

An algorithm for the complete solution of the quartic eigenvalue problem

ZLATKO DRMAČ, IVANA ŠAIN GLIBIĆ

Abstract

Quartic eigenvalue problem $(\lambda^4 A + \lambda^3 B + \lambda^2 C + \lambda D + E)x = \mathbf{0}$ naturally arises in a plethora of applications, e.g. when solving the Orr–Sommerfeld equation in the stability analysis of the Poiseuille flow, in theoretical analysis and experimental design of locally resonant phononic plates, modeling a robot with electric motors in the joints, calibration of catadioptric vision system, or e.g. computation of the guided and leaky modes of a planar waveguide. This paper proposes a new numerical method for the full solution (all eigenvalues and all left and right eigenvectors) that, starting with a suitable linearization, uses an initial, structure preserving, reduction designed to reveal and deflate certain number of zero and infinite eigenvalues before the final linearization is forwarded to the QZ algorithm. The backward error in the reduction phase is bounded column wise in each coefficient matrix, which is advantageous if the coefficient matrices are graded. Numerical examples show that the proposed algorithm is capable of computing the eigenpairs with small residuals, and that it is competitive with the available state of the art methods.

1 Introduction and preliminaries

We propose a new method for numerical solution of the *quartic eigenvalue problem*

$$(\lambda^4 A + \lambda^3 B + \lambda^2 C + \lambda D + E)x = \mathbf{0}, \quad (1)$$

where the coefficient matrices $A, B, C, D, E \in \mathbb{C}^{n \times n}$ are assumed general, with no particular structure (such as symmetry, sparsity). We are interested in the full solution, i.e. computation of all eigenvalues with the corresponding (left and/or right) eigenvectors, and our ultimate goal is to provide a robust mathematical software that can be used in ever increasing number of applications in applied sciences and engineering.

The quartic eigenvalue problem naturally arises in solving the Orr – Sommerfeld equation which appears in the hydrodynamic analysis of the stability of the Poiseuille flow by eliminating the pressure from the linearized Navier-Stokes equation. Other applications include e.g. theoretical analysis and experimental design of locally resonant phononic plates [38], finite element analysis of two dimensional phononic crystals [34], modeling a robot with electric motors in the joints [25], computing deformation modes of thin-walled structures [35], or e.g. computation of the guided and leaky modes of a planar waveguide [28], or solving an optical waveguiding problem involving atomically thick 2D materials [27]. In these examples the matrix eigenvalue problem is the result of discretization of differential operators and thus (depending on the discretization method) the coefficient matrices are sparse and usually only some eigenvalues are needed – those may be prescribed by specifying e.g. a region of interest in the complex plane.

In such cases, methods for large sparse problems such as e.g. NLFEAST [18], [26], [39], [7] will find a subspace that contains eigenvectors of interest, and then Rayleigh–Ritz extraction uses the projected problem in which the Rayleigh quotients are medium size dense matrices. Reliable solution of the projected problem is important both for the convergence of the iterations towards the wanted part of the spectrum (e.g. for robust implementation of locking and purging) and for the accuracy of the computed solution.

These examples illustrate the wide spectrum of important applications of the quartic eigenvalue problem, and justify, even demand, development of methods specialized for (1). Yet, to the best of our knowledge, there is no published custom-built solver with a supporting analysis that would provide certain level of confidence/guarantee that is comparable e.g. to the currently available solvers of the quadratic eigenvalue problem such as [20], [14]. Instead, (1) is usually numerically solved by a standard linearization and deployment of the solvers such as `polyeig` in Matlab. On the other hand, numerical difficulties in solving nonlinear eigenvalue problems become nontrivial even in the simplest case of the polynomial quadratic problem, which is at the core of the theory and applications of mechanical systems. More carefully designed custom-made algorithm often proves much better than a generic solver – an excellent example is the quadratic eigenvalue problem, where the algorithms `quadeig` [20] and `kvadeig` [14] outperform `polyeig`, in particular when the spectrum contains multiple infinite eigenvalues.

1.1 Backward stability and scaling

Numerical algorithms for computing the eigenvalues and the corresponding right and left eigenvectors of a general regular matrix polynomial $P_k(\lambda) = \sum_{i=0}^k \lambda^i A_i$ usually consist of three main steps: (i) linearization, i.e. definition of an equivalent linear generalized eigenvalue problem for a suitably constructed $kn \times kn$ pencil $\mathcal{A} - \lambda\mathcal{B}$; (ii) computation of the eigenpairs of the linearization $\mathcal{A} - \lambda\mathcal{B}$, using e.g. the QZ algorithm; (iii) reconstruction of the eigenpairs of the original problem. Some algorithms include a preprocessor that transforms the linearization (i) into a form that reveals some canonical (e.g. spectral) structure and that is in some numerical sense better input to a particular software implementation of the QZ algorithm in (ii) (e.g. scaling). For a general framework of this scheme we refer the reader to [32], [33].

1.1.1 Weak and strong norm-wise backward stability

In order to be used with confidence in applications, the eigenvalues and eigenvectors computed in finite precision arithmetic are justified by proving that they are exact spectral elements of a nearby polynomial

$$\tilde{P}_k(\lambda) = \sum_{i=0}^k \lambda^i (A_i + \Delta A_i).$$

If the sizes of the perturbations (backward errors) ΔA_i are appropriately small (e.g. of the same order of magnitude as initial uncertainties in the coefficients A_i , as measured in a matrix norm) then the computation is usually deemed backward stable.

In a numerical algorithm, the canonical structure revealing steps [32], [33] and the QZ algorithm are based on unitary transformations, so that the entire process is backward stable – the computed result corresponds exactly to a linear pencil

$$\mathcal{A} + \Delta\mathcal{A} - \lambda(\mathcal{B} + \Delta\mathcal{B}), \quad \|\Delta\mathcal{A}\|_F \leq \xi\|\mathcal{A}\|_F, \quad \|\Delta\mathcal{B}\|_F \leq \xi\|\mathcal{B}\|_F, \quad 0 \leq \xi \ll 1.$$

We refer to this as *strong norm-wise backward stability* of the solution of the linearized problem. However, this statement must be carefully interpreted. The QZ algorithm is oblivious to the

underlying structure of the linear pencil, and $\mathcal{A} + \Delta\mathcal{A} - \lambda(\mathcal{B} + \Delta\mathcal{B})$ most likely will not have the structure of the linearization of a matrix polynomial, and the backward stability cannot be stated in terms of the original polynomial eigenproblem.

Relating the pencil $\mathcal{A} + \Delta\mathcal{A} - \lambda(\mathcal{B} + \Delta\mathcal{B})$ with a matrix polynomial close to $P_k(\lambda)$ requires an additional theoretical construction in the error analysis. For large classes of linearizations, there is an equivalence transformation

$$\mathcal{A} + \Delta\mathcal{A} - \lambda(\mathcal{B} + \Delta\mathcal{B}) \longrightarrow (I + E)[\mathcal{A} + \Delta\mathcal{A} - \lambda(\mathcal{B} + \Delta\mathcal{B})](I + F), \quad \|E\|_F \leq \epsilon_1, \quad \|F\|_F \leq \epsilon_2,$$

such that the new pencil is the linearization of $\tilde{P}_k(\lambda) = \sum_{i=0}^k \lambda^i (A_i + \Delta A_i)$, where, under certain assumptions,

$$\|(\Delta A_0 \quad \Delta A_1 \quad \dots \quad \Delta A_k)\|_F \leq \epsilon_3 \| (A_0 \quad A_1 \quad \dots \quad A_k) \|_F, \quad (2)$$

with some $0 \leq \epsilon_1, \epsilon_2, \epsilon_3 \ll 1$ that depend on the roundoff unit, dimensions of the problem and algorithmic details. This form of *weak norm-wise backward stability* bounds each $\|\Delta A_i\|_F$ relative to the norm of the coefficients array $(A_0 \quad A_1 \quad \dots \quad A_k)$. For detailed in depth discussion see [33, §4], [22], [11], [12].

Note that (2) may be difficult to interpret in an application where the coefficient matrices carry information of different physical nature (e.g. mass, damping and stiffness) expressed in appropriate physical units. Unfortunately, except in some particular cases, (2) cannot be strengthened into *strong norm-wise backward stability*¹ estimate

$$\|\Delta A_i\|_F \leq \epsilon \|A_i\|_F, \quad i = 0, \dots, k. \quad (3)$$

Construction of an algorithm with the backward error (3) may not be feasible without the framework of mixed error analysis, see [23] where this is shown for the case $n = 1, k = 2$. Constructing an algorithm that has guaranteed (provable) strong norm-wise backward/mixed stability is a challenging open problem.

1.1.2 Backward stability of individual eigenpairs. Residual

Another way to measure the quality of an approximate eigenpair (an eigenvalue with a corresponding right or left eigenvector) *a posteriori* is through the residual, which we discuss next, in terms of our original problem (1). If (λ, x) is a computed eigenpair with the right eigenvector² x , then the minimal size of backward error of the type (3) that makes (λ, x) an exact eigenpair of a backward perturbed quartic eigenvalue problem, i.e.

$$\min\{\epsilon : (\lambda^4(A + \Delta A) + \lambda^3(B + \Delta B) + \lambda^2(C + \Delta C) + \lambda(D + \Delta D) + (E + \Delta E))x = \mathbf{0}, \\ \|\Delta A\|_F \leq \epsilon \|A\|_F, \|\Delta B\|_F \leq \epsilon \|B\|_F, \|\Delta C\|_F \leq \epsilon \|C\|_F, \|\Delta D\|_F \leq \epsilon \|D\|_F, \|\Delta E\|_F \leq \epsilon \|E\|_F\},$$

can be explicitly computed as the normalized residual [29, §2.2]

$$\eta(\lambda, x) = \frac{\|(\lambda^4 A + \lambda^3 B + \lambda^2 C + \lambda D + E)x\|_2}{(|\lambda|^4 \|A\|_2 + |\lambda|^3 \|B\|_2 + |\lambda|^2 \|C\|_2 + |\lambda| \|D\|_2 + \|E\|_2) \|x\|_2}, \quad \lambda \neq \infty; \\ \eta(\infty, x) = \frac{\|Ax\|_2}{\|A\|_2 \|x\|_2}. \quad (4)$$

The key difficulty is that an eigenpair obtained by this procedure may have large residual (norm-wise backward error (4)), although the norm-wise backward error for the eigenpair (with

¹An anonymous referee suggested the term *coefficient-wise backward stability*.

²For the sake of brevity, here we omit an analogous discussion with the left eigenvector.

the same λ) of the corresponding linearization³ is acceptably small. This phenomenon is further analyzed in [21], and it is proven that this kind of variation in the backward errors is due to the fact that the norms of the coefficient matrices of the original problem are not equilibrated. As a result, the backward stability of the linearization is not inherited by the transformation to the original polynomial. The problem can be alleviated by parameter scaling, which we review next.

1.1.3 Parameter scaling

Parameter scaling is a useful, powerful, albeit not omnipotent, tool for stabilizing polynomial eigensolvers; the techniques vary from simple heuristics to sophisticated concepts from tropical polynomial algebra. Here we briefly review the scaling used in this paper; other scalings can be easily incorporated.

Write $\lambda = \gamma\nu$, where $\gamma > 0$ is parameter to be determined. Using an additional free parameter $\theta > 0$, define the scaled quartic polynomial as

$$\theta E + \nu(\gamma\theta D) + \nu^2(\gamma^2\theta C) + \nu^3(\gamma^3\theta B) + \nu^4(\gamma^4\theta A) \equiv \widehat{E} + \nu\widehat{D} + \nu^2\widehat{C} + \nu^3\widehat{B} + \nu^4\widehat{A}.$$

The parameters γ and θ are defined so that the norms of the new coefficient matrices do not vary much, and are close to a target value. This can be done by adapting the Fan, Lin and Van Dooren's scaling [17]. For γ , we choose $\gamma = \sqrt[4]{\frac{\|E\|_F}{\|A\|_F}}$, which is the optimal γ for minimizing the factor $\max(1, \|\widehat{A}\|_F, \|\widehat{B}\|_F, \|\widehat{C}\|_F, \|\widehat{D}\|_F, \|\widehat{E}\|_F)^2 / \min(\|\widehat{E}\|_F, \|\widehat{A}\|_F)$ in the backward error ratio bounds [1], and for θ , we choose, following [6], $\theta = 4/(\|E\|_F + \gamma\|D\|_F + \gamma^2\|C\|_F + \gamma^3\|B\|_F)$.

Although simple and easy to implement in any polynomial eigensolver, scaling can achieve good results in controlling the growth factor of the backward error. For instance, [40] showed that the quadratic eigenvalue problem can be solved with small backward errors (4) in all eigenpairs by carefully examining six linearizations. Detailed analysis showed that backward stability in all eigenpairs was achieved using two linearizations.

Note that in this interpretation of backward stability, the optimal backward errors (4) are constructed separately for each eigenpair, which is different from the backward stability discussed in §1.1.1.

1.1.4 Graded matrices

We should keep in mind that, in addition to different magnitudes of the norms $\|A\|_F, \|B\|_F, \dots, \|E\|_F$, the entries inside each coefficient matrix may be on different scales of magnitude, where some small entries may be important parameters. If, for instance, A has graded columns whose norms vary over several orders of magnitude, and ΔA is small perturbation that satisfies $\|\Delta A\|_F \leq \epsilon\|A\|_F$ with a small $\epsilon > 0$, then some small columns of A may be completely wiped out by ΔA . (In an engineering application the columns of the coefficient matrices may be scaled to properly interpret the norm of the eigenvector x whose components are of different physical nature.) Parameter scaling cannot counteract differently scaled columns in a particular coefficient matrix.

Remark 1. *In addition to parameter scaling, we can use diagonal scaling matrices Δ_ℓ and Δ_r , for scaling all coefficients from the left with Δ_ℓ and from the right with Δ_r . The goal is to equilibrate the absolute values of all matrix entries. These scaling matrices can be computed by an extension of the scheme described in [14, §4.2].*

³The norm-wise backward error of λ as an approximate eigenvalue of the linearization is computed analogously to (4).

1.2 A quandary about the infinite eigenvalues

The presence of infinite eigenvalues, indicated by the rank deficiency of A , may cause difficulties in the QZ algorithm, which is usually deployed for solving the linearized problem; infinite eigenvalues may not be identified correctly, they may have negative impact on the accuracy of the computed finite eigenvalues, see e.g. [32, Example 2]. It is then advantageous to remove infinite eigenvalues by a deflation and proceed with a problem of smaller dimension, with only finite eigenvalues. This framework, introduced in [20], proved much better than direct solution of the linearized problem.

In some cases, certain number of infinite eigenvalues of (1) can be identified and removed already during the problem formulation. An illustrative example is given in the eigenvalue problem for the channel and Blasius boundary layer in semi-infinite domain [8]. The Orr–Sommerfeld differential equation is discretized using the Chebyshev collocation matrix method, and the boundary conditions are imposed in E ; the remaining matrix coefficients A , B , C , D have the corresponding last four rows equal to zero. In the case of linearly independent boundary conditions, by a clever column permutation, four infinite eigenvalues can be separated and deflated, see [8] for technical details.

The structure of the infinite eigenvalue (the number and the dimensions of blocks in the Kronecker Canonical Form (KCF)) cannot be inferred by only inspecting the rank of A . Rank deficiency in A reveals only certain number of infinite eigenvalues and further steps are necessary to either confirm that there are no more infinite eigenvalues or to reveal more blocks carrying $\lambda = \infty$ in the KCF. For more details see [32] and [33].

Removals of infinite and zero eigenvalues involve decisions on the numerical ranks of some intermediate matrices that have been contaminated by the roundoff noise from the previous steps. If the data is not well scaled, and if the computation cannot be interpreted as backward stable in terms of the original coefficients that may have an initial uncertainty from the very problem formulation, then there may be quite a few spurious eigenvalues with large absolute values. The backward stability of the beginning steps that carry the critical responsibility of removing as many as possible infinite eigenvalues must be as much as possible in terms of the initial coefficient matrices, and it has to be as much as possible structured, e.g. column-wise small (backward error in each column small relative to that column's norm) instead of only small in matrix norm.

As we pointed out in [14], the goal of the pre-processing is to remove many zero and infinite eigenvalues before calling the QZ algorithm, and to use QZ software optimized for a given computing machinery. The reason is that handling infinities numerically in QZ is a delicate issue with many fine details [37], and the development of optimized software often uses techniques, such as e.g. block-oriented formulations and parallelization, that accept speed-accuracy trade-offs.

Another possibility to deal with infinite eigenvalues of matrix polynomials, pursued in [31], [30], is, after a suitable linearization and scaling, to modify software implementation of the QZ algorithm by lowering the original threshold for setting small numbers to zero in the part of the algorithm that can create infinite eigenvalues. This change of a critical parameter was compensated by increasing the maximal allowed number of iterations. This method performed well as measured by the residual of refined eigenvectors $\eta_P(\lambda) = \min_{x \neq 0} \eta(\lambda, x)$.

Although the goal to safely deflate eigenvalues that may be difficult for QZ iterations is the same, our approach of deflation, based as much as possible on initial data, is conceptually different, as we describe next.

2 A new approach to the quartic eigenvalue problem

In our recent paper [14], we built upon the `quadeig` algorithm of [20] and [32], [33] and constructed an algorithm (designated as `kvadeig`) for the quadratic eigenvalue problem that makes several reduction steps toward the KCF. One of distinctive features of that reduction is that the backward error in the coefficient matrices is bounded on a finer-scale, e.g. column-wise. Although such a column-wise error bound does not extend to the entire algorithm (the QZ algorithm enjoys only the norm-wise bound), it may be of critical importance in the beginning steps when decisions about the zero and infinite eigenvalues have to be made, in particular if the matrix coefficients are graded, as discussed in §1.1.4. Since this issue, tackled in `kvadeig`, is separate from parameter scaling, the approach introduced in `kvadeig` can benefit from any good scaling, so that it can be combined e.g. with the strategy introduced in [40].

In this paper we extend the techniques of `quadeig` and `kvadeig` to the quartic eigenvalue problem (1). A direct connection of (1) with the quadratic problem is quadratification. In §2.1, we briefly review quadratification by companion forms of grade 2, and then we discuss practical advantages and shortcomings of this approach to the quartic eigenvalue problem. Then, in §2.2 we present the main idea of the paper – a linearization based on the quadratification provides a two-level structure that allows for a generalization of the scheme used in `kvadeig`.

2.1 Quadratification

Both `quadeig` and `kvadeig` outperform the general solvers such as e.g. `polyeig` from Matlab. In addition, [40] provides a backward stable algorithm (in the sense of residuals reviewed in §1.1.2, albeit at double cost) whose semi-tropical scaling can be used in `quadeig/kvadeig` as well. Hence, good quadratic solvers supported by numerical analysis are available. If one has these quadratic solvers implemented in a reliable software, then it makes sense to use *quadratification* [9] to reduce the quartic problem to a quadratic one and use the off the shelf quadratic code.

