arithmetic. *Int J. High Perform. Comput. Appl.*, 35(4):344–369, 2021.
- [3] Pierre Ablin and Gabriel Peyré. Fast and accurate optimization on the orthogonal manifold without retraction, 2021. arXiv:2102.07432.
- [4] P.-A. Absil, R. Mahony, and R. Sepulchre. *Optimization algorithms on matrix manifolds*. Princeton University Press, Princeton, NJ, 2008. With a foreword by Paul Van Dooren.
- [5] P.-A. Absil, R. Mahony, R. Sepulchre, and P. Van Dooren. A Grassmann-Rayleigh quotient iteration for computing invariant subspaces. *SIAM Rev.*, 44(1):57–73, 2002.
- [6] Uri M. Ascher, Hongsheng Chin, and Sebastian Reich. Stabilization of DAEs and invariant manifolds. *Numer. Math.*, 67(2):131–149, 1994.
- [7] Z. Bai, D. Day, J. W. Demmel, and J. J. Dongarra. A test matrix collection for non-Hermitian eigenvalue problems (release 1.0). Technical Report CS-97-355, Department of Computer Science, University of Tennessee, Knoxville, TN, USA, March 1997. Also available online from <http://math.nist.gov/MatrixMarket>.
- [8] W.-J. Beyn, W. Kleß, and V. Thümmler. Continuation of low-dimensional invariant subspaces in dynamical systems of large dimension. In *Ergodic theory, analysis, and efficient simulation of dynamical systems*, pages 47–72. Springer, Berlin, 2001.

- [9] D. S. Bindel, J. W. Demmel, and M. Friedman. Continuation of invariant subspaces in large bifurcation problems. *SIAM J. Sci. Comput.*, 30(2):637–656, 2008.
- [10] Nicolas Boumal. An introduction to optimization on smooth manifolds. To appear with Cambridge University Press, Jan 2022.
- [11] K. Braman, R. Byers, and R. Mathias. The multishift QR algorithm. I. Maintaining well-focused shifts and level 3 performance. *SIAM J. Matrix Anal. Appl.*, 23(4):929–947, 2002.
- [12] K. Braman, R. Byers, and R. Mathias. The multishift QR algorithm. II. Aggressive early deflation. *SIAM J. Matrix Anal. Appl.*, 23(4):948–973, 2002.
- [13] F. Chatelin. Simultaneous Newton’s iteration for the eigenproblem. In *Defect correction methods (Oberwolfach, 1983)*, volume 5 of *Comput. Suppl.*, pages 67–74. Springer, Vienna, 1984.
- [14] Roy O. Davies and J. J. Modi. A direct method for completing eigenproblem solutions on a parallel computer. *Linear Algebra Appl.*, 77:61–74, 1986.
- [15] J. W. Demmel. Three methods for refining estimates of invariant subspaces. *Computing*, 38:43–57, 1987.
- [16] P. J. Eberlein. A Jacobi-like method for the automatic computation of eigenvalues and eigenvectors of an arbitrary matrix. *J. Soc. Indust. Appl. Math.*, 10:74–88, 1962.
- [17] Bin Gao, Xin Liu, Xiaojun Chen, and Ya-xiang Yuan. A new first-order algorithmic framework for optimization problems with orthogonality constraints. *SIAM J. Optim.*, 28(1):302–332, 2018.
- [18] Gene H. Golub and Charles F. Van Loan. *Matrix computations*. Johns Hopkins Studies in the Mathematical Sciences. Johns Hopkins University Press, Baltimore, MD, fourth edition, 2013.
- [19] J. Greenstadt. A method for finding roots of arbitrary matrices. *Math. Tables Aids Comput.*, 9:47–52, 1955.
- [20] N. J. Higham. *Accuracy and Stability of Numerical Algorithms*. SIAM, Philadelphia, PA, second edition, 2002.
- [21] N. J. Higham. *Functions of matrices*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2008.
- [22] Roger A. Horn and Charles R. Johnson. *Matrix analysis*. Cambridge University Press, Cambridge, second edition, 2013.