To that end, define matrix polynomials $B_1(\lambda) = \lambda^2 C + \lambda D + E$, $B_2(\lambda) = \lambda^2 A + \lambda B$. The first and the second companion form of grade 2 are then defined, respectively, as

$$C_1^2(\lambda) = \begin{pmatrix} B_2(\lambda) & B_1(\lambda) \\ -\mathbb{I}_n & \lambda^2 \mathbb{I}_n \end{pmatrix} = \begin{pmatrix} \lambda^2 A + \lambda B & \lambda^2 C + \lambda D + E \\ -\mathbb{I}_n & \lambda^2 \mathbb{I}_n \end{pmatrix} = \lambda^2 \begin{pmatrix} A & C \\ \mathbf{0} & \mathbb{I}_n \end{pmatrix} + \lambda \begin{pmatrix} B & D \\ \mathbf{0} & \mathbf{0} \end{pmatrix} + \begin{pmatrix} \mathbf{0} & E \\ -\mathbb{I}_n & \mathbf{0} \end{pmatrix}.$$

$$C_2^2(\lambda) = \begin{pmatrix} B_2(\lambda) & -\mathbb{I}_n \\ B_1(\lambda) & \lambda^2 \mathbb{I}_n \end{pmatrix} = \lambda^2 \begin{pmatrix} A & \mathbf{0} \\ C & \mathbb{I}_n \end{pmatrix} + \lambda \begin{pmatrix} B & \mathbf{0} \\ D & \mathbf{0} \end{pmatrix} + \begin{pmatrix} \mathbf{0} & -\mathbb{I}_n \\ E & \mathbf{0} \end{pmatrix} = \lambda^2 \mathbb{M} + \lambda \mathbb{C} + \mathbb{K}. \quad (5)$$

Both $C_1^2(\lambda)$ and $C_2^2(\lambda)$ are strong quadratifications, see [Theorem 5.3, Theorem 5.4][9], [10].

2.1.1 Why quadratification alone is not enough?

While solving the eigenvalue problem for (5) by a reliable quadratic eigensolver is a viable approach, its straightforward implementation has a significant limitation in light of the discussions in §1.1 and §1.2. Namely, just as the QZ algorithm is oblivious to the structure of the linearization, the quadratic eigensolver will be oblivious to the fact that the quadratic pencil represents a quadratification of the quartic problem and that the matrices \mathbb{M} , \mathbb{C} , \mathbb{K} in (5) have block structure defined by the matrices of the original problem. As a result, in the preprocessing phase the algorithm will not make the critical decisions (such as numerical rank revealing) based on the original coefficients of the quartic problem (1).

On the other hand, with a suitable choice of the quadratification and its linearization, we can generalize the framework of `kvadeig` by zooming into the block structure of the matrices

constructed in the quadratification, and thus work with the original coefficients. This is the key idea in this work, and in the rest of this section we set the scene and present the structure of the paper. We will use the second companion form of grade 2 because its structure is compatible with the deflation scheme of `kvadeig`.

2.2 A generalization of the deflation scheme from `KVADEIG`

The starting point of the development of the proposed algorithm is the quadratic polynomial (5). It can be further linearized using e.g. the second companion form. In that case, the final matrix pencil of size $4n \times 4n$, that represents a linearization of the quartic problem 1, reads

$$\mathbb{A} - \lambda \mathbb{B} = \begin{pmatrix} \mathbb{C} & -\mathbb{I}_{2n} \\ \mathbb{K} & \mathbf{0}_{2n} \end{pmatrix} - \lambda \begin{pmatrix} -\mathbb{M} & \mathbf{0}_{2n} \\ \mathbf{0}_{2n} & -\mathbb{I}_{2n} \end{pmatrix} = \left(\begin{array}{c|c|c|c} B & \mathbf{0}_n & -\mathbb{I}_n & \mathbf{0}_n \\ \hline D & \mathbf{0}_n & \mathbf{0}_n & -\mathbb{I}_n \\ \hline \mathbf{0}_n & -\mathbb{I}_n & \mathbf{0}_n & \mathbf{0}_n \\ \hline E & \mathbf{0}_n & \mathbf{0}_n & \mathbf{0}_n \end{array} \right) - \lambda \left(\begin{array}{c|c|c|c} -A & \mathbf{0}_n & \mathbf{0}_n & \mathbf{0}_n \\ \hline -C & -\mathbb{I}_n & \mathbf{0}_n & \mathbf{0}_n \\ \hline \mathbf{0}_n & \mathbf{0}_n & -\mathbb{I}_n & \mathbf{0}_n \\ \hline \mathbf{0}_n & \mathbf{0}_n & \mathbf{0}_n & -\mathbb{I}_n \end{array} \right). \quad (6)$$

Notice that (6) is actually a block Kronecker linearization up to simultaneous interchanges of block rows 2 and 3, and block columns 2 and 3; hence backward error analysis of type (2) applies, see [13, Theorem 5.22, Corollary 5.24]. However, we find the interpretation via quadratification (5) more intuitive and more natural for generalization of the scheme from [14].

Now, we can follow the structure of `quadeig/kvadeig`, attempting to deflate the infinite eigenvalues of $\lambda^2 \mathbb{M} + \lambda \mathbb{C} + \mathbb{K}$. Even if that is expected to perform better than a straightforward companion type linearization followed by `polyeig`, it is not the best one can do, see §2.1.1. Instead, the goal is to implement the beginning critical steps, including deflations, with small (hopefully to some extent structured) backward error in the original coefficient matrices A , B , C , D , E . Therefore, on the global level, we follow the strategy of `kvadeig` [14, Algorithm 3.1] to bring (6) to an upper triangular Kronecker Canonical Form (KCF), but the elementary steps are rewritten in terms of the original matrices whenever feasible. We will keep the two-level partitioning throughout the paper, and try to use transformations that respect the block structure as much as possible and, if deflation is needed, the structure of the linearization should be considered when defining the transformation matrices.

2.2.1 Outline of the paper

The new algorithm, designated as `kvarteig`, is described in detail in §3. Recovering the eigenvectors of (1) from those of (6) is discussed in detail in §4, where we propose using least squares regularization, and suggest algorithmic details for an efficient software implementation. In §5 we provide detailed backward error analysis of the first two steps that are critical for removing zero and infinite eigenvalues. We clearly identify moments in the algorithm where scaling of the data plays the key role in keeping the backward error in the initial data small. The numerical experiments, presented in §6, show the advantage of the new framework, both of `kvarteig` and `kvadeig` (applied to the quadratification). The strong point of the proposed algorithm is the deflation process. The biggest differences in the results, as compared to other algorithms, can be seen in the element-wise backward errors and, in particular, in the examples with zero and/or infinite eigenvalues. Altogether, the numerical examples clearly demonstrate the importance of both scaling (including balancing) and deflation in the pre-processing phase.

The material of this paper should be considered as the second part of [14], and numerical results in §6 once more illustrate the power of the approach introduced in `quadeig` and `kvadeig`.

3 The kvarteig algorithm

We now describe the main ideas of the proposed procedure for deflating infinite and/or zero eigenvalues. As outlined in §2.2, our plan is to adapt the deflation scheme from `quadeig`/`kvadeig` to the linearization (6). Although the structure of the proposed algorithm is inspired by our recent quadratic eigensolver and its connection to the problem (1) via quadratification, we stress again that the new algorithm is not a composition of quadratification and quadratic eigensolver. We recall our discussion in §2.1.1 that applying a robust quadratic solver to the quadratification blindly (i.e. by ignoring the origin of the quadratic problem) is not satisfactory.

The first immediate problem is revealing the numerical rank of the coefficient matrices. In §3.1, we briefly review this issue and use it to illustrate the two-level approach to the linearization (6) – at the level of the 2×2 partition, the algorithm mimics the structure of `kvadeig`, but all operations are adapted to the 4×4 block structure of the linearization and then, in §3.2, further tailored for the quartic problem.

3.1 Numerical rank and block-structure

Revealing infinite and zero eigenvalues in presence of perturbations is a delicate task because it depends on the numerical ranks [19] of matrices that are either initial coefficients (possibly polluted by noise) or intermediate results in finite precision computation. An additional difficulty is the underlying structure of the involved matrices, that should preferably be preserved in a backward stability interpretation of the computed results. This is one of the reasons why applying `quadeig`/`kvadeig` directly to a quadratification is numerically not optimal.

Namely, applying a rank revealing decomposition (such as the SVD or the pivoted QR factorization) to \mathbb{M} would mean looking for a small perturbation $\delta\mathbb{M}$ such that $\mathbb{M} + \delta\mathbb{M}$ has lower rank that cannot be further reduced by a small perturbation. Such a construction does not respect the block structure of \mathbb{M} , and better way is to think at this step in terms of the numerical rank with constrained perturbation. If J denotes the first n columns of \mathbb{I}_{2n} then the allowed perturbation might be $\delta\mathbb{M} = J\delta AJ^T$ with an $n \times n$ δA . Similarly, the numerical rank of \mathbb{K} will be determined under the constraint that only $\mathbb{K}(n+1 : 2n, 1 : n) = E$ is allowed to change. For a systematic treatment of the general case using the generalized SVD, see [41].

Since we have the natural block structure, we can formulate the rank revealing steps directly, in terms of the original coefficients. This defines the first step of the procedure whose details are explained in §3.2.

Let $r_A = \text{rank}(A)$, $r_E = \text{rank}(E)$ and let

$$A\Pi_A = Q_A R_A, \quad R_A = \begin{pmatrix} \widehat{R}_A \\ \mathbf{0}_{n-r_A, n} \end{pmatrix}, \quad E\Pi_E = Q_E R_E, \quad R_E = \begin{pmatrix} \widehat{R}_E \\ \mathbf{0}_{n-r_E, n} \end{pmatrix}, \quad (7)$$

be the rank revealing QR factorizations for A and E , computed as in [5], [15]. Note that (7) yields a structure preserving rank revealing decomposition of the matrix $\mathbb{M} = \begin{pmatrix} A & \mathbf{0} \\ C & \mathbb{I}_n \end{pmatrix}$ as

$$\mathbb{M}\Pi_M = Q_M R_M, \quad Q_M = \left(\begin{array}{c|c} \mathbf{0} & Q_A \\ \mathbb{I}_n & \mathbf{0} \end{array} \right), \quad \Pi_M = \left(\begin{array}{c|c} \mathbf{0} & \Pi_A \\ \mathbb{I}_n & \mathbf{0} \end{array} \right), \quad R_M = \left(\begin{array}{c|c} \mathbb{I}_n & C\Pi_A \\ \mathbf{0} & R_A \end{array} \right). \quad (8)$$

The truncation of R_M is done by truncating R_A , and the truncation can be pushed back into a backward perturbation of A ; see [14, §2.1, §2.3].

Similarly, the rank revealing factorization of the matrix $\mathbb{K} = \begin{pmatrix} \mathbf{0} & -\mathbb{I}_n \\ E & \mathbf{0} \end{pmatrix}$ is

$$\mathbb{K}\Pi_K = Q_K R_K, \quad Q_K = \left(\begin{array}{c|c} \mathbb{I}_n & \mathbf{0} \\ \mathbf{0} & Q_E \end{array} \right), \quad \Pi_K = \left(\begin{array}{c|c} \mathbf{0} & \Pi_E \\ \mathbb{I}_n & \mathbf{0} \end{array} \right), \quad R_K = \left(\begin{array}{c|c} -\mathbb{I}_n & \mathbf{0} \\ \mathbf{0} & R_E \end{array} \right). \quad (9)$$

Notice that the permutation of the column blocks only ensures that the matrix R_K is upper triangular. If this structure is not important for the process, we can skip the permutation step and just make the following transformation

$$\begin{pmatrix} \mathbb{I}_n & \mathbf{0} \\ \mathbf{0} & Q_E^* \end{pmatrix} \mathbb{K} \begin{pmatrix} \Pi_E & \mathbf{0} \\ \mathbf{0} & \mathbb{I}_n \end{pmatrix} = \begin{pmatrix} \mathbf{0} & -\mathbb{I}_n \\ R_E & \mathbf{0} \end{pmatrix}. \quad (10)$$

Remark 2. *To determine the numerical rank using the rank revealing QR factorization, we use the thresholding strategies as in [14, §2.3.1]. For a softer thresholding we look for a drop-off of absolute values of two consecutive diagonal entries in the upper triangular form. In general, determination of the numerical rank (thresholding strategy and thresholds for truncating the triangular factor) should take into account the size and the structure of the initial uncertainty in the data. Such an additional information is application specific.*

3.2 The decision tree of kvarteig

The algorithm is designed to remove zero eigenvalues; the infinities are removed by switching to the reversed pencil. Similarly as in⁴ [14], the deflation process is an adaptation of the algorithm by [32] for computing the structure of the eigenvalues 0 and ∞ . The first two steps are modified using the structure of the linearization (6), and for possible additional steps the algorithm proceeds with the rank revealing QR factorizations and carefully implemented URV decompositions.

As in the `kvadeig`, there are three main cases: both A and E regular; only one of A and E is singular; and both A and E are singular.

3.2.1 Both matrices A and E regular

If both matrices A and E are regular, we can use the factorization (8) to reduce the matrix \mathbb{B} from (6) to upper triangular form, since this is already the first step of the QZ algorithm.

$$\begin{aligned} & \begin{pmatrix} Q_M^* & \mathbf{0} \\ \mathbf{0} & \mathbb{I}_{2n} \end{pmatrix} \left\{ \begin{pmatrix} C & -\mathbb{I}_{2n} \\ K & \mathbf{0} \end{pmatrix} - \lambda \begin{pmatrix} -M & \mathbf{0} \\ \mathbf{0} & -\mathbb{I}_{2n} \end{pmatrix} \right\} \begin{pmatrix} \Pi_M & \mathbf{0} \\ \mathbf{0} & \mathbb{I}_{2n} \end{pmatrix} \\ = & \left(\begin{array}{c|c|c} \mathbf{0} & D\Pi_A & \mathbf{0} & -\mathbb{I}_n \\ \mathbf{0} & Q_A^* B \Pi_A & -Q_A^* & \mathbf{0} \\ \hline -\mathbb{I}_n & \mathbf{0} & & \\ \mathbf{0} & E \Pi_A & & \mathbf{0}_{2n} \end{array} \right) - \lambda \left(\begin{array}{c|c|c} -\mathbb{I}_n & -C \Pi_A & \mathbf{0}_{2n} \\ \mathbf{0} & -R_A & \\ \hline \mathbf{0}_{2n} & & -\mathbb{I}_{2n} \end{array} \right). \quad (11) \end{aligned}$$

The rest of the computation depends on the QZ algorithm. Note that the special structure of the pencil (11) can be exploited for designing a more efficient Hessenberg-triangular decomposition. This is a separate issue that we will not tackle in this work.

3.2.2 Only one matrix is singular

Assume first that E is singular, $r_E < n$, and thus there are at least $n - r_E$ zero eigenvalues which can be deflated. If our setup is to remove only the block of zero eigenvalues that is revealed by the null space of E , then we can achieve that and, at the same time, transform the

⁴Here, some familiarity with the reduction/deflation in the `kvadeig` algorithm is helpful for understanding the details of `kvarteig`.

matrix \mathbb{B} to upper triangular form by the equivalence transformation

$$\begin{aligned} & \begin{pmatrix} Q_M^* & \mathbf{0} \\ \mathbf{0} & Q_K^* \end{pmatrix} \left\{ \begin{pmatrix} \mathbb{C} & -\mathbb{I}_{2n} \\ \mathbb{K} & \mathbf{0} \end{pmatrix} - \lambda \begin{pmatrix} -\mathbb{M} & \mathbf{0} \\ \mathbf{0} & -\mathbb{I}_{2n} \end{pmatrix} \right\} \begin{pmatrix} \Pi_M & \mathbf{0} \\ \mathbf{0} & Q_K \end{pmatrix} \\ = & \left(\begin{array}{c|cc|cc} \mathbf{0} & D\Pi_A & \mathbf{0} & -Q_E \\ \mathbf{0} & Q_A^* B \Pi_A & -Q_A^* & \mathbf{0} \\ \hline -\mathbb{I}_n & \mathbf{0} & & \\ \hline \mathbf{0} & \widehat{R}_E \Pi_E^* \Pi_A & & \mathbf{0}_{2n} \\ \mathbf{0} & \mathbf{0} & & \end{array} \right) - \lambda \left(\begin{array}{c|cc|cc} -\mathbb{I}_n & -C\Pi_A & & \mathbf{0}_{2n} \\ \mathbf{0} & -R_A & & \\ \hline & & \mathbf{0}_{2n} & \\ \hline & & & -\mathbb{I}_{2n} \end{array} \right). \quad (12) \end{aligned}$$

The $n - r_E$ zero eigenvalues are now deflated implicitly by working with the leading $(3n + r_E) \times (3n + r_E)$ sub-pencil of (12). If we want to check for the existence of further blocks corresponding to $\lambda = 0$, then it is convenient to use the following transformation:

$$\begin{aligned} & \begin{pmatrix} Q_K^* & \mathbf{0} \\ \mathbf{0} & Q_K^* \end{pmatrix} \left\{ \begin{pmatrix} \mathbb{C} & -\mathbb{I}_{2n} \\ \mathbb{K} & \mathbf{0} \end{pmatrix} - \lambda \begin{pmatrix} -\mathbb{M} & \mathbf{0} \\ \mathbf{0} & -\mathbb{I}_{2n} \end{pmatrix} \right\} \begin{pmatrix} \mathbb{I}_{2n} & \mathbf{0} \\ \mathbf{0} & Q_K \end{pmatrix} \\ = & \left(\begin{array}{c|cc|cc} B & \mathbf{0} & & -\mathbb{I}_{2n} \\ \hline Q_E^* D & \mathbf{0} & & \\ \hline 0 & -\mathbb{I}_n & & \\ \hline \widehat{R}_E \Pi_E^* & \mathbf{0} & & \mathbf{0}_{2n} \\ \mathbf{0} & \mathbf{0} & & \end{array} \right) - \lambda \left(\begin{array}{c|cc|cc} -A & \mathbf{0} & & \mathbf{0}_{2n} \\ \hline -Q_E^* C & -Q_E^* & & \\ \hline & & \mathbf{0}_{2n} & \\ \hline & & & -\mathbb{I}_{2n} \end{array} \right). \quad (13) \end{aligned}$$

The deflated pencil of order $3n + r_E$ reads

$$\mathbb{A}_{22} - \lambda \mathbb{B}_{22} = \left(\begin{array}{c|cc|cc|cc} B & \mathbf{0} & & -\mathbb{I}_n & & \\ \hline Q_{E,1}^* D & \mathbf{0} & & & -\mathbb{I}_{r_E} & \\ \hline Q_{E,2}^* D & \mathbf{0} & & \mathbf{0} & & \mathbf{0} \\ \hline \mathbf{0} & -\mathbb{I}_n & & & & \\ \hline \widehat{R}_E \Pi_E^* & \mathbf{0} & & & & \mathbf{0}_{n+r_E} \end{array} \right) - \lambda \left(\begin{array}{c|cc|cc} -A & \mathbf{0} & & \mathbf{0}_{(2n) \times (n+r_E)} \\ \hline -Q_E^* C & -Q_E^* & & \\ \hline & & \mathbf{0}_{(n+r_E) \times (2n)} & \\ \hline & & & -\mathbb{I}_{n+r_E} \end{array} \right), \quad (14)$$

where $Q_{E,1}^* = Q_E^*(1 : r_E, \cdot)$ and $Q_{E,2}^* = Q_E^*(r_E + 1 : n, \cdot)$. Note that $\mathbb{A}_{22} - \lambda \mathbb{B}_{22}$ is the block at the position (1, 1) of a block-upper triangular pencil (13); the block position (2, 2) corresponds to the deflated $n - r_E$ zeros. Denote the left and the right transformation matrices from (13) with \mathbf{P}_1 and \mathbf{Q}_1 respectively, and the linearization pencil with $\mathbb{A} - \lambda \mathbb{B} = \mathbb{A}_{11} - \lambda \mathbb{B}_{11}$. After the first deflation step we have⁵

$$\mathbf{P}_1 (\mathbb{A}_{11} - \lambda \mathbb{B}_{11}) \mathbf{Q}_1 = \begin{pmatrix} \mathbb{A}_{22} - \lambda \mathbb{B}_{22} & \spadesuit \\ \mathbf{0} & -\lambda \check{\mathbb{B}}_{11} \end{pmatrix}, \quad \check{\mathbb{B}}_{11} = -\mathbb{I}_{n-r_E}. \quad (15)$$

The next step in the deflation process is to determine the rank of the matrix \mathbb{A}_{22} . From the structure of the matrix, we conclude that the rank of \mathbb{A}_{22} is equal to $2n + r_E +$ "the rank of the $n \times n$ matrix $\left(\frac{Q_{E,2}^* D}{\widehat{R}_E \Pi_E^*} \right)$ ", which is defined in terms of the coefficient matrices D and E of the original problem. So, we compute the rank revealing factorization

$$\left(\frac{Q_{E,2}^* D}{\widehat{R}_E \Pi_E^*} \right) \Pi_{A_{22}} = Q_{A_{22}} R_{A_{22}}. \quad (16)$$

If (16) is of full rank n , then \mathbb{A}_{22} is regular, there are no more zeros in the spectrum, and the single deflation step is done by removing the trailing $n - r_E$ rows and columns in (12).

⁵See [14, §5.2] for more details.

If, on the other hand, (16) is rank deficient with $\text{rank}(R_{A_{22}}) = r_2 < n$, the corresponding number of $n - r_2$ zero eigenvalues can be deflated. To that end, note that $R_{A_{22}} = \begin{pmatrix} \widehat{R}_{A_{22}} \\ \mathbf{0}_{n-r_2, n} \end{pmatrix}$ and transform the pencil (14) to get zero rows at the bottom of \mathbb{A}_{22} . This is done by the permutation $\pi = (1 : n + r_E, 2n + 1 : 3n, n + r_E + 1 : 2n, 3n + 1 : 3n + r_E)$. If Π is the corresponding row permutation matrix, and if we set

$$\widehat{P}_2 = \begin{pmatrix} \mathbb{I}_{2n+r_E} & \\ & Q_{A_{22}}^* \end{pmatrix} \Pi, \quad (17)$$

then the transformed pencil is

$$\widehat{P}_2 \mathbb{A}_{22} = \left(\begin{array}{c|c|c|c} B & \mathbf{0} & -\mathbb{I}_n & \\ \hline Q_{E,1}^* D & \mathbf{0} & & -\mathbb{I}_{r_E} \\ \hline \mathbf{0} & -\mathbb{I}_n & \mathbf{0} & \mathbf{0} \\ \hline \widehat{R}_{A_{22}} \Pi_{A_{22}}^T & \mathbf{0} & & \\ \hline \mathbf{0} & \mathbf{0} & & \mathbf{0}_{n \times (n+r_E)} \end{array} \right), \quad \widehat{P}_2 \mathbb{B}_{22} = \left(\begin{array}{c|c|c|c} -A & \mathbf{0} & \mathbf{0}_{n+r_E} & \\ \hline -Q_{E,1}^* C & -Q_{E,1}^* & \mathbf{0}_{r_E \times (n+r_E)} & \\ \hline \mathbf{0}_n & \mathbf{0}_n & -\mathbb{I}_n & \mathbf{0}_{n \times r_E} \\ \hline -N_{[1]} & -N_{[2]} & N_{[3]} & N_{[4]} \end{array} \right). \quad (18)$$

To deflate the additional $n - r_2$ zeros, we reduce the trailing $n - r_2$ rows of the blocks $-N_{[1]}$, $-N_{[2]}$ and $N_{[3]}$ to zero. This is done by the complete orthogonal decomposition

$$(\widehat{P}_2 \mathbb{B}_{22})(2n + r_E + r_2 + 1 : 3n + r_E, :) = U_{BB} R_{BB} V_{BB}^*, \quad (19)$$

so that $(\widehat{P}_2 \mathbb{B}_{22})(2n + r_E + r_2 + 1 : 3n + r_E, :) V_{BB} = (\mathbf{0} \quad \check{\mathbb{B}}_{22})$. Finally, the deflated pencil is

$$\widehat{P}_2 \mathbb{A}_{22} V_{BB} - \lambda \widehat{P}_2 \mathbb{B}_{22} V_{BB} = \begin{pmatrix} \mathbb{A}_{33} - \lambda \mathbb{B}_{33} & \blacksquare \\ \mathbf{0} & -\lambda \check{\mathbb{B}}_{22} \end{pmatrix}. \quad (20)$$

This reduction process continues by forwarding $\mathbb{A}_{33} - \lambda \mathbb{B}_{33}$ to the next step of reduction toward an upper triangular KCF, as described in [14].

Remark 3. For a more structured backward error in case of graded matrices, the complete orthogonal (URV) decomposition (19) should be computed as in [14, §2.2].

Remark 4. If the matrix A is rank deficient, and E is full rank, we process the reversed problem $(\mu^4 E + \mu^3 D + \mu^2 C + \mu B + A)x = \mathbf{0}$, $\mu = 1/\lambda$, and the corresponding truncated linearization pencil of order $3n + r_A$ reads

$$\mathbb{A}_{22} - \lambda \mathbb{B}_{22} = \left(\begin{array}{c|c|c|c} D & \mathbf{0} & -\mathbb{I}_n & \\ \hline Q_{A,1}^* B & \mathbf{0} & & -\mathbb{I}_{r_A} \\ \hline Q_{A,2}^* B & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \hline \mathbf{0} & -\mathbb{I}_n & & \\ \hline \widehat{R}_A P_A^* & \mathbf{0} & & \mathbf{0}_{n+r_A} \end{array} \right) - \lambda \left(\begin{array}{c|c|c|c} -E & \mathbf{0} & \mathbf{0}_{(2n) \times (n+r_A)} & \\ \hline -Q_A^* C & -Q_A^* & & \\ \hline \mathbf{0}_{(n+r_A) \times (2n)} & & & -\mathbb{I}_{n+r_A} \end{array} \right), \quad (21)$$

and the rank of matrix \mathbb{A}_{22} is now $2n + r_A +$ the rank of the $n \times n$ matrix $\begin{pmatrix} Q_{A,2}^* B \\ \widehat{R}_A P_A^* \end{pmatrix}$.

3.2.3 Both matrices A and E are singular

When both matrices A and E are rank deficient, then, following the discussion from §3.2.2, the key information is in the numerical ranks of the matrices

$$\Phi = \begin{pmatrix} Q_{A,2}^* B \\ \widehat{R}_A \Pi_A^* \end{pmatrix}, \quad \Psi = \begin{pmatrix} Q_{E,2}^* D \\ \widehat{R}_E \Pi_E^* \end{pmatrix}. \quad (22)$$

Both Φ and Ψ are full rank In this case, in the KCF the zero and the infinite eigenvalue occupy single block each, induced by the rank deficiency of E and A . The deflation process starts by creating $n-r_E$ and $n-r_A$ zero rows in the coefficients of the corresponding linearization as follows:

$$\begin{aligned}
& \begin{pmatrix} Q_M^* & \mathbf{0} \\ \mathbf{0} & Q_K^* \end{pmatrix} \left\{ \begin{pmatrix} \mathbb{C} & -\mathbb{I}_{2n} \\ \mathbb{K} & \mathbf{0} \end{pmatrix} - \lambda \begin{pmatrix} -\mathbb{M} & \mathbf{0} \\ \mathbf{0} & -\mathbb{I}_{2n} \end{pmatrix} \right\} \begin{pmatrix} \mathbb{I}_{2n} & \mathbf{0} \\ \mathbf{0} & Q_K \end{pmatrix} \\
= & \left(\begin{array}{c|c|c|c|c} \mathbf{0}_n & D & \mathbf{0}_n & -Q_E(:, 1:r_E) & -Q_E(:, r_E+1:n) \\ \hline \mathbf{0}_{n \times r_A} & Q_A^*(1:r_A, :)B & -Q_A^*(1:r_A, :) & \mathbf{0}_{r_A \times r_E} & \mathbf{0}_{r_A \times (n-r_E)} \\ \mathbf{0}_{n \times (n-r_A)} & Q_A^*(r_A+1:n, :)B & -Q_A^*(r_A+1:n, :) & \mathbf{0}_{(n-r_A) \times r_E} & \mathbf{0}_{(n-r_A) \times (n-r_E)} \\ \hline -\mathbb{I}_n & \mathbf{0}_n & \mathbf{0}_n & \mathbf{0}_{n \times r_E} & \mathbf{0}_{n \times (n-r_E)} \\ \hline \mathbf{0}_{r_E \times n} & \widehat{R}_E \Pi_E^* & \mathbf{0}_{r_E \times n} & \mathbf{0}_{r_E} & \mathbf{0}_{r_E \times (n-r_E)} \\ \mathbf{0}_{(n-r_E) \times n} & \mathbf{0}_{(n-r_E) \times n} & \mathbf{0}_{(n-r_E) \times n} & \mathbf{0}_{(n-r_E) \times r_E} & \mathbf{0}_{(n-r_E)} \end{array} \right) \\
- \lambda & \left(\begin{array}{c|c|c|c|c} -\mathbb{I}_n & C & \mathbf{0}_n & \mathbf{0}_{n \times r_E} & \mathbf{0}_{n \times (n-r_E)} \\ \hline \mathbf{0}_{r_A \times n} & -\widehat{R}_A \Pi_A^* & \mathbf{0}_{r_A \times n} & \mathbf{0}_{r_A \times r_E} & \mathbf{0}_{r_A \times (n-r_E)} \\ \mathbf{0}_{(n-r_A) \times n} & \mathbf{0}_{(n-r_A) \times n} & \mathbf{0}_{(n-r_A) \times n} & \mathbf{0}_{(n-r_A) \times r_E} & \mathbf{0}_{(n-r_A) \times (n-r_E)} \\ \hline \mathbf{0}_n & \mathbf{0}_n & -\mathbb{I}_n & \mathbf{0}_{n \times r_E} & \mathbf{0}_{n \times (n-r_E)} \\ \hline \mathbf{0}_{r_E \times n} & \mathbf{0}_{(n-r_E) \times n} & \mathbf{0}_{r_E \times n} & -\mathbb{I}_{r_E} & \mathbf{0}_{r_E \times (n-r_E)} \\ \mathbf{0}_{(n-r_E) \times n} & \mathbf{0}_{(n-r_E) \times n} & \mathbf{0}_{(n-r_E) \times n} & \mathbf{0}_{(n-r_E) \times r_E} & -\mathbb{I}_{n-r_E} \end{array} \right). \tag{23}
\end{aligned}$$

The next step is to compute the complete orthogonal decomposition

$$\begin{pmatrix} Q_A^*(r_A+1:n, :)B & Q_A^*(r_A+1:n, :) & \mathbf{0}_{(n-r_A) \times r_E} \end{pmatrix} = Q_X \begin{pmatrix} R_X & \mathbf{0}_{(n-r_A) \times (n+r_E+r_A)} \end{pmatrix} Z_X^*, \tag{24}$$

and permute the first $(n-r_A)$ and the last $(n+r_E+r_A)$ columns to get

$$\begin{aligned}
& Q_X^* \begin{pmatrix} Q_A^*(r_A+1:n, :)B & Q_A^*(r_A+1:n, :) & \mathbf{0}_{(n-r_A) \times r_E} \end{pmatrix} Z_X^* \begin{pmatrix} \mathbf{0} & \mathbb{I}_{n-r_E} \\ \mathbb{I}_{n+r_A+r_E} & \mathbf{0} \end{pmatrix} \\
& = \begin{pmatrix} \mathbf{0}_{(n-r_A) \times (n+r_E+r_A)} & R_X \end{pmatrix}.
\end{aligned}$$

Finally, to complete the deflation process, the following left and right transformation matrices must be applied on the pencil (23):

$$\begin{pmatrix} \mathbb{I}_{n+r_A} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbb{I}_{r_E} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbb{I}_n & \mathbf{0} & \mathbf{0} \\ \hline \mathbf{0} & Q_X^* & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbb{I}_{r_E} \end{pmatrix}, \quad \left(\begin{array}{c|c} Z_X^* \begin{pmatrix} \mathbf{0} & \mathbb{I}_{n-r_E} \\ \mathbb{I}_{n+r_A+r_E} & \mathbf{0} \end{pmatrix} & \mathbf{0} \\ \hline \mathbf{0} & \mathbb{I}_{2n-r_E} \end{array} \right).$$

After the transformation step, the deflation is finished by removing the last $2n-r_E-r_A$ rows and columns from the obtained pencil. The resulting pencil of dimension $2n+r_A+r_E$ is forwarded to the QZ algorithm.

Only one matrix in (22) is singular This means that there are at least two KCF blocks for the zero (if Ψ is singular) or the infinite (if Φ is singular) eigenvalue. In either case, we deflate two blocks for the zero eigenvalue using the structure described in §3.2.2 (see also [14, §5.2, §6.1]), meaning that the reversed problem is considered if there are more blocks for the infinite eigenvalues.

After deflating two blocks of zero eigenvalues, we obtain the pencil (20). Now, the existence of additional zero eigenvalues depends on the rank of the matrix \mathbb{A}_{22} . To deflate possible additional zeros, the pencil $\mathbb{A}_{22} - \lambda\mathbb{B}_{22}$ is forwarded to the algorithm for computing the KCF [14, §3.2]. As the output we get the pencil $\mathbb{A}_{\ell+1,\ell+1} - \lambda\mathbb{B}_{\ell+1,\ell+1}$ and transformation matrices Q_p and P_p , with $\mathbb{A}_{\ell+1,\ell+1}$ regular. Denote with $n_{\ell+1}$ the dimension of the resulting pencil.

Finally, we have to deflate one block of infinite eigenvalues, which have been detected at the beginning. This is done by forwarding the reversed pencil $\mathbb{B}_{\ell+1,\ell+1} - \lambda\mathbb{A}_{\ell+1,\ell+1}$ to the procedure described in [14, §3.2]. As the input to the algorithm we supply the information that there is only one block to be deflated, so that only one step of the algorithm is needed. In addition, we also send the number of infinite eigenvalues so that the rank determination of the matrix $\mathbb{B}_{\ell+1,\ell+1}$ is omitted. As an output, we get the pencil $\mathbb{A}_{\ell+1,\ell+1} - \lambda\mathbb{B}_{\ell+1,\ell+1}$ with both $\mathbb{A}_{\ell+1,\ell+1}$ and $\mathbb{B}_{\ell+1,\ell+1}$ regular, and the corresponding transformation matrices P_{p1} and Q_{p1} . The final transformation matrices Q and P are

$$Q = \begin{pmatrix} \mathbb{I}_{2n} & \mathbf{0} \\ \mathbf{0} & Q_{\mathbb{K}} \end{pmatrix} \begin{pmatrix} \mathbb{I}_{2n} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & V_{BB}^* P_{BB} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbb{I}_{n-r_E} \end{pmatrix} \begin{pmatrix} Q_p & \mathbf{0} \\ \mathbf{0} & \mathbb{I}_{4n-r_E-r_2} \end{pmatrix} \begin{pmatrix} Q_{p1} & \mathbf{0} \\ \mathbf{0} & \mathbb{I}_{4n-n_{\ell+1}} \end{pmatrix}$$

$$P = \begin{pmatrix} P_{p1} & \mathbf{0} \\ \mathbf{0} & \mathbb{I}_{4n-n_{\ell+1}} \end{pmatrix} \begin{pmatrix} P_p & \mathbf{0} \\ \mathbb{I}_{4n-r_E-r_2} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbb{I}_{n+r_E} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & Q_{A_{22}}^* & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbb{I}_{2n-r_E} \end{pmatrix} \begin{pmatrix} Q_{\mathbb{K}}^* & \mathbf{0} \\ \mathbf{0} & Q_{\mathbb{K}}^* \end{pmatrix}.$$

Both matrices in (22) are singular This case is analogous to the previous one. The only difference is that, when we call the algorithm on the reversed pencil $B_{\ell+1,\ell+1} - \lambda A_{\ell+1,\ell+1}$, we provide additional information that there are least two steps of deflation ahead, as well as the dimensions of the first two blocks which were previously determined by the rank revealing decompositions of A and Φ .

3.2.4 On making more reduction steps

After all detected zero and/or infinite eigenvalues have been deflated, as described above, we check the ranks of the matrices in the resulting pencil in order to determine whether there are more blocks of these eigenvalues. If these matrices are rank deficient, then another step of deflation must take place. Unfortunately, with that step, the structure of the linearization is lost, so we use the standard deflation process for generalized eigenvalue problem as in [32] and [14]. It remains an interesting problem to determine an equivalence transformation to restore the structure for more steps, while working on an equivalent representation of the original problem.

There is, of course, a trade-off between this increased numerical robustness and computational cost (complexity), and, in a software implementation, the number of reduction/deflation steps will be limited. But, even with these few steps we can make some critical decisions on the zero and infinite eigenvalues, with backward error in terms of the coefficients of the original problem; see §5 and §6.

3.2.5 An illustrative example

Let us illustrate the action of the additional reduction steps toward the KCF. We use the `mirror` example from the NLEVP library [2]; it originates from the calibration of catadioptric vision system [42]. The problem is of order $n = 9$.

Both A and E are rank deficient, with the rank $r_E = r_A = 2$, which means that there are at least 7 zero and 7 infinite eigenvalues. They were correctly identified and deflated in the preprocessing in `quadeig`⁶; in the next step, the QZ algorithm found an additional zero eigenvalue, and two more infinite eigenvalues. On the other hand, `polyeig`⁷ identified in total only 2 zero and 9 infinite eigenvalues. This shows the advantage of the preprocessing introduced in `quadeig` for early revealing of zeros and infinities. These numbers of computed zero and infinite eigenvalues were independent of whether the parameter scaling was on or off before calling `quadeig` and `polyeig`.

On the other hand, the preprocessing in both `kvadeig` and `kvarteig` found additional two zero and two infinite eigenvalues, making the total of 9 zero and 9 infinite eigenvalues deflated before calling the QZ. Again, the same numbers of 9 zero and 9 infinite eigenvalues were found with and without parameter scaling. This almost agrees with the result of `quadeig`, up to one zero eigenvalue.

Next, we check the norm-wise backward errors (4) for all computed eigenpairs (for all four algorithms). The details of computing the eigenvectors in `kvarteig` are given in §4. The computed residuals, shown in Figure 1, seem to indicate that all results are acceptable up to small norm-wise backward errors (separate for each eigenpair) of the order of machine precision. (The eigenvalues are indexed in non-decreasing absolute values.)

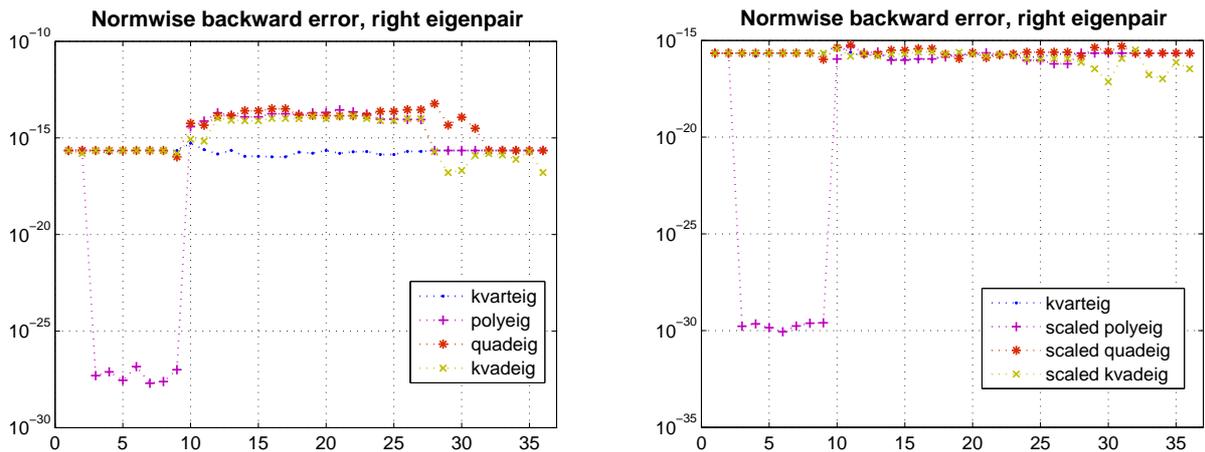


Figure 1: Norm-wise backward errors for all eigenvalues with the corresponding right eigenvectors in the `mirror` NLEVP benchmark example. *Left panel:* No parameter scaling in `polyeig` and the quadratic solvers. *Right panel:* The coefficients of the problem are scaled as described in §1.1.3.