- [23] K. Hüper and P. Van Dooren. New algorithms for the iterative refinement of estimates of invariant subspaces. *Journal Future Generation Computer Systems*, 19:1231–1242, 2003.
- [24] I. C. F. Ipsen. Computing an eigenvector with inverse iteration. *SIAM Rev.*, 39(2):254–291, 1997.
- [25] H. A. Jahn. Improvement of an approximate set of latent roots and modal columns of a matrix by methods akin to those of classical perturbation theory. *Quart. J. Mech. Appl. Math.*, 1:131–144, 1948.
- [26] I. Jonsson and B. Kågström. Recursive blocked algorithm for solving triangular systems. I. one-sided and coupled Sylvester-type matrix equations. *ACM Trans. Math. Software*, 28(4):392–415, 2002.
- [27] I. Jonsson and B. Kågström. Recursive blocked algorithm for solving triangular systems. II. Two-sided and generalized Sylvester and Lyapunov matrix equations. *ACM Trans. Math. Software*, 28(4):416–435, 2002.
- [28] W. Kahan. `Refineig`: a program to refine eigensystems, 2007. Available from <https://people.eecs.berkeley.edu/~wkahan/Math128/RefinEig.pdf>.
- [29] M. M. Konstantinov, P. Hr. Petkov, and N. D. Christov. Nonlocal perturbation analysis of the Schur system of a matrix. *SIAM J. Matrix Anal. Appl.*, 15(2):383–392, 1994.
- [30] Marko Lange and Siegfried M. Rump. An alternative approach to Ozaki’s scheme for error-free transformation of matrix multiplication, 2022. Presentation at International Workshop on Reliable Computing and Computer-Assisted Proofs (ReCAP 2022).
- [31] E. Lundström and L. Eldén. Adaptive eigenvalue computations using Newton’s method on the Grassmann manifold. *SIAM J. Matrix Anal. Appl.*, 23(3):819–839, 2002.
- [32] C. Mehl. On asymptotic convergence of nonsymmetric Jacobi algorithms. *SIAM J. Matrix Anal. Appl.*, 30(1):291–311, 2008.
- [33] Jorge Nocedal and Stephen J. Wright. *Numerical optimization*. Springer Series in Operations Research and Financial Engineering. Springer, New York, second edition, 2006.
- [34] Takeshi Ogita and Kensuke Aishima. Iterative refinement for symmetric eigenvalue decomposition. *Japan Journal of Industrial and Applied Mathematics*, 35:1007–1035, 2018.

- [35] Takeshi Ogita and Kensuke Aishima. Iterative refinement for symmetric eigenvalue decomposition II: clustered eigenvalues. *Jpn. J. Ind. Appl. Math.*, 36(2):435–459, 2019.
- [36] Dianne P. O’Leary and G. W. Stewart. Data-flow algorithms for parallel matrix computations. *Comm. ACM*, 28(8):840–853, 1985.
- [37] Katsuhisa Ozaki, Takeshi Ogita, Shin’ichi Oishi, and Siegfried M. Rump. Error-free transformations of matrix multiplication by using fast routines of matrix multiplication and its applications. *Numerical Algorithms*, 59(1):95–118, January 2012.
- [38] G. W. Stewart. A Jacobi-like algorithm for computing the Schur decomposition of a non-Hermitian matrix. *SIAM J. Sci. Statist. Comput.*, 6(4):853–864, 1985.
- [39] J.-G. Sun. Perturbation bounds for the Schur decomposition. Report UMINF-92.20, Department of Computing Science, Umeå University, Umeå, Sweden, 1992.
- [40] J.-G. Sun. Perturbation bounds for the generalized Schur decomposition. *SIAM J. Matrix Anal. Appl.*, 16(4):1328–1340, 1995.
- [41] H. Wielandt. Beiträge zur mathematischen behandlung komplexer eigenwertprobleme, teil v: Bestimmung höherer eigenwerte durch gebrochene iteration. Bericht b 44/j/37, Aerodynamische Versuchsanstalt Göttingen, Germany, 1944.
- [42] James H. Wilkinson. *The Algebraic Eigenvalue Problem*. Clarendon Press, Oxford, 1965.