Hence, with all backward errors at the level of the round-off, and with different numbers of zero and infinite eigenvalues computed by different algorithms, how can we tell which one is correct? What assurance is given in a particular algorithm concerning the existence of infinite eigenvalues of a perturbed matrix polynomial in a vicinity of the given one? The difficulty is best illustrated in [32, Example 2], which actually contains the key idea pursued in `quadeig`, `kvadeig` and `kvarteig`.

If we look at the structure of the matrices A and E for this particular problem, we see that their ranks can be determined exactly because each has 7 zero columns, and the independence

⁶For the purpose of testing and comparisons, we apply quadratic solvers to the quadratification (5) of the quartic problem.

⁷We use `polyeig` from Matlab, version 7.11.0.584 (R2010b).

of the remaining columns is easy to check. Further, the block matrices (22), which are used to determine the existence of more than one block for zero and infinite eigenvalues, also have two zero columns each, and the remaining 9×7 submatrices are well conditioned. Thus we can argue that `kvarteig` has determined the correct numbers of zero and infinite eigenvalues.

We call the reader to revisit this example after reading Example 5.

4 Computing the eigenvectors

In the computation of the eigenvectors, we have two main computational tasks: (i) restore the eigenvectors of the quartic problem from the eigenvectors of its *linearization* (§4.1); (ii) assemble the eigenvectors of the linearization from the eigenvectors of the deflated (linearization) pencil, using the transformation matrices (§4.2).

4.1 Quartic eigenvectors from the eigenvectors of the linearization

For an eigenvalue λ , the eigenvectors of the original problem (1) and the final linearization pencil (6) can be related using explicit formulas. For the reader's convenience, we briefly outline the crux of this connection.

We use $z \in \mathbb{C}^{4n}$ and $w \in \mathbb{C}^{4n}$ to denote the right and the left eigenvector for the linearization, and $x \in \mathbb{C}^n$, $y \in \mathbb{C}^n$ to denote the right and the left eigenvector for the original problem. The eigenvalue $\lambda \in \mathbb{C}$ is now fixed as assumed nonzero and finite.

Let $z = (z_1^T \ z_2^T \ z_3^T \ z_4^T)^T \in \mathbb{C}^{4n}$, $z_i \in \mathbb{C}^n$, $i = 1, 2, 3, 4$ be a right eigenvector for the eigenvalue λ ($0 < |\lambda| < \infty$) of the linearized problem, i.e. $(\mathbb{A} - \lambda\mathbb{B})z = \mathbf{0}$:

$$(\mathbb{A} - \lambda\mathbb{B})z = \left\{ \left(\begin{array}{cc|cc} B & \mathbf{0} & -\mathbb{I} & \mathbf{0} \\ D & \mathbf{0} & \mathbf{0} & -\mathbb{I} \\ \hline \mathbf{0} & -\mathbb{I} & \mathbf{0} & \mathbf{0} \\ E & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{array} \right) - \lambda \left(\begin{array}{cc|cc} -A & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ -C & -\mathbb{I} & \mathbf{0} & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{0} & -\mathbb{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & -\mathbb{I} \end{array} \right) \right\} \begin{pmatrix} z_1 \\ z_2 \\ z_3 \\ z_4 \end{pmatrix} = \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \\ \mathbf{0} \\ \mathbf{0} \end{pmatrix}. \quad (25)$$

By equating the corresponding block components on the left and on the right we get

$$Bz_1 - z_3 + \lambda Az_1 = \mathbf{0} \Leftrightarrow z_3 = (\lambda A + B)z_1, \quad (26)$$

$$Dz_1 - z_4 + \lambda Cz_1 + \lambda z_2 = \mathbf{0} \Leftrightarrow Dz_1 + (1/\lambda)Ez_1 + \lambda Cz_1 + \lambda^2(\lambda A + B)z_1 = \mathbf{0}, \quad (27)$$

$$-z_2 + \lambda z_3 = \mathbf{0} \Leftrightarrow z_2 = \lambda z_3, \quad (28)$$

$$Ez_1 + \lambda z_4 = \mathbf{0} \Leftrightarrow \lambda z_4 = -Ez_1. \quad (29)$$

It follows immediately that $z_1 \neq \mathbf{0}$; if $\det(\lambda A + B) \neq 0$, then, in addition, $z_3 \neq \mathbf{0}$ and $z_2 \neq \mathbf{0}$; if $\det(E) \neq 0$, then also $z_4 \neq \mathbf{0}$. Using (27) we easily check that $x = z_1/\lambda$ is an eigenvector of the original quartic problem. Further, (26) implies that x satisfies $z_3 = \lambda(\lambda A + B)x$, and (28) yields $z_2 = \lambda^2(\lambda A + B)x$, and finally from (29) it follows that $z_4 = -Ex$. Similarly, if we initially assume that x is an eigenvector of the quartic problem, these formulas for the z_i 's give an eigenvector of (25).

An analogous computation reveals a left eigenvector y , using the partitioned left eigenvector of the linearization, as $w = (w_1^T \ w_2^T \ w_3^T \ w_4^T)^T$, $w_i \in \mathbb{C}^n$, $i = 1, 2, 3, 4$. Altogether, we obtain the following relations between the two sets of eigenvectors:

$$z = \begin{pmatrix} z_1 \\ z_2 \\ z_3 \\ z_4 \end{pmatrix} = \begin{pmatrix} \lambda x \\ \lambda^2(\lambda A + B)x \\ \lambda(\lambda A + B)x \\ -Ex \end{pmatrix}, \quad w = \begin{pmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \end{pmatrix} = \begin{pmatrix} \lambda^3 y \\ \lambda^2 y \\ \lambda y \\ y \end{pmatrix}. \quad (30)$$

For both the right and the left eigenvector there are four choices to recover x and y . Reconstruction of the left eigenvector seems easier. We just choose one of the block components w_1, w_2, w_3 or w_4 and rescale appropriately.

For the right eigenvector we can choose $z_1, (\lambda A + B)^{-1}z_2, (\lambda A + B)^{-1}z_3$ or $E^{-1}z_4$. Notice that, for the last three choices we have to solve system of linear equations in order to compute the wanted vector. Given all the difficulties in numerical solution of nonlinear eigenvalue problem, we ought to use all alternatives in order to obtain better output – in this case, for instance, we can solve all systems and select the vector with smallest residual.

Remark 5. *If $\lambda = 0$, for the corresponding right eigenvector we have $Ex = \mathbf{0}$ and $\mathbb{A}x = \mathbf{0}$. By the same reasoning as above, we conclude the following connection $z = \begin{pmatrix} x^T & \mathbf{0} & (Bx)^T & (Dx)^T \end{pmatrix}^T$.*

4.1.1 Computing $(\lambda A + B)^{-1}z_2, (\lambda A + B)^{-1}z_3$ or $E^{-1}z_4$ multiple times

Inverting E (assuming $\det(E) \neq 0$) multiple times can be done by reusing initially computed LU decomposition. On the other hand computing $(\lambda A + B)^{-1}z_2, (\lambda A + B)^{-1}z_3$ for $4n$ values of λ is not that simple because the coefficients of the linear system change with λ ; $O(n^3)$ flops per eigenvalue to compute the corresponding eigenvector is prohibitive complexity. Fortunately, this can be reduced using a bag of tricks for solving shifted linear systems. In particular, this problem is similar to evaluating the transfer function of a descriptor LTI dynamical system at multiple frequencies [24, §4].

We can compute the triangular-Hessenberg form of (A, B) , i.e. a unitary Q , an upper triangular T and an upper Hessenberg matrix H can be constructed in $O(n^3)$ time so that $A = QTQ^*$, $B = QHQ^*$. Hence, for any vector v

$$(\lambda A + B)^{-1}v = Q[(\lambda T + H)^{-1}(Q^*v)],$$

which has $O(n^2)$ complexity because $\lambda T + H$ is upper Hessenberg. This means that the total work (for all $4n$ eigenvalues) of choosing the eigenvectors with smallest residuals remains $O(n^3)$. (Here, the tacit assumption is that $A + \lambda B$ is nonsingular and well conditioned with respect to inversion.) These details can be taken into account for a development of an optimized software for multicore architectures; for more information see [3], [4].

Remark 6. *In some applications, such as e.g. computing deformation modes of thin-walled structures [35], the cubic term is zero, $B = \mathbf{0}$, so that the shifted systems can be replaced with linear system matrix A for all λ 's. Other details include e.g. the case of real data and using the complex conjugate eigenpairs to save unnecessary computation. Here we omit those details and leave them for the detailed description of a software implementation, which is a subject of our future work.*

4.1.2 Least squares reconstruction of the eigenvectors

Since in a finite precision computation the computed eigenvector z is only an approximation (thus noisy), and since $A + \lambda B$ is not guaranteed to be well conditioned, it makes sense to turn the conditions (30) into a least squares problem, but keeping in mind than we may have to solve it $4n$ times (i.e. we may take e.g. only two conditions to form the least squares problem). So, for instance, we can compute x by solving the least squares problem

$$\left\| \begin{pmatrix} \lambda \mathbb{I}_n \\ E \end{pmatrix} x - \begin{pmatrix} z_1 \\ -z_4 \end{pmatrix} \right\|_2 \longrightarrow \min \quad (\text{or e.g. } \left\| \begin{pmatrix} \mathbb{I}_n \\ E \end{pmatrix} x - \begin{pmatrix} z_1/\lambda \\ -z_4 \end{pmatrix} \right\|_2 \longrightarrow \min) \quad (31)$$

Actually, the second (and any additional) condition serves as a regularization that can be given a positive weight. Such a strategy can be used for a deeper study of selected eigenpairs. We omit the details for the sake of brevity.

If the data is well scaled, then for some eigenvalues (semi-)normal equations can be used. In general, the least squares problem (31) can be solved efficiently for any eigenvalue $\lambda \neq 0$ by pre-computing the SVD $E = U_E \Sigma_E V_E^*$ (which actually may be available if we used it for a strong rank revealing of E) and then, for each triple λ, z_1, z_4 , solving in $O(n^2)$ flops the equivalent problem

$$\left\| \begin{pmatrix} \lambda \mathbb{I}_n \\ \Sigma_E \end{pmatrix} V_E^* x - \begin{pmatrix} V_E^* z_1 \\ -U_E^* z_4 \end{pmatrix} \right\|_2 \longrightarrow \min. \quad (\text{or} \quad \left\| \begin{pmatrix} \mathbb{I}_n \\ \Sigma_E \end{pmatrix} V_E^* x - \begin{pmatrix} V_E^* z_1 / \lambda \\ -U_E^* z_4 \end{pmatrix} \right\|_2 \longrightarrow \min.) \quad (32)$$

If $\lambda = 0$, then, based on Remark 5, the corresponding eigenvector can be found from either of the following least squares problems

$$\left\| \begin{pmatrix} \mathbb{I}_n \\ B \end{pmatrix} x - \begin{pmatrix} z_1 \\ z_3 \end{pmatrix} \right\|_2 \longrightarrow \min, \quad \left\| \begin{pmatrix} \mathbb{I}_n \\ D \end{pmatrix} x - \begin{pmatrix} z_1 \\ z_4 \end{pmatrix} \right\|_2 \longrightarrow \min, \quad (33)$$

which can be efficiently solved for all eigenvectors z of $\lambda = 0$, using one of the approaches discussed above. Other possibilities include e.g. using the bidiagonalization instead of the SVD (of B or D).

4.2 Assembling the eigenvectors of the linearization

Let \tilde{z} and \tilde{w} be the computed right and left eigenvector for the linearization pencil (6). Both right and left eigenvectors will have $4n$ elements if no deflation occurred, otherwise the number of elements will be $4n - d$, where d is the total number of zero and infinite eigenvalues deflated. $4n - d$ is also the dimension of the truncated pencil $\tilde{\mathbb{A}}_{22} - \lambda \tilde{\mathbb{B}}_{22} = P(\mathbb{A} - \lambda \mathbb{B})Q$ which is passed to the QZ algorithm for computation of finite nonzero eigenvalues.

4.2.1 Case 1: No deflation has occurred

Let \tilde{z} and \tilde{w} be the right and the left eigenvector of the transformed pencil $P(\mathbb{A} - \lambda \mathbb{B})Q$. The corresponding right and the left eigenvectors for the original linearization pencil are $z = Q\tilde{z}$ and $w = P^T \tilde{w}$. The right and the left eigenvector for the quartic problem are computed as described in §4.1.

4.2.2 Case 2: Deflation has occurred

Let $n_{\ell+1}$ be the dimension of the deflated linearization $A_{\ell+1, \ell+1} - \lambda B_{\ell+1, \ell+1}$, i.e. both $A_{\ell+1, \ell+1}$ and $B_{\ell+1, \ell+1}$ are regular. Let $\tilde{z} \in \mathbb{C}^{n_{\ell+1}}$ and $\tilde{w} \in \mathbb{C}^{n_{\ell+1}}$ be the right and the left eigenvector for a finite nonzero eigenvalue λ .

To recover eigenvectors of the initial linearization, we must lift \tilde{z} and \tilde{w} to the $4n$ -dimensional space. For the right eigenvector this is easy; we just append $4n - n_{\ell+1}$ zeros to \tilde{z} to get $z = Q \begin{pmatrix} \tilde{z}^T & \mathbf{0}_{1 \times (4n - n_{\ell+1})} \end{pmatrix}^T$.

For the left eigenvector, let $\tilde{w}_2 \in \mathbb{C}^{n - n_{\ell+1}}$ be the vector satisfying $(\tilde{w}^T \quad \tilde{w}_2^T) P(\mathbb{A} - \lambda \mathbb{B})Q = \mathbf{0}$. From

$$(\tilde{w}^T \quad \tilde{w}_2^T) P(\mathbb{A} - \lambda \mathbb{B})Q = (\tilde{w}^T \quad \tilde{w}_2^T) \begin{pmatrix} A_{\ell+1, \ell+1} - \lambda B_{\ell+1, \ell+1} & X \\ \mathbf{0} & Y \end{pmatrix} \quad (34)$$

we conclude that $\tilde{w}_2^T = -\tilde{w}^T XY^{-1}$. (Note that it follows from the procedure in §3.2.2 that Y is nonsingular. Namely, Y is a block upper triangular matrix with the diagonal blocks that are by construction nonsingular.) Now, the left eigenvector for the original linearization is $w = P^T (\tilde{w}^T \quad \tilde{w}_2^T)^T$.

The right eigenvectors for the zero eigenvalue span the nullspace of the matrix E . The corresponding basis is computed for the orthogonal complement of the range of E^* . To compute this basis we can use the already computed QR factorization of E (7) as follows. First compute the QR factorization of $\Pi_E \hat{R}_E^* = Q_{\hat{R}_E^*} R_{\hat{R}_E^*}$. Now, the last $n - r_E$ columns of $Q_{\hat{R}_E^*}$ represent the basis for the nullspace of the matrix E . Similarly, the right eigenvectors of the infinite eigenvalue span the nullspace of the matrix A . The basis is computed using the already computed QR factorization (7). Again, compute the QR factorization of $\Pi_A \hat{R}_A^* = Q_{\hat{R}_A^*} R_{\hat{R}_A^*}$, and the last $n - r_A$ columns of $Q_{\hat{R}_A^*}$ represent the basis for the nullspace of A .

The left eigenvectors for the zero eigenvalue are determined as the last $n - r_E$ columns of the unitary matrix Q_E from the corresponding QR factorization, and the left eigenvectors for the infinite eigenvalue are selected as the last $n - r_A$ columns of the unitary matrix Q_A from the QR factorization of A .

5 Backward error analysis

As we discussed in §1.1 and §1.2, it is important that the computational errors in the pre-processing phase correspond (in a backward error sense) to small perturbations of the initial coefficients, i.e. that we have strong norm-wise backward stability. Moreover, if the initial coefficient matrices have graded columns, then it is advantageous to have backward error in each column to be small relative to the size of that column (instead of relative to the norm of the whole matrix). In this section, we analyze the backward errors in the proposed preprocessor, where we use a framework of mixed error analysis. As expected, such an analysis is rather technical.

We provide a backward/mixed error analysis for the first two steps of the deflation procedure described in §3.2.2 (only one of A and E is singular), which includes the key details and principles of the analysis. Extending this further to the first two steps of the general case §3.2.3 (both A and E singular) would follow the same steps and it is omitted for the sake of brevity.

The following proposition deals with the first step, that is, the deflation of the first batch of $n - r_E$ zero eigenvalues.

Proposition 1. *Let $E\tilde{\Pi}_E \approx \tilde{Q}_E \begin{pmatrix} R_E \\ \mathbf{0} \end{pmatrix}$ be the computed rank revealing QR factorization of E , and let \tilde{r}_E be the computed numerical rank of E . Further, let $\tilde{X} = \text{computed}(\tilde{Q}_E^* D)$, $\tilde{Y} = \text{computed}(\tilde{Q}_E^* C)$. Let*

$$\tilde{\mathbb{A}}_{22} - \lambda \tilde{\mathbb{B}}_{22} = \left(\begin{array}{cc|cc} B & \mathbf{0} & -\mathbb{I}_n & \mathbf{0} \\ \tilde{X} & \mathbf{0} & \mathbf{0} & -\mathbb{I}_{\tilde{r}_E} \\ \hline \mathbf{0} & -\mathbb{I}_n & & \\ \tilde{R}_E \tilde{\Pi}_E^T & \mathbf{0} & & \mathbf{0}_{n+\tilde{r}_E} \end{array} \right) - \lambda \left(\begin{array}{cc|c} -A & \mathbf{0} & \mathbf{0}_{2n \times (n+\tilde{r}_E)} \\ -\tilde{Y} & -\tilde{Q}_E^* & \\ \hline \mathbf{0}_{(n+\tilde{r}_E) \times 2n} & & -\mathbb{I}_{n+\tilde{r}_E} \end{array} \right) \quad (35)$$

be the computed reduced pencil (14), extracted from the transformed linearization (13). There

exists small structured perturbation

$$\delta\tilde{\mathbb{B}}_{22} = \left(\begin{array}{c|c|c} \mathbf{0} & \mathbf{0} & \\ \hline \mathbf{0} & -\delta Q_E^* & \mathbf{0}_{2n \times (n+\tilde{r}_E)} \\ \hline \mathbf{0}_{(n+\tilde{r}_E) \times 2n} & & \mathbf{0}_{n+\tilde{r}_E} \end{array} \right),$$

such that $\tilde{\mathbb{A}}_{22} - \lambda(\tilde{\mathbb{B}}_{22} + \delta\tilde{\mathbb{B}}_{22})$ corresponds to an exact reduced linearization of a quartic matrix polynomial

$$\lambda^4 A + \lambda^3 B + \lambda^2(C + \delta C) + \lambda(D + \delta D) + (E + \delta E + \Delta E)$$

with at least $n - \tilde{r}_E$ zero eigenvalues, where, for all $i = 1, \dots, n$,

$$\|\delta C(:, i)\|_2 \leq \epsilon_C \|C(:, i)\|_2, \quad \|\delta D(:, i)\|_2 \leq \epsilon_D \|D(:, i)\|_2, \quad \|\delta E(:, i)\|_2 \leq \epsilon_{qr} \|E(:, i)\|_2, \quad (36)$$

and the truncation error from the determination of the numerical rank of E is⁸

$$\max_{j=1:n-k} \|(\Delta E)\tilde{\Pi}_E(:, k+j)\|_2 \leq \tau \min_{i=1:k} \|(E + \delta E)\tilde{\Pi}_E(:, i)\|_2; \quad (\Delta E)\tilde{\Pi}_E(:, 1:k) = \mathbf{0}_{n,k}. \quad (37)$$

Here $\epsilon_C, \epsilon_D, \epsilon_{qr}$ are bounded by a moderate function of n times the machine precision ϵ , and τ is prescribed threshold parameter.

Proof. The computed QR factorization of E , $E\tilde{\Pi}_E \approx \tilde{Q}_E \begin{pmatrix} \tilde{R}_E \\ \mathbf{0} \end{pmatrix}$ can be represented as $(E + \delta E + \Delta E)\tilde{\Pi}_E = \hat{Q}_E \begin{pmatrix} \hat{R}_E \\ \mathbf{0} \end{pmatrix}$, where \hat{Q}_E is exactly unitary and $\|\tilde{Q}_E - \hat{Q}_E\|_F \leq \epsilon_{qr}$; the backward error δE is induced by rounding errors during the factorization, and ΔE is the truncation error from the numerical rank. If we set $\delta Q_E = \tilde{Q}_E - \hat{Q}_E$, then $\tilde{Q}_E = \hat{Q}_E(\mathbb{I}_n + \hat{Q}_E^* \delta Q_E) = (\mathbb{I}_n + \delta Q_E \hat{Q}_E^*)\hat{Q}_E$. We can also write $(E + \delta_1 E + \Delta E)\tilde{\Pi}_E = \tilde{Q}_E \begin{pmatrix} \tilde{R}_E \\ \mathbf{0} \end{pmatrix}$, where $\delta_1 E = \delta E + \delta Q_E \begin{pmatrix} \tilde{R}_E \\ \mathbf{0} \end{pmatrix}$, and thus $\begin{pmatrix} \tilde{R}_E \\ \mathbf{0} \end{pmatrix} \tilde{\Pi}_E^T = \hat{Q}_E^* (E + \delta E + \Delta E) = \tilde{Q}_E^{-1} (E + \delta_1 E + \Delta E)$.

There is an important subtlety here, and it is instructive to discuss it in more detail. In the actually computed matrix $\tilde{\mathbb{B}}_{22}$, stored in the computer memory, one of its blocks is the numerically computed numerically orthogonal \tilde{Q}_E . The backward stability of the QR factorization is usually stated in terms of an exactly unitary matrix \hat{Q}_E , which is an inaccessible object as it is artificially constructed in the proof of backward stability. This is motivated by the desire to be able to say that we have computed the exact QR factorization of a nearby matrix. The matrices \tilde{X} and \tilde{Y} are computed by using the floating point matrix \tilde{Q}_E , possibly implicitly as in the LAPACK subroutine `xORMQR`, or by explicit matrix multiply (`xGEMM` from BLAS) using explicitly formed \tilde{Q}_E , using `xORGQR` (LAPACK). The computed \tilde{X}, \tilde{Y} can be represented as

$$\begin{aligned} \text{computed}(\tilde{Q}_E^* D) &= \tilde{Q}_E^* D + \delta_0 D = \tilde{Q}_E^* (D + \delta D), \quad \delta D = \tilde{Q}_E^{-*} \delta_0 D, \quad |\delta_0 D| \leq \epsilon |\tilde{Q}_E^*| |D|, \\ &= \hat{Q}_E^* (D + \Delta D), \quad \Delta D = \delta D + \hat{Q}_E (\delta Q_E)^* D + \hat{Q}_E (\delta Q_E)^* \delta D, \quad \epsilon \leq O(n)\epsilon; \\ \text{computed}(\tilde{Q}_E^* C) &= \tilde{Q}_E^* C + \delta_0 C = \tilde{Q}_E^* (C + \delta C), \quad \delta C = \tilde{Q}_E^{-*} \delta_0 C, \quad |\delta_0 C| \leq \epsilon |\tilde{Q}_E^*| |C|, \\ &= \hat{Q}_E^* (C + \Delta C), \quad \Delta C = \delta C + \hat{Q}_E (\delta Q_E)^* C + \hat{Q}_E (\delta Q_E)^* \delta C. \end{aligned}$$

On the other hand, the unit blocks $\mathbb{I}_n \oplus \mathbb{I}_{\tilde{r}_E}$ in $\tilde{\mathbb{A}}_{22}$ and $\mathbb{I}_{n+\tilde{r}_E}$ in $\tilde{\mathbb{B}}_{22}$ assume exact orthogonality of \tilde{Q}_E , which is not feasible in finite precision arithmetic. If we set $\Delta_{\Sigma_1} E = \delta_1 E + \Delta E$, then

⁸See Remark 2.

we can represent the computed linearization (13) as

$$\begin{aligned} & \begin{pmatrix} \mathbb{I}_n & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \tilde{Q}_E^* & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbb{I}_n & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \tilde{Q}_E^{-1} \end{pmatrix} \left\{ \begin{pmatrix} B & \mathbf{0}_n & -\mathbb{I}_n & \mathbf{0} \\ D+\delta D & \mathbf{0}_n & \mathbf{0}_n & -\mathbb{I}_n \\ 0 & -\mathbb{I}_n & \mathbf{0}_n & \mathbf{0}_n \\ E+\Delta_{\Sigma_1} E & \mathbf{0}_n & \mathbf{0}_n & \mathbf{0}_n \end{pmatrix} - \lambda \begin{pmatrix} -A & \mathbf{0}_n & \mathbf{0}_n & \mathbf{0}_n \\ -(C+\delta C) & -\mathbb{I}_n & \mathbf{0}_n & \mathbf{0}_n \\ \mathbf{0}_n & \mathbf{0}_n & -\mathbb{I}_n & \mathbf{0}_n \\ \mathbf{0}_n & \mathbf{0}_n & \mathbf{0}_n & -\mathbb{I}_n \end{pmatrix} \right\} \begin{pmatrix} \mathbb{I}_n & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbb{I}_n & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbb{I}_n & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \tilde{Q}_E \end{pmatrix} \\ &= \left(\begin{array}{cc|cc} \frac{B}{\tilde{X}} & \mathbf{0} & \frac{-\mathbb{I}_n}{\mathbf{0}} & \frac{\mathbf{0}}{-\mathbb{I}_n} \\ \hline 0 & -\mathbb{I}_n & & \\ \hline \frac{R_E \tilde{\Pi}_E^T}{\mathbf{0}} & \mathbf{0} & & \mathbf{0}_{2n} \end{array} \right) + \left(\begin{array}{cc|cc} \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{0} & \mathbf{0} & \Xi \\ \hline \mathbf{0} & \mathbf{0} & & \mathbf{0}_{2n} \\ \hline \mathbf{0} & \mathbf{0} & & \mathbf{0}_{2n} \end{array} \right) - \lambda \left(\begin{array}{cc|cc} -A & \mathbf{0} & & \\ \hline -\tilde{Y} & -\tilde{Q}_E^* & & \mathbf{0}_{2n} \\ \hline \mathbf{0}_{2n} & & & -\mathbb{I}_{2n} \end{array} \right). \end{aligned}$$

Now we see at the block position (2,4) in the left matrix, $\tilde{Q}_E^*(-\mathbb{I}_n)\tilde{Q}_E = -\mathbb{I}_n + \Xi \neq -\mathbb{I}_n$. Hence, (35) can be justified by a mixed stability scenario – if the computed pencil is changed by $\|\Xi\|_2 \leq \epsilon_{qr}$ to restore identity at the (2,4) position in the left matrix, then it can be interpreted as an exact transformation of a slightly changed initial pencil.

Alternatively, we can set $\Delta_{\Sigma} E = \delta E + \Delta E$ and model (13) as

$$\begin{aligned} & \begin{pmatrix} \mathbb{I}_n & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \tilde{Q}_E^* & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbb{I}_n & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \tilde{Q}_E \end{pmatrix} \left\{ \begin{pmatrix} B & \mathbf{0}_n & -\mathbb{I}_n & \mathbf{0} \\ D+\Delta D & \mathbf{0}_n & \mathbf{0}_n & -\mathbb{I}_n \\ 0 & -\mathbb{I}_n & \mathbf{0}_n & \mathbf{0}_n \\ E+\Delta_{\Sigma} E & \mathbf{0}_n & \mathbf{0}_n & \mathbf{0}_n \end{pmatrix} - \lambda \begin{pmatrix} -A & \mathbf{0}_n & \mathbf{0}_n & \mathbf{0}_n \\ -(C+\Delta C) & -\mathbb{I}_n & \mathbf{0}_n & \mathbf{0}_n \\ \mathbf{0}_n & \mathbf{0}_n & -\mathbb{I}_n & \mathbf{0}_n \\ \mathbf{0}_n & \mathbf{0}_n & \mathbf{0}_n & -\mathbb{I}_n \end{pmatrix} \right\} \begin{pmatrix} \mathbb{I}_n & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbb{I}_n & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbb{I}_n & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \tilde{Q}_E \end{pmatrix} \\ &= \left(\begin{array}{cc|cc} \frac{B}{\tilde{X}} & \mathbf{0} & \frac{-\mathbb{I}_n}{\mathbf{0}} & \frac{\mathbf{0}}{-\mathbb{I}_n} \\ \hline 0 & -\mathbb{I}_n & & \\ \hline \frac{R_E \tilde{\Pi}_E^T}{\mathbf{0}} & \mathbf{0} & & \mathbf{0}_{2n} \end{array} \right) - \lambda \left\{ \left(\begin{array}{cc|cc} -A & \mathbf{0} & & \\ \hline -\tilde{Y} & -\tilde{Q}_E^* & & \mathbf{0}_{2n} \\ \hline \mathbf{0}_{2n} & & & -\mathbb{I}_{2n} \end{array} \right) + \left(\begin{array}{cc|cc} \mathbf{0} & \mathbf{0} & & \\ \hline \mathbf{0} & -\delta Q_E^* & & \mathbf{0}_{2n} \\ \hline \mathbf{0}_{2n} & & & -\mathbf{0}_{2n} \end{array} \right) \right\} \quad (38) \end{aligned}$$

In this case, the (2,2) block in the right matrix in (35) should be changed from $-\tilde{Q}_E^*$ to $-\hat{Q}_E^*$, by adding δQ_E^* , to establish exact equivalence with a slightly perturbed initial pencil. \square

Remark 7. *The forward error introduced in (38) (thus making the model of the analysis of mixed forward-backward type) is due to the fact that in finite precision computation unitarity/orthogonality cannot be guaranteed.⁹ Note that this error is localized to one block of the linearization; its structure can be easily seen from the backward analysis of the e.g. Householder QR factorization.*

We now consider the first two steps and show that the algorithm remains mixed stable. The proof is technically more involved, but it is important to see how the reduced linear pencil after small forward modification exactly corresponds to a quartic pencil with backward errors in the initial coefficient matrices. Also, the proof nicely illustrates the benefits of well scaled data.

Theorem 1. *Assume the notation of Proposition 1, and let*

$$\widetilde{\hat{P}_2 A_{22}} = \left(\begin{array}{cc|cc} B & \mathbf{0} & \frac{-\mathbb{I}_n}{\mathbf{0}} & \frac{\mathbf{0}}{-\mathbb{I}_{\tilde{r}_E}} \\ \hline \hat{Q}_{E,1}^*(D+\Delta D) & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \hline \mathbf{0} & -\mathbb{I}_n & & \\ \hline \frac{R_{A_{22}} \tilde{\Pi}_{A_{22}}^T}{\mathbf{0}} & \mathbf{0} & & \frac{\mathbf{0}_{\tilde{r}_2 \times (n+\tilde{r}_E)}}{\mathbf{0}_{(n-\tilde{r}_2) \times (n+\tilde{r}_E)}} \end{array} \right), \quad (39)$$

⁹For that reason the QR factorization can only be mixed stable, and in general it is not backward stable.

$$\widetilde{\widehat{P}_2\mathbb{B}_{22}} = \left(\begin{array}{c|c|c} -A & \mathbf{0} & \mathbf{0}_{n+\tilde{r}_E} \\ \hline -\widehat{Q}_{E,1}^*(C + \Delta C) & -\widehat{Q}_{E,1}^* & \frac{\mathbf{0}_{\tilde{r}_E \times (n+\tilde{r}_E)}}{-\mathbb{I}_n} \\ \hline \mathbf{0}_n & \mathbf{0}_n & \mathbf{0}_{n \times \tilde{r}_E} \\ \hline -\widetilde{N}_{[1]} & -\widetilde{N}_{[2]} & \widetilde{N}_{[3]} \mid \widetilde{N}_{[4]} \end{array} \right) \quad (40)$$

be the computed version of (18). There exist small structured forward perturbation

$$\mathcal{F}_{\mathbb{B}_{22}} = \left(\begin{array}{c|c|c} \mathbf{0} & \mathbf{0} & \mathbf{0}_{n+\tilde{r}_E} \\ \hline \mathbf{0} & -\delta\widehat{Q}_{E,1}^* & \frac{\mathbf{0}_{\tilde{r}_E \times (n+\tilde{r}_E)}}{\mathbf{0}_n \mid \mathbf{0}_{n \times \tilde{r}_E}} \\ \hline \mathbf{0}_n & \mathbf{0}_n & \mathbf{0}_{n \times \tilde{r}_E} \\ \hline \Delta\widetilde{N}_{[1]} & \Delta\widetilde{N}_{[2]} & \mathbf{0} \mid \Delta\widetilde{N}_{[4]} \end{array} \right), \quad \|\delta\widetilde{Q}_{E,1}^*\|_2 \leq \epsilon_{qr}, \quad \|\Delta\widetilde{N}_{[4]}\|_2 \leq \epsilon_{qr},$$

of $\widehat{P}_2\mathbb{B}_{22}$, and backward errors ΔC , $\Delta_\Sigma D$, $\Delta_\Sigma E$ such that $\widetilde{\widehat{P}_2\mathbb{A}_{22}} - \lambda(\widetilde{\widehat{P}_2\mathbb{B}_{22}} + \mathcal{F}_{\mathbb{B}_{22}})$ corresponds to an exactly reduced linearization of a quartic matrix polynomial

$$\lambda^4 A + \lambda^3 B + \lambda^2(C + \Delta C) + \lambda(D + \Delta_\Sigma D) + (E + \Delta_\Sigma E),$$

with the exact transformation given in (46) and (53) below. Under mild technical assumption (on the size of $n\epsilon$), we have, for each column index i ,

$$\|\Delta\widetilde{N}_{[1]}(:, i)\|_2 \leq f_1(n)\epsilon\|\widetilde{N}_{[1]}(:, i)\|_2, \quad \|\Delta\widetilde{N}_{[2]}\|_2 \leq f_2(n)\epsilon\|\widetilde{N}_{[2]}\|_2, \quad (41)$$

where $f_1(n)$, $f_2(n)$ are mildly growing functions. Further, with ΔC , ΔD , δE , ΔE as in Proposition 1, it holds that $\Delta_\Sigma D = \Delta D + \widehat{Q}_{E,2}\Gamma_1$ and $\Delta_\Sigma E = \delta E + \Delta E + \widehat{Q}_{E,1}\Gamma_2$, where

$$\left\| \begin{pmatrix} \Gamma_1 \\ \Gamma_2 \end{pmatrix}(:, i) \right\|_2 \leq \epsilon_{qr} \left\| \begin{pmatrix} (D + \Delta D) \\ E + \delta E + \Delta E \end{pmatrix}(:, i) \right\|_2. \quad (42)$$

The latter shows the benefits of well scaled and balanced D and E (§1.1.3, §1.1.4, Remark 1).

Proof. We continue based on the details and the notation from the proof of Proposition 1. The next step is computation of the rank revealing factorization of the block matrix $\begin{pmatrix} \widehat{Q}_{E,2}^*(D + \delta D) \\ \widehat{R}_E \widetilde{\Pi}_E^T \end{pmatrix}$. It is convenient to consider the left matrix in (35) with the relevant blocks already swapped (see (18))

$$\left(\begin{array}{c|c|c} B & \mathbf{0} & -\mathbb{I}_n \mid \mathbf{0} \\ \hline \widehat{Q}_{E,1}^*(D + \Delta D) & \mathbf{0} & \mathbf{0} \mid -\mathbb{I}_{\tilde{r}_E} \\ \widehat{Q}_{E,2}^*(D + \Delta D) & \mathbf{0} & \mathbf{0} \mid \mathbf{0} \\ \hline \mathbf{0}_n & -\mathbb{I}_n & \mathbf{0}_{n+\tilde{r}_E} \\ \hline \widehat{R}_E \widetilde{\Pi}_E^T & \mathbf{0} & \end{array} \right) \rightarrow \left(\begin{array}{c|c|c} B & \mathbf{0} & -\mathbb{I}_n \mid \mathbf{0} \\ \hline \widehat{Q}_{E,1}^*(D + \Delta D) & \mathbf{0} & \mathbf{0} \mid -\mathbb{I}_{\tilde{r}_E} \\ \mathbf{0} & -\mathbb{I}_n & \mathbf{0} \mid \mathbf{0} \\ \hline \widehat{Q}_{E,2}^*(D + \Delta D) & \mathbf{0} & \mathbf{0}_{n \times (n+\tilde{r}_E)} \\ \hline \widehat{R}_E \widetilde{\Pi}_E^T & \mathbf{0} & \end{array} \right).$$

For the computed factors $\widetilde{\Pi}_{A_{22}}$, $\widetilde{Q}_{A_{22}}$, $\widetilde{R}_{A_{22}}$ it holds that

$$\left[\begin{pmatrix} \widehat{Q}_{E,2}^*(D + \delta D) \\ \widehat{R}_E \widetilde{\Pi}_E^T \end{pmatrix} + \begin{pmatrix} \Gamma_1 \\ \Gamma_2 \end{pmatrix} \right] \widetilde{\Pi}_{A_{22}} \equiv \left[\begin{pmatrix} \widehat{Q}_{E,2}^*(D + \Delta D) \\ \widehat{R}_E \widetilde{\Pi}_E^T \end{pmatrix} + \begin{pmatrix} \Gamma_1 \\ \Gamma_2 \end{pmatrix} \right] \widetilde{\Pi}_{A_{22}} = \widehat{Q}_{A_{22}} \begin{pmatrix} \widetilde{R}_{A_{22}} \\ \mathbf{0} \end{pmatrix}, \quad (43)$$

where $\widehat{Q}_{A_{22}}$ is exactly unitary and $\widehat{Q}_{A_{22}} \approx \widetilde{Q}_{A_{22}}$, $\widetilde{R}_{A_{22}}$ is $\tilde{r}_2 \times n$ of full row rank,¹⁰ and $\Gamma = \begin{pmatrix} \Gamma_1 \\ \Gamma_2 \end{pmatrix}$ is the backward error of the QR factorization that can be estimated by

$$\left\| \begin{pmatrix} \Gamma_1 \\ \Gamma_2 \end{pmatrix}(:, i) \right\|_2 \leq \epsilon_{qr} \left\| \begin{pmatrix} \widehat{Q}_{E,2}^*(D + \Delta D) \\ \widehat{R}_E \widetilde{\Pi}_E^T \end{pmatrix}(:, i) \right\|_2. \quad (44)$$

¹⁰The zero block beneath of $\widetilde{R}_{A_{22}}$ may be void.

We can push Γ_1 and Γ_2 backward in D and E , respectively, as follows. First, $\widehat{Q}_{E,2}^*(D + \Delta D) + \Gamma_1 = \widehat{Q}_{E,2}^*(D + \Delta D + \widehat{Q}_{E,2}\Gamma_1)$ and

$$\widehat{Q}_E^*(D + \Delta D + \widehat{Q}_{E,2}\Gamma_1) = \begin{pmatrix} \widehat{Q}_{E,1}^*(D + \Delta D) \\ \widehat{Q}_{E,2}^*(D + \Delta D) + \Gamma_1 \end{pmatrix}.$$

If D and E are so scaled that their norms are nearly of the same order, then Γ_1, Γ_2 will be, respectively, their relatively small perturbations. Further, an analogous conclusion holds also column-wise, which motivates scaling the initial data by diagonal matrices to equilibrate on the matrix elements level.¹¹ Hence, if the additive perturbation ΔD is replaced with $\Delta_\Sigma D = \Delta D + \widehat{Q}_{E,2}\Gamma_1$, \widetilde{X}_1 remains unchanged, and \widetilde{X}_2 is precisely as in (43). (Here $\widetilde{X} = \begin{pmatrix} \widetilde{X}_1 \\ \widetilde{X}_2 \end{pmatrix}$.) Similarly, using $\delta E, \Delta E$ from Proposition 1, we have

$$\widehat{Q}_E^*(E + \delta E + \Delta E + \widehat{Q}_{E,1}\Gamma_2) = \begin{pmatrix} \widetilde{R}_E \widetilde{\Pi}_E^T + \Gamma_2 \\ \mathbf{0} \end{pmatrix}. \quad (45)$$

Now define $\widetilde{P}_2 = (\mathbb{I}_{2n+\widetilde{r}_E} \oplus \widehat{Q}_{A_{22}}^*) \widetilde{\Pi}$ analogously to (17). The matrix in (39) can be interpreted as an exact transformation of type (38), followed by the transformation of type (18) with \widetilde{P}_2 , but with initial matrices that are changed as $D \rightsquigarrow D + \Delta_\Sigma D$; $E \rightsquigarrow E + \Delta_\Sigma E$, $\Delta_\Sigma E = \delta E + \Delta E + \widehat{Q}_{E,1}\Gamma_2$. The transformation reads:

$$\begin{pmatrix} \mathbb{I}_{2n+\widetilde{r}_E} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \widehat{Q}_{A_{22}}^* & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbb{I}_{n-\widetilde{r}_E} \end{pmatrix} \begin{pmatrix} \widetilde{\Pi} & \mathbf{0} \\ \mathbf{0} & \mathbb{I}_{n-\widetilde{r}_E} \end{pmatrix} \begin{pmatrix} \mathbb{I}_n & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \widehat{Q}_E^* & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbb{I}_n & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \widehat{Q}_E^* \end{pmatrix} \begin{pmatrix} B & \mathbf{0}_n & -\mathbb{I}_n & \mathbf{0} \\ D + \Delta_\Sigma D & \mathbf{0}_n & \mathbf{0}_n & -\mathbb{I}_n \\ \mathbf{0} & -\mathbb{I}_n & \mathbf{0}_n & \mathbf{0}_n \\ E + \Delta_\Sigma E & \mathbf{0}_n & \mathbf{0}_n & \mathbf{0}_n \end{pmatrix} \begin{pmatrix} \mathbb{I}_n & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbb{I}_n & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbb{I}_n & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \widehat{Q}_E \end{pmatrix} \quad (46)$$

Consider now the second coefficient. The block swapping on the right matrix (pre-multiplication of the rows by the $(3n + \widetilde{r}_E) \times (3n + \widetilde{r}_E)$ permutation matrix $\widetilde{\Pi}$) reads

$$\left(\begin{array}{cc|c} -A & \mathbf{0} & \\ \hline -\widehat{Q}_{E,1}^*(C + \Delta C) & -\widehat{Q}_{E,1}^* & \mathbf{0}_{2n \times (n+\widetilde{r}_E)} \\ -\widehat{Q}_{E,2}^*(C + \Delta C) & -\widehat{Q}_{E,2}^* & \\ \hline \mathbf{0}_{(n+\widetilde{r}_E) \times 2n} & & -\mathbb{I}_{n+\widetilde{r}_E} \end{array} \right) \rightarrow \left(\begin{array}{cc|c|c} -A & \mathbf{0} & \mathbf{0}_{n+\widetilde{r}_E} & \\ \hline -\widehat{Q}_{E,1}^*(C + \Delta C) & -\widehat{Q}_{E,1}^* & \mathbf{0}_{\widetilde{r}_E \times (n+\widetilde{r}_E)} & \\ \mathbf{0}_n & \mathbf{0}_n & -\mathbb{I}_n & \mathbf{0}_{n \times \widetilde{r}_E} \\ \hline -\widehat{Q}_{E,2}^*(C + \Delta C) & -\widehat{Q}_{E,2}^* & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & -\mathbb{I}_{\widetilde{r}_E} \end{array} \right)$$

Recall, $\widetilde{Y} = \text{computed}(\widetilde{Q}_E^* C) = \widehat{Q}_E^*(C + \Delta C)$; introduce block-row partition $\widetilde{Y} = \begin{pmatrix} \widetilde{Y}_1 \\ \widetilde{Y}_2 \end{pmatrix}$ with $\widetilde{Y}_2 = \widehat{Q}_{E,2}^*(C + \Delta C)$. Similarly, introduce block-column partitions $\widetilde{Q}_{A_{22}}^* = (\widetilde{\Omega}_1 \ \widetilde{\Omega}_2)$, $\widehat{Q}_{A_{22}}^* = (\widehat{\Omega}_1 \ \widehat{\Omega}_2)$. The last column block $\widetilde{N}_{[4]}$ in the matrix

$$\widetilde{N} = \text{computed}(\widetilde{Q}_{A_{22}}^* \begin{pmatrix} -\widehat{Q}_{E,2}^*(C + \Delta C) & -\widehat{Q}_{E,2}^* & \mathbf{0}_{(n-\widetilde{r}_E) \times n} & \mathbf{0} \\ \mathbf{0}_{\widetilde{r}_E \times n} & \mathbf{0}_{\widetilde{r}_E \times n} & \mathbf{0}_{\widetilde{r}_E \times n} & -\mathbb{I}_{\widetilde{r}_E} \end{pmatrix}) = \begin{pmatrix} -\widetilde{N}_{[1]} & -\widetilde{N}_{[2]} & \widetilde{N}_{[3]} & \widetilde{N}_{[4]} \end{pmatrix} \quad (47)$$

is simply $-\widetilde{\Omega}_2$. Since we used $\widehat{Q}_{A_{22}}$ in the backward error analysis of the left-hand matrix, here we will have to use a mixed error analysis: $-\widetilde{\Omega}_2$ will be changed by a forward error into $-\widehat{\Omega}_2$. Recall that our model of the analysis (using exactly unitary instead of the computed numerically unitary matrices) will also require small forward perturbation to change $-\widehat{Q}_{E,1}^*$ into $-\widehat{Q}_{E,1}^*$.

¹¹See Example 6.

Consider now the first two blocks in \tilde{N} .

$$\tilde{N}_{[1]} = \text{computed}(\tilde{\Omega}_1 \tilde{Y}_2) = \tilde{\Omega}_1 \tilde{Y}_2 + \delta \tilde{N}_{[1]} = \hat{\Omega}_1 \tilde{Y}_2 + \delta \tilde{\Omega}_1 \tilde{Y}_2 + \delta \tilde{N}_{[1]}, \quad |\delta \tilde{N}_{[1]}| \leq \epsilon |\tilde{\Omega}_1| |\tilde{Y}_2|. \quad (48)$$

(Here ϵ estimates the backward error for matrix multiplication, $0 \leq \epsilon \leq O(n)\epsilon$.) In this block too, we will commit a forward error and replace it with

$$\hat{Q}_{A22}^* \begin{pmatrix} \hat{Q}_{E,2}^*(C + \Delta C) \\ \mathbf{0}_{\tilde{r}_E \times n} \end{pmatrix} = \hat{\Omega}_1 \tilde{Y}_2 = \tilde{N}_{[1]} - \Delta \tilde{N}_{[1]}, \quad \Delta \tilde{N}_{[1]} = \delta \tilde{\Omega}_1 \tilde{Y}_2 + \delta \tilde{N}_{[1]}. \quad (49)$$

To estimate this forward change we first note that, for each column index i ,

$$\|\tilde{Y}_2(:, i)\|_2 \leq \frac{\|\tilde{N}_{[1]}(:, i)\|_2}{1 - \|\delta \tilde{\Omega}_1\|_2 - \epsilon \|\tilde{\Omega}_1\|_2}.$$

Hence

$$\|\delta \tilde{\Omega}_1 \tilde{Y}_2(:, i)\|_2 \leq \frac{\|\delta \tilde{\Omega}_1\|_2 \|\tilde{N}_{[1]}(:, i)\|_2}{1 - \|\delta \tilde{\Omega}_1\|_2 - \epsilon \|\tilde{\Omega}_1\|_2}, \quad \|\delta \tilde{N}_{[1]}(:, i)\|_2 \leq \frac{\epsilon \|\tilde{\Omega}_1\|_2 \|\tilde{N}_{[1]}(:, i)\|_2}{1 - \|\delta \tilde{\Omega}_1\|_2 - \epsilon \|\tilde{\Omega}_1\|_2}, \quad (50)$$

and we conclude that $\Delta \tilde{N}_{[1]}$ is a column-wise small perturbation of $\tilde{N}_{[1]}$. Computation of $\tilde{N}_{[2]} = \text{computed}(\tilde{\Omega}_1 \tilde{Q}_{E,2}^*)$ is analogous, but for the purpose of mixed stability interpretation, $\tilde{Q}_{E,2}^*$ has to be replaced with $\hat{Q}_{E,2}^* = \tilde{Q}_{E,2}^* - \delta \tilde{Q}_{E,2}^*$, which yields

$$\tilde{N}_{[2]} = \hat{\Omega}_1 \hat{Q}_{E,2}^* + \hat{\Omega}_1 \delta \tilde{Q}_{E,2}^* + \underbrace{\delta \tilde{\Omega}_1 \hat{Q}_{E,2}^* + \delta \tilde{\Omega}_1 \delta \tilde{Q}_{E,2}^*}_{\delta \tilde{\Omega}_1 \tilde{Q}_{E,2}^*} + \delta \tilde{N}_{[2]}, \quad |\delta \tilde{N}_{[2]}| \leq \epsilon |\tilde{\Omega}_1| |\tilde{Q}_{E,2}^*|. \quad (51)$$

Hence, if the computed right-hand matrix is changed by a forward perturbation as

$$\begin{aligned} & \left(\begin{array}{c|c|c} -A & \mathbf{0} & \mathbf{0}_{n+\tilde{r}_E} \\ \hline -\hat{Q}_{E,1}^*(C + \Delta C) & -\hat{Q}_{E,1}^* & \mathbf{0}_{\tilde{r}_E \times (n+\tilde{r}_E)} \\ \mathbf{0}_n & \mathbf{0}_n & -\mathbb{I}_n \mid \mathbf{0}_{n \times \tilde{r}_E} \\ \hline -\tilde{N}_{[1]} & -\tilde{N}_{[2]} & \tilde{N}_{[3]} \mid \tilde{N}_{[4]} \end{array} \right) + \left(\begin{array}{c|c|c} \mathbf{0} & \mathbf{0} & \mathbf{0}_{n+\tilde{r}_E} \\ \hline \mathbf{0} & -\delta \tilde{Q}_{E,1}^* & \mathbf{0}_{\tilde{r}_E \times (n+\tilde{r}_E)} \\ \mathbf{0}_n & \mathbf{0}_n & \mathbf{0}_n \mid \mathbf{0}_{n \times \tilde{r}_E} \\ \hline \Delta \tilde{N}_{[1]} & \Delta \tilde{N}_{[2]} & \mathbf{0} \mid \Delta \tilde{N}_{[4]} \end{array} \right) \\ & = \left(\begin{array}{c|c|c} -A & \mathbf{0} & \mathbf{0}_{n+\tilde{r}_E} \\ \hline -\hat{Q}_{E,1}^*(C + \Delta C) & -\hat{Q}_{E,1}^* & \mathbf{0}_{\tilde{r}_E \times (n+\tilde{r}_E)} \\ \mathbf{0}_n & \mathbf{0}_n & -\mathbb{I}_n \mid \mathbf{0}_{n \times \tilde{r}_E} \\ \hline \hat{Q}_{A22}^* \begin{pmatrix} -\hat{Q}_{E,2}^*(C + \Delta C) \\ \mathbf{0}_{\tilde{r}_E \times n} \end{pmatrix} & \hat{Q}_{A22}^* \begin{pmatrix} -\hat{Q}_{E,2}^* \\ \mathbf{0}_{\tilde{r}_E \times n} \end{pmatrix} & \hat{Q}_{A22}^* \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \end{pmatrix} \mid \hat{Q}_{A22}^* \begin{pmatrix} \mathbf{0} \\ -\mathbb{I}_{\tilde{r}_E} \end{pmatrix} \end{array} \right), \quad (52) \end{aligned}$$

the resulting matrix is the $(3n + \tilde{r}_E) \times (3n + \tilde{r}_E)$ main submatrix of

$$\begin{pmatrix} \mathbb{I}_{2n+\tilde{r}_E} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \hat{Q}_{A22}^* & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbb{I}_{n-\tilde{r}_E} \end{pmatrix} \begin{pmatrix} \tilde{\Pi} & \mathbf{0} \\ \mathbf{0} & \mathbb{I}_{n-\tilde{r}_E} \end{pmatrix} \begin{pmatrix} \mathbb{I}_n & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \hat{Q}_E^* & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbb{I}_n & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \hat{Q}_E^* \end{pmatrix} \begin{pmatrix} -A & \mathbf{0}_n & \mathbf{0}_n & \mathbf{0}_n \\ -(C+\Delta C) & -\mathbb{I}_n & \mathbf{0}_n & \mathbf{0}_n \\ \mathbf{0}_n & \mathbf{0}_n & -\mathbb{I}_n & \mathbf{0}_n \\ \mathbf{0}_n & \mathbf{0}_n & \mathbf{0}_n & -\mathbb{I}_n \end{pmatrix} \begin{pmatrix} \mathbb{I}_n & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbb{I}_n & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbb{I}_n & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \hat{Q}_E \end{pmatrix}. \quad (53)$$

□

6 Numerical examples

In this section, we present numerical examples and compare our new algorithm `kvarteig` with `polyeig` from MATLAB, `quadeig` [20] and `kvadeig` (including `balanced_kvadeig`¹²) [14] applied to the linearization (6) and the quadratification (5), respectively.

Our goal is to illustrate the potential of the techniques introduced in `kvadeig` and `kvarteig`, and to motivate further development. We in particular stress the benefits of additional balancing of the coefficient matrices, that is used in combination with the well developed parameter scaling. The analysis of backward errors in §5 indicates, and numerical experiments in this section provide empirical evidence of importance of well balanced data. Balancing is applicable, *mutatis mutandis*, to other solvers as well.

All test examples are taken from the NLEVP benchmark collection [2]. In all examples and figures shown in this section, the eigenpairs are indexed so that the eigenvalues are sorted increasingly in modulus.

Example 1. We first test `kvarteig` on three examples with the default input values: `butterfly` ($n = 64$); `orr_sommerfeld` ($n = 64$); `planar waveguide` ($n = 129$). The results are tested using the norm-wise (residual) backward error (4).

In the first run of the experiment, `polyeig`, `quadeig`, and both variants of `kvadeig` worked with raw data A, B, C, D, E , and the parameter scaling is applied in `quadeig`, `kvadeig` only to the quadratic pencil $\lambda^2\mathbb{M} + \lambda\mathbb{C} + \mathbb{K}$ from (5). In `balanced_kvadeig`, parameter scaling is combined with diagonal balancing of $\mathbb{M}, \mathbb{C}, \mathbb{K}$. For the sake of the experiment, `kvarteig` worked in two modes: (i) with parameter scaling (designated as `kvarteig`); and (ii) without parameter scaling, but with the diagonal balancing switched on (designated as `balanced_kvarteig (-s)`).

Switching the scaling off may simulate the case of an unsuccessful parameter scaling of the initial matrix coefficients. Also, this may serve as a simulation of a genuine quadratic problem in which the coefficients are composed of blocks with different parameter dependencies, possibly on different scales – then parameter scaling cannot resolve different scales inside $\mathbb{M}, \mathbb{C}, \mathbb{K}$. Hence, this is primarily a test of the quadratic solvers as potential tools for quadratification based solution of quartic problems. We refer the reader to §1.1.4, Remark 1, §2.1.1, and [14]. Further, with `balanced_kvarteig (-s)` we want to check whether diagonal balancing can make a difference in the absence (or failure) of the parameter scaling.

The extreme values of η over all computed right eigenpairs are given in Table 1:

Table 1: Comparison of backward errors for `polyeig`, `quadeig`, `kvadeig` and `kvarteig`

Algorithm	butterfly		orr_sommerfeld		planar waveguide	
	min η	max η	min η	max η	min η	max η
<code>polyeig</code>	2.04e-016	8.61e-016	1.36e-017	8.01e-006	1.60e-016	3.08e-012
<code>quadeig</code>	6.56e-017	2.03e-015	6.11e-015	4.07e-004	4.99e-016	2.03e-009
<code>kvadeig</code>	6.56e-017	2.03e-015	6.25e-021	2.12e-007	4.75e-016	1.67e-009
<code>balanced_kvadeig</code>	6.56e-017	2.03e-015	2.81e-021	2.06e-012	1.49e-016	2.32e-012
<code>balanced_kvarteig (-s)</code>	3.18e-017	9.11e-016	3.40e-021	5.25e-008	3.20e-016	5.16e-012
<code>kvarteig</code>	5.84e-017	1.13e-015	6.37e-021	1.76e-015	4.32e-016	1.75e-013

Note that `quadeig` and `kvadeig` had relatively large relative errors for `orr_sommerfeld` and `planar waveguide`, while `balanced_kvadeig` performed well despite the fact that it received unscaled original matrices. Although contrived, this example illustrates the main point well –

¹²`balanced_kvadeig` denotes the `kvadeig` algorithm enhanced with balancing by diagonal scaling matrices, see Remark 1.

parameter scaling (here applied to $\lambda^2 M + \lambda C + \mathbb{K}$) combined with diagonal balancing is better than parameter scaling alone.

We complete this experiment with the initial parameter scaling included in all methods. It performed well and, as a result, all measured backward errors were in all five methods at most $O(10^{-12})$. This is of the order of the machine precision multiplied by a low order polynomial of the dimension n of the problem.

Example 2. Structured backward errors provide a better insight into the numerical quality of the computed solutions. For the data of `orr_sommerfeld` in Example 1, we compute for each right eigenpair λ, x the component-wise backward error

$$\begin{aligned} \omega(\lambda, x) &= \min\{\epsilon : (\lambda^4 \tilde{A} + \lambda^3 \tilde{B} + \lambda^2 \tilde{C} + \lambda \tilde{D} + \tilde{E})x = \mathbf{0}, |\delta A| \leq \epsilon |A|, \dots, |\delta E| \leq \epsilon |E|\} \\ &= \max_{i=1:n} \frac{|(\lambda^4 A + \lambda^3 B + \lambda^2 C + \lambda D + E)x|_i}{(|\lambda|^4 |A| + |\lambda|^3 |B| + |\lambda|^2 |C| + |\lambda| |D| + |E|)|x|_i}, \quad \tilde{A} = A + \delta A, \dots, \tilde{E} = E + \delta E. \end{aligned} \quad (54)$$

The corresponding error $\omega'(\lambda, y)$ for a left eigenpair λ, y is defined analogously. We examine the component-wise backward errors in the `orr_sommerfeld` example with two sets of defining parameters. Recall, the function from the NLEVP library for generating this quartic eigenvalue problem has three optional input arguments, n , ω and R : n represents the dimension of the problem, ω is the frequency, and R is the Reynolds number. The default values are: $n = 64$, $\omega = 0.26943$ and $R = 5772$ (these values are used in Table 1). In the first test, we use these default values.

The values of $\omega(\lambda, x)$ and $\omega'(\lambda, y)$ are shown for all computed eigenpairs in Figure 2.

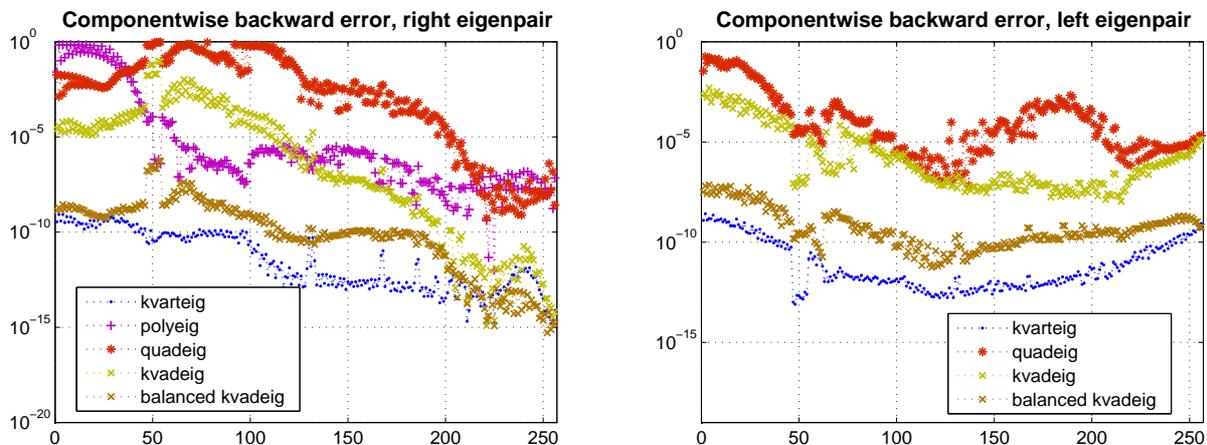


Figure 2: (Example 2.) Component-wise backward errors (54) for all computed eigenpairs of the `orr_sommerfeld` example with $n = 64$, $\omega = 0.26943$ and $R = 5772$.

In the second run of the test, we increase the Reynolds number to $R = 10000$. The norm-wise and the component-wise backward errors for all eigenvalues, are shown in Figure 3 and Figure 4, respectively. Note how the backward error for `kvar eig` in Figure 3 remains nearly flat at the roundoff level, and how `kvadeig` also performs well (even with structured backward error for the right eigenpairs), despite being oblivious to the underlying structure of the quadratification and receiving unscaled original coefficients.

A conclusion of this and Example 1 is that quadratic solver equipped with parameter scaling and diagonal balancing might work reasonably well on a quadratification of the quartic problem, even when the scaling of the coefficients of the original quartic problem is omitted or unsuccessful.

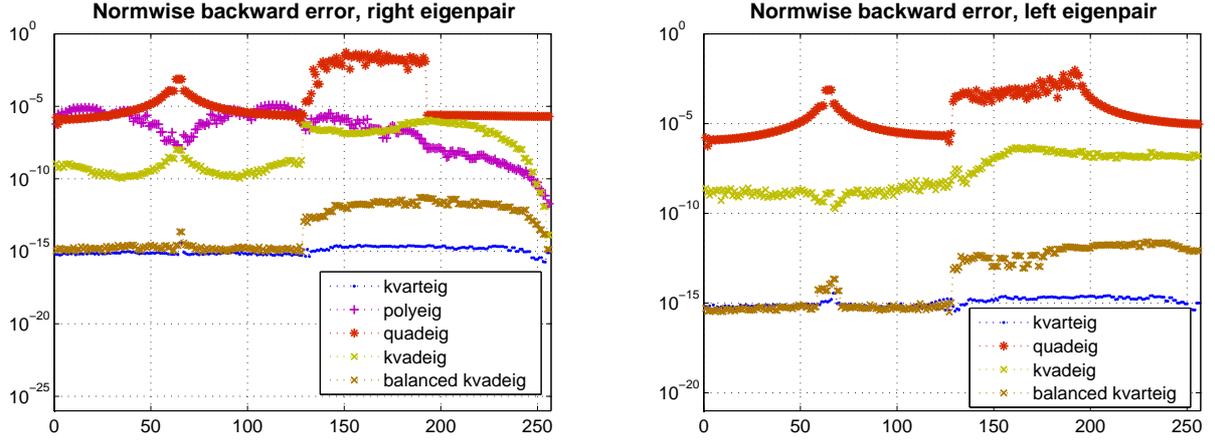


Figure 3: (Example 2.) The norm-wise backward errors for all computed eigenpairs of the `orr_sommerfeld` example with $n = 64$, $\omega = 0.26943$ and $R = 10000$.

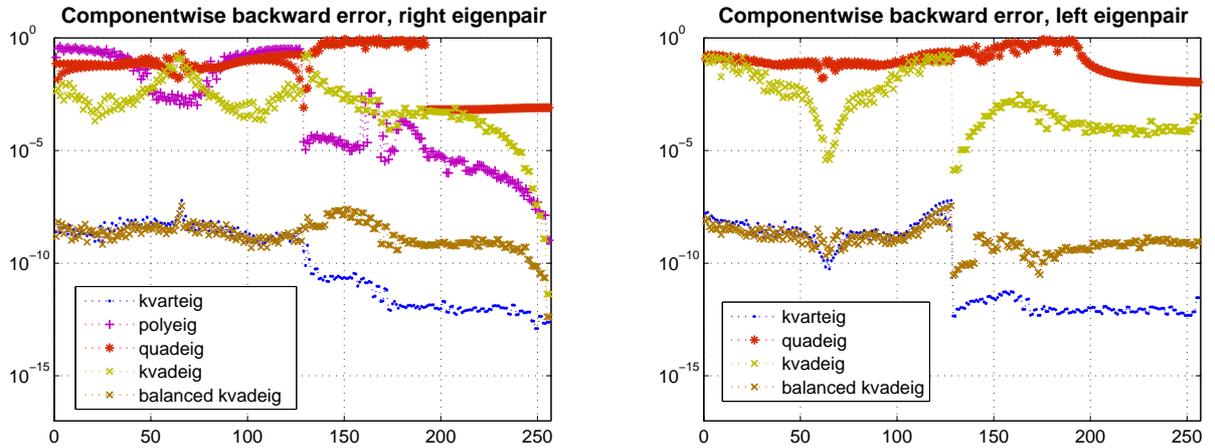
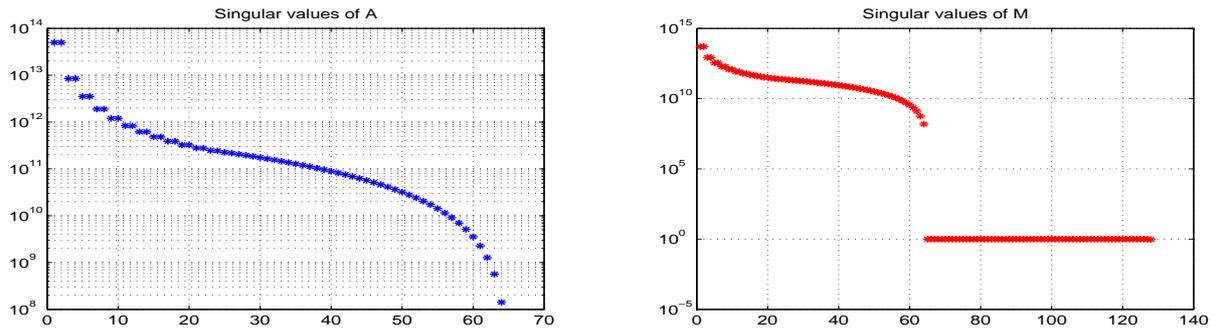


Figure 4: (Example 2.) The component-wise backward errors of all computed eigenpairs of the `orr_sommerfeld` example with $n = 64$, $\omega = 0.26943$ and $R = 10000$.

Remark 8. *The results of this experiment, with the computed backward errors shown in Figure 3 and Figure 4, are instructive. First, in this example `quadeig` deflated 64 infinite eigenvalues of the quadratic pencil $\lambda^2\mathbb{M} + \lambda\mathbb{C} + \mathbb{K}$ (see (5)), because in the preprocessing stage of the algorithm, the numerical rank of the matrix $\mathbb{M} = \begin{pmatrix} A & \mathbf{0} \\ \mathbf{C} & \mathbb{I}_n \end{pmatrix}$ of order 128 is computed as 64. On the other hand, the existence of infinite eigenvalues in the original quartic eigenvalue problem depends on the rank of the leading coefficient matrix A . If we inspect the singular values¹³ $\sigma_i(\mathbb{M})$ of \mathbb{M} and $\sigma_i(A)$ of A , then, as clearly shown in Figure 5, A is numerically of full rank (its condition number is below 10^6 , so `kvarteig` safely removed the possibility of infinite eigenvalues). On the other hand, $\sigma_{65}(\mathbb{M})/\sigma_1(\mathbb{M})$ is at the level of dimension of \mathbb{M} times machine precision and in many algorithms this is the truncation threshold for numerical rank deficiency. Note that parameter scaling of the quadratic pencil (5) cannot remove this problem.*

On the other hand, `kvadeig` (applied to the same $\lambda^2\mathbb{M} + \lambda\mathbb{C} + \mathbb{K}$) declared the matrix \mathbb{M} non-singular, and thus no infinite eigenvalues were deflated nor found by the QZ algorithm. This is because `kvadeig` uses more local truncation strategy; it truncates at index i if $\sigma_{i+1}(\mathbb{M})/\sigma_i(\mathbb{M})$ is estimated to be small; see Remark 2 and Example 4. Good results by `balanced_kvadeig` are

¹³Singular values are indexed in non-increasing order, $\sigma_i(\cdot) \geq \sigma_{i+1}(\cdot)$.



(a) Leading coefficient A of the quartic problem (b) Leading coefficient M of the quadratification

Figure 5: (Example 2.) Singular values of the leading coefficient matrices of the original quartic problem (1) and the quadratification (5). Note that $\sigma_{\max}(A)/\sigma_{\min}(A) = O(10^6)$, $\sigma_{65}(M)/\sigma_1(M) = O(n)\varepsilon = O(10^{-14})$, $\sigma_{65}(M)/\sigma_{64}(M) = O(\sqrt{\varepsilon})$. Here $\varepsilon \approx 2.2 \cdot 10^{-16}$ is the machine precision.

due to the additional balancing [14, §4.2], and this example once more justifies our approach in *kvadeig* (using local truncation strategy and balancing in combination with parameter scaling).

Now, we turn on the parameter scaling, which is a necessary tool for numerical stability of a polynomial eigensolver. Although the scaling described in §1.1.3 is a simple combination of the existing and well known formulas, it seems that it works well. In particular, in many cases it works well for *polyeig*, as we already showed in Example 1. This is illustrated in the next two numerical experiments with the *orr_sommerfeld* example of dimensions $n = 64$ and $n = 1000$.

Example 3. We use the same benchmark problem as in the second part of Example 2 (*orr_sommerfeld* example with $n = 64$, $\omega = 0.26943$ and $R = 10000$.), but initially we scale the matrices as described in §1.1.3, so that all algorithms start with scaled data. This example has no infinite eigenvalues. The results of all algorithms depend on the *QZ* algorithm, thus the similar results. (It seems that *polyeig* and *kvarteig* are a little bit better than *kvadeig* and *quadeig*, which makes sense because both work on the original coefficients, while the quadratic solvers work on M , \mathbb{C} , \mathbb{K} from the quadratification.)

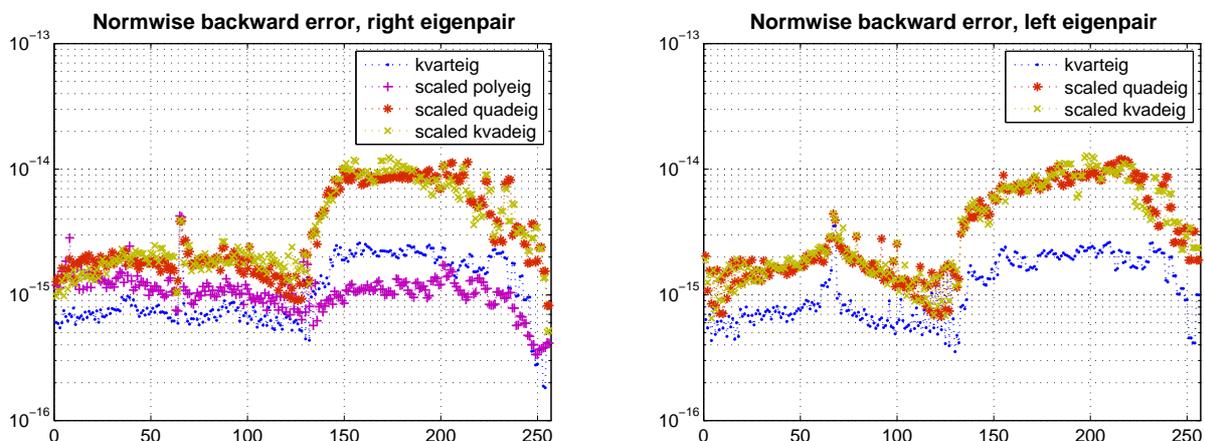


Figure 6: (Example 3.) The norm-wise backward errors for all computed eigenpairs for the *orr_sommerfeld* example with $n = 64$, $\omega = 0.26943$ and $R = 10000$.

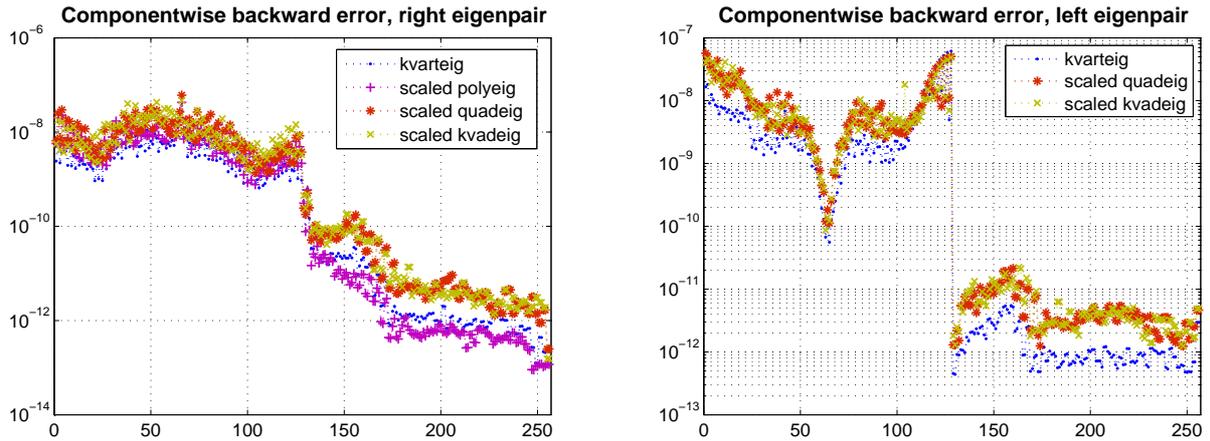


Figure 7: (Example 3.) The component-wise backward errors for all computed eigenpairs for the `orr_sommerfeld` example with $n = 64$, $\omega = 0.26943$ and $R = 10000$.

Example 4. We continue experimenting with the `orr_sommerfeld` example; we choose the default values of the Reynolds number R and the frequency ω , but increase the dimension to $n = 1000$, and compute all 4000 eigenpairs. The matrix coefficients are scaled using the same strategy as in the previous example.¹⁴

An application of `quadeig` to the quadratification (5) returned 282 infinite eigenvalues. With `kvadeig` and the same quadratification, 31 infinite eigenvalues are detected. On the other hand, if we use balancing (`balanced_kvadeig`), the leading coefficient matrix is declared regular, and no infinite eigenvalues are detected. The difference is mainly due to the softer drop-off truncation in the rank revealing QR factorization. The result of `kvarteig` also depends on the truncation strategy. If the truncation of the pivoted QR factorization is done relative to the norm of A , the numerical rank is 988, meaning that 12 infinite eigenvalues are deflated immediately in the preprocessing phase. In the case of drop-off strategy, the matrix A is not numerically rank deficient.

The norm-wise and component-wise backward errors for the computed right and left eigenpairs are shown in Figures 8, 9, 10, 11.

Example 5. In this example, we use transposed matrices from the `mirror` example¹⁵ and scale them as described in §1.1.3. The number of zero and infinite eigenvalues found by the four algorithms were: `polyeig` (no zeros and 5 infinities); `quadeig` (7 zeros and 9 infinities); `kvadeig` and `kvarteig` 9 zeros and 9 infinite eigenvalues. The component-wise and the norm-wise backward errors are given in Figure 12.

Example 6. In this example we illustrate potential benefits of equilibration of the coefficient matrices on the element level, mentioned in Remark 1. Such diagonal scalings balance the absolute values of nonzero entries over all matrices.

We take the `butterfly` example and pre-multiply its coefficient matrices with diagonal matrix Δ with randomly permuted powers 2^i , $i = 1, \dots, n = 64$ on the diagonal. This is an entirely artificial step to simulate a situation with ill-conditioning caused by removable scaling

¹⁴Without parameter scaling of the initial data, the Matlab function `polyeig` failed completely – all computed eigenvalues were of the form $\pm \text{Inf} \pm \text{Inf}i$.

¹⁵Recall, we analyzed this example in §3.2.5, where we argued that there are nine zero and nine infinite eigenvalues.

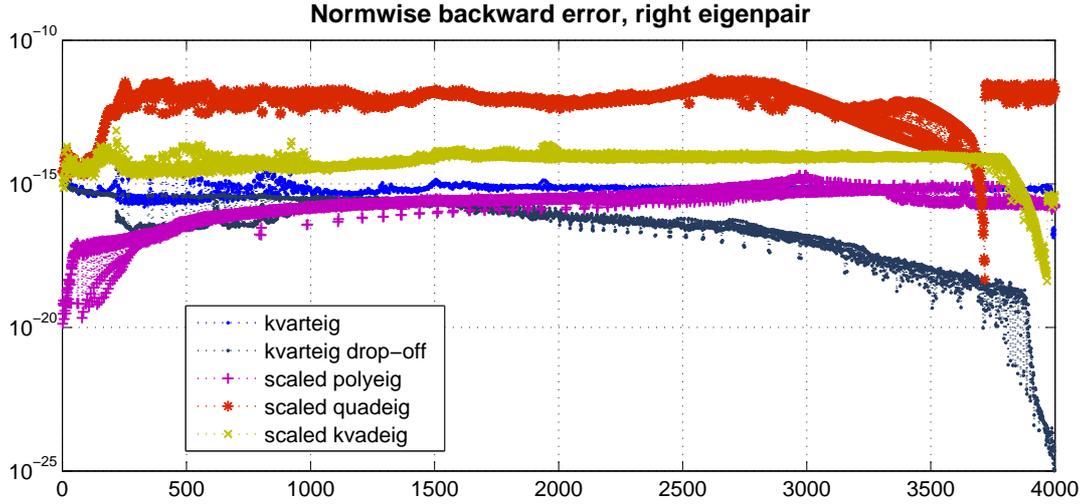


Figure 8: (Example 4.) Norm-wise backward errors for the right eigenpairs.

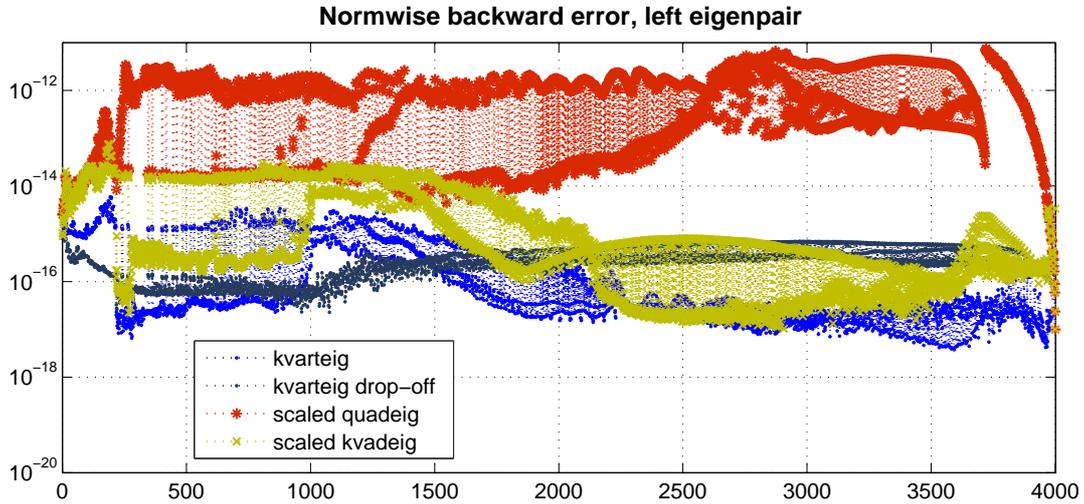


Figure 9: (Example 4) Norm-wise backward errors for the left eigenpairs.

(that may originate in an inappropriate scale of physical units). We obtain an equivalent problem, but numerical algorithms may be more or less sensitive to this change of representation.

Then, we compute balancing matrices Δ_ℓ , Δ_r (see Remark 1) and examine how this pre-processing $((A, B, C, D, E) \rightsquigarrow \Delta_\ell(A, B, C, D, E)\Delta_r)$ influences the numerical accuracy of the algorithms under study. The computed component-wise backward errors, shown in Figure 13, clearly demonstrate the impact of the balancing $(A, B, C, D, E) \rightsquigarrow \Delta_\ell(A, B, C, D, E)\Delta_r$. Note also that without balancing *kvadeig* still performs well, much better than *polyeig* and *quadeig* under the same conditions.

Example 7. In our last example, we checked the least squares approach to recovering the eigenvectors, as described in §4.1.2. The computed backward errors in all tested cases were comparable with the method of selecting the vector with smallest residual. We believe that this least squares approach could be useful for getting good eigenvectors for selected eigenvalues that are of particular importance in some applications.

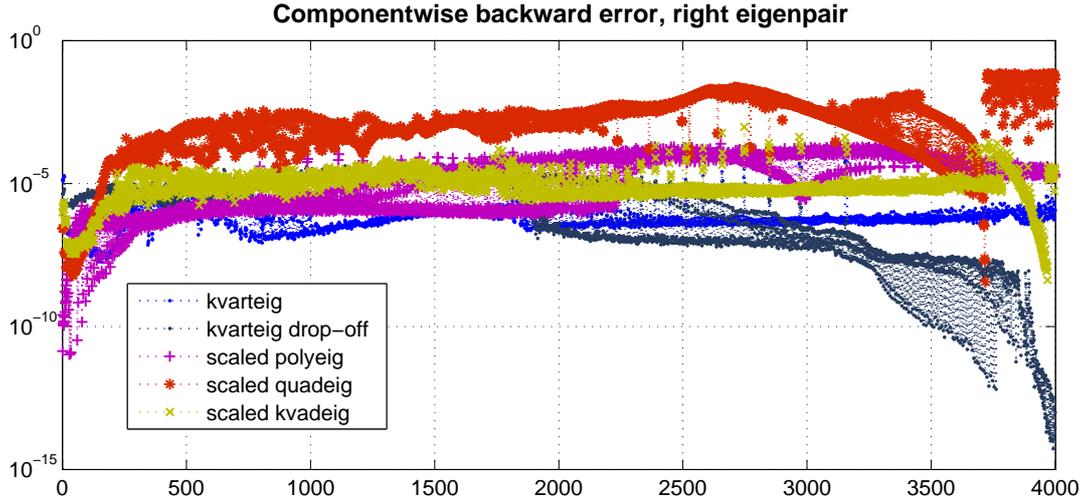


Figure 10: (Example 4.) Component-wise backward errors for the right eigenpairs.

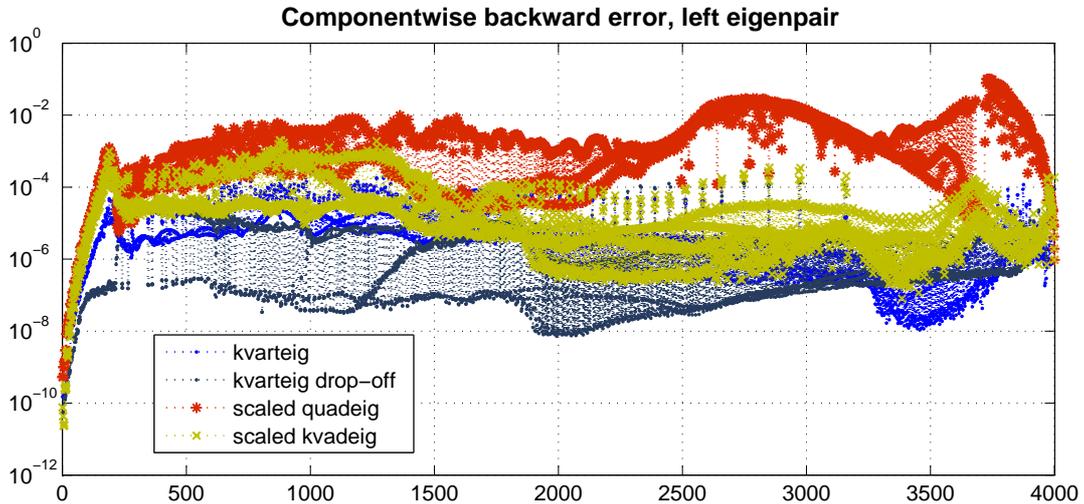


Figure 11: (Example 4.) Component-wise backward errors for the left eigenpairs.

7 Concluding remarks

We have shown that the proposed algorithm `kvarteig` for solving quartic eigenvalue problems is a useful contribution that fills the gap in the toolbox for the polynomial eigenvalue problems, both for the full solution of medium size non-structured problems and for solving the projected problems in subspace based methods for large scale structured/sparse problems. Numerical experiments with the benchmark examples from the NLEVP collection show that `kvarteig` is superior to `polyeig` from Matlab, or `quadeig` applied to a quadratification of the original quartic problem. Further, the numerical performances of `kvadeig` on the quadratification of the quartic problem additionally justify the modifications that underpinned the development both in [14] and in this paper.

Given the wide spectrum of applications of the quartic eigenvalue problem, we are certain that our proposed algorithm will prove useful in many computational tasks in applied sciences and engineering. Further, the presented techniques can be adapted for other methods and suitable linearizations of polynomial eigenvalue problems.

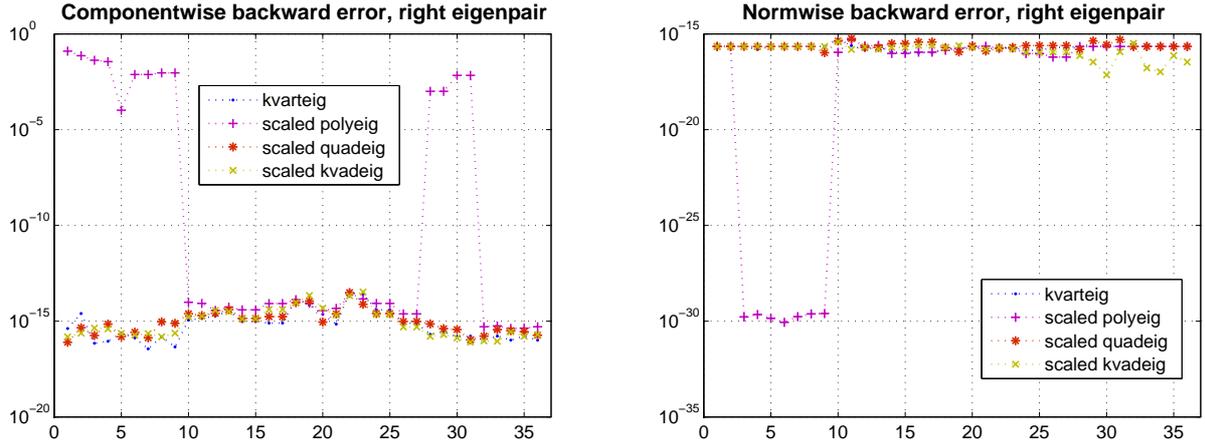


Figure 12: (Example 5) *Left panel*: Component-wise backward errors for the transposed `mirror` example. *Right panel*: Norm-wise backward error for the original `mirror` example (This is the right panel from Figure 1, here given for comparison.)

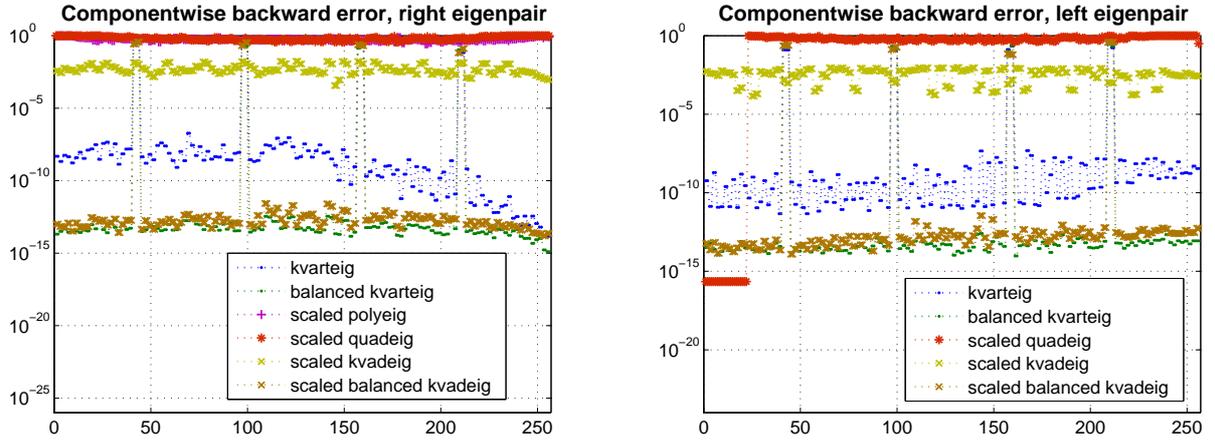


Figure 13: (Example 6, `butterfly`.) Component-wise backward errors for the modified `butterfly` example, where the coefficients are premultiplied by a diagonal matrix Δ , $(A, B, C, D, E) \rightsquigarrow \Delta(A, B, C, D, E)$.

An early version of this work is available at [16].

Acknowledgement

This research has been supported by the Croatian Science Foundation (CSF) grants IP-2019-04-6268, and in part (the second author) UIP-2019-04-5200. Parts of this work originate from the second author's thesis [36]. The authors thank to Serkan Gugercin (Virginia Tech, Blacksburg), Luka Grubišić and Zvonimir Bujanović (University of Zagreb) for valuable comments, and in particular to the three anonymous referees for their constructive criticism and detailed reports.

References

- [1] Timo Betcke. Optimal scaling of generalized and polynomial eigenvalue problems. *SIAM Journal on Matrix Analysis and Applications*, 30(4):1320–1338, 2009.
- [2] Timo Betcke, Nicholas J. Higham, Volker Mehrmann, Christian Schröder, and Françoise Tisseur. NLEVP: A collection of nonlinear eigenvalue problems. *ACM Transactions on Mathematical Software (TOMS)*, 39(2):1–28, 2013.
- [3] Nela Bosner, Zvonimir Bujanović, and Zlatko Drmač. Efficient generalized Hessenberg form and applications. *ACM Trans. Math. Softw.*, 39(3), May 2013.
- [4] Nela Bosner, Zvonimir Bujanović, and Zlatko Drmač. Parallel solver for shifted systems in a hybrid CPU–GPU framework. *SIAM Journal on Scientific Computing*, 40(4):C605–C633, 2018.
- [5] Peter Businger and Gene H. Golub. Linear least squares solutions by Householder transformations. *Numerische Mathematik*, 7(3):269–276, 1965.
- [6] Carmen Campos and Jose E. Roman. Parallel Krylov solvers for the polynomial eigenvalue problem in SLEPc. *SIAM Journal on Scientific Computing*, 38(5):S385–S411, 2016.
- [7] Hongjia Chen, Akira Imakura, and Tetsuya Sakurai. Improving backward stability of Sakurai–Sugiura method with balancing technique in polynomial eigenvalue problem. *Applications of Mathematics*, 62(4):357–375, 2017.
- [8] Gökhan Danabasoglu and Sedat Biringen. A Chebyshev matrix method for the spatial modes of the Orr—Sommerfeld equation. *International Journal for Numerical Methods in Fluids*, 11:1033 – 1037, 11 1990.
- [9] Fernando De Terán, Froilán M Dopico, and D Steven Mackey. Spectral equivalence of matrix polynomials and the index sum theorem. *Linear Algebra and its Applications*, 459:264–333, 2014.
- [10] Fernando De Terán, Froilán M. Dopico, and Paul Van Dooren. Constructing strong ℓ -ifications from dual minimal bases. *Linear Algebra and its Applications*, 495:344 – 372, 2016.
- [11] Andrii Dmytryshyn and Froilán M. Dopico. Generic complete eigenstructures for sets of matrix polynomials with bounded rank and degree. *Linear Algebra and its Applications*, 535:213–230, 2017.
- [12] Andrii Dmytryshyn, Stefan Johansson, Bo Kågström, and Paul van Dooren. Geometry of matrix polynomial spaces. *Found Comput Math*, 20(3):423–450, 2020.
- [13] Froilán M Dopico, Piers W Lawrence, Javier Pérez, and Paul Van Dooren. Block Kronecker linearizations of matrix polynomials and their backward errors. *Numerische Mathematik*, 140(2):373–426, 2018.
- [14] Zlatko Drmač and Ivana Šain Glibić. New numerical algorithm for deflation of infinite and zero eigenvalues and full solution of quadratic eigenvalue problems. *ACM Transactions on Mathematical Software (TOMS)*, 46(4):1–32, 2020.

- [15] Zlatko Drmač and Zvonimir Bujanović. On the failure of rank revealing QR factorization software – a case study. *ACM Transactions on Mathematical Software (TOMS)*, 35(2):1–28, 2008.
- [16] Zlatko Drmač and Ivana Šain Glibić. An algorithm for the complete solution of the quartic eigenvalue problem. *arXiv e-prints*, page arXiv:1905.07013, May 2019.
- [17] Hung-Yuan Fan, Wen-Wei Lin, and Paul Van Dooren. Normwise scaling of second order polynomial matrices. *SIAM Journal on Matrix Analysis and Applications*, 26(1):252–256, 2004.
- [18] Brendan Gavin, Agnieszka Miedlar, and Eric Polizzi. FEAST eigensolver for nonlinear eigenvalue problems. *Journal of Computational Science*, 27:107–117, jul 2018.
- [19] Gene Golub, Virginia Klema, and Gilbert W Stewart. Rank degeneracy and least squares problems. Technical report, Computer Science Department, Stanford University, 1976.
- [20] Sven Hammarling, Christopher J. Munro, and Françoise Tisseur. An algorithm for the complete solution of quadratic eigenvalue problems. *ACM Transactions on Mathematical Software (TOMS)*, 39(3):1–19, 2013.
- [21] Nicholas J. Higham, Ren-Cang Li, and Françoise Tisseur. Backward error of polynomial eigenproblems solved by linearization. *SIAM Journal on Matrix Analysis and Applications*, 29(4):1218–1241, 2008.
- [22] Stefan Johansson, Bo Kågström, and Paul Van Dooren. Stratification of full rank polynomial matrices. *Linear Algebra and its Applications*, 439(4):1062–1090, 2013. 17th Conference of the International Linear Algebra Society, Braunschweig, Germany, August 2011.
- [23] Nicola Mastronardi and Paul Van Dooren. Revisiting the stability of computing the roots of a quadratic polynomial. *CoRR*, abs/1409.8072, 2014.
- [24] Volker Mehrmann and Tatjana Stykel. Descriptor systems: A general mathematical framework for modelling, simulation and control (Deskriptorsysteme: Ein allgemeines mathematisches Konzept für Modellierung, Simulation und Regelung). *Automatisierungstechnik*, 54(8):405 – 415, 2006.
- [25] Volker Mehrmann and David S Watkins. Polynomial eigenvalue problems with Hamiltonian structure. *Electronic Transactions on Numerical Analysis*, 13:106–118, 2002.
- [26] Tetsuya Sakurai and Hiroshi Sugiura. A projection method for generalized eigenvalue problems using numerical integration. *Journal of Computational and Applied Mathematics*, 159(1):119 – 128, 2003. 6th Japan-China Joint Seminar on Numerical Mathematics; In Search for the Frontier of Computational and Applied Mathematics toward the 21st Century.
- [27] Jung Heon Song, Matthias Maier, and Mitchell Luskin. Nonlinear eigenvalue problems for coupled Helmholtz equations modeling gradient-index graphene waveguides. *Journal of Computational Physics*, 423:109871, 2020.
- [28] David Stowell and Johannes Tausch. Variational formulation for guided and leaky modes in multilayer dielectric waveguides. *Communications in Computational Physics*, 7(3):564, 01 2010.

- [29] Françoise Tisseur. Backward error and condition of polynomial eigenvalue problems. *Linear Algebra and its Applications*, 309(1-3):339–361, 2000.
- [30] Françoise Tisseur and Marc Van Barel. Min-max elementwise backward error for roots of polynomials and a corresponding backward stable root finder. *arXiv preprint arXiv:2001.05281*, 2020.
- [31] Marc Van Barel and Françoise Tisseur. Polynomial eigenvalue solver based on tropically scaled Lagrange linearization. *Linear Algebra and its Applications*, 542:186–208, 2018. Proceedings of the 20th ILAS Conference, Leuven, Belgium 2016.
- [32] Paul Van Dooren. The computation of Kronecker’s canonical form of a singular pencil. *Linear Algebra and its Applications*, 27:103–140, 1979.
- [33] Paul Van Dooren and Patrick Dewilde. The eigenstructure of an arbitrary polynomial matrix: computational aspects. *Linear Algebra and its Applications*, 50:545–579, 1983.
- [34] Istvan A. Veres, Thomas Berer, and Osamu Matsuda. Complex band structures of two dimensional phononic crystals: Analysis by the finite element method. *Journal of Applied Physics*, 114(8):083519, 2013.
- [35] R.F. Vieira, F.B. Virtuoso, and E.B.R. Pereira. A higher order model for thin-walled structures with deformable cross-sections. *International Journal of Solids and Structures*, 51(3-4):575 – 598, 2014.
- [36] I. Šain Glibić. *Robust numerical methods for nonlinear eigenvalue problems*. PhD thesis, University of Zagreb, Faculty of Science, Department of Mathematics, 12 2018.
- [37] David S. Watkins. Performance of the QZ algorithm in the presence of infinite eigenvalues. *SIAM Journal on Matrix Analysis and Applications*, 22(2):364–375, 2000.
- [38] Yong Xiao, Jihong Wen, Lingzhi Huang, and Xisen Wen. Analysis and experimental realization of locally resonant phononic plates carrying a periodic array of beam-like resonators. *Journal of Physics D: Applied Physics*, 47:045307, 01 2013.
- [39] Shinnosuke Yokota and Tetsuya Sakurai. A projection method for nonlinear eigenvalue problems using contour integrals. *JSIAM Letters*, 5:41–44, 2013.
- [40] Linghui Zeng and Yangfeng Su. A backward stable algorithm for quadratic eigenvalue problems. *SIAM Journal on Matrix Analysis and Applications*, 35(2):499–516, 2014.
- [41] H. Zha. The restricted singular value decomposition of matrix triplets. *SIAM Journal on Matrix Analysis and Applications*, 12(1):172–194, 1991.
- [42] B. Zhang and Y. F. Li. A method for calibrating the central catadioptric camera via homographic matrix. In *2008 International Conference on Information and Automation*, pages 972–977, June 2008.