# The Effect of Aggressive Early Deflation on the Convergence of the QR Algorithm

Daniel Kressner[*]

June 23, 2008

### Abstract

Aggressive early deflation has proven to significantly enhance the convergence of the QR algorithm for computing the eigenvalues of a nonsymmetric matrix. One purpose of this paper is to point out that this deflation strategy is equivalent to extracting converged Ritz vectors from certain Krylov subspaces. As a special case, the single-shift QR algorithm enhanced with aggressive early deflation corresponds to a Krylov subspace method whose starting vector undergoes a Rayleigh-quotient iteration. It is shown how these observations can be used to derive improved convergence bounds for the QR algorithm.

## 1 Introduction

Let $A$ be a complex $n \times n$ matrix. The aim of the QR algorithm is to compute a Schur decomposition $S = Q^{\mathsf{H}} A Q$, where $Q \in \mathbb{C}^{n \times n}$ is unitary and $S \in \mathbb{C}^{n \times n}$ is upper triangular. The QR algorithm, as introduced by Francis [13] and Kublanovskaya [19], is an iterative process that generates a sequence of unitarily similar matrices $A_0 \leftarrow A, A_1, A_2, \dots$. Before each iteration, $m$ so-called shifts $\sigma_1, \dots, \sigma_m \in \mathbb{C}$ are skillfully chosen, defining the shift polynomial $p_i(\lambda) = (\lambda - \sigma_1) \cdots (\lambda - \sigma_m)$. The QR decomposition of $p_i(A_{i-1})$ determines the unitary similarity transformation that yields the next iterate:

$$
\begin{aligned}
p_i(A_{i-1}) &= Q_i R_i, \qquad \text{(QR decomposition)} \\
A_i &\leftarrow Q_i^{\mathsf{H}} A_{i-1} Q_i.
\end{aligned}
\tag{1}
$$

Any practically viable implementation of (1) contains at least two further ingredients: initial reduction to condensed form and deflation.

In the following, we assume that $A$ is already in upper Hessenberg form [14]. It is well known that this condensed form is preserved during the iteration (1) and greatly helps reduce its computational cost. To be more precise, the implicit Q theorem [14] implies that if the Hessenberg matrix $A_{i-1}$ is unreduced (all subdiagonal entries are different from zero) then (1) is equivalent to reducing $V^{\mathsf{H}} A_{i-1} V$ back to Hessenberg form, where $V$ is a unitary matrix that maps the first column of $p_i(A_{i-1})$ to a scalar multiple of the first unit vector $e_1$. Such an

implicit shifted QR iteration requires only $O(mn^2)$ flops (floating point operations), compared to $O(mn^3)$ flops needed by a literal implementation of (1).

As the QR algorithm proceeds, one or more subdiagonal entries of $A_i$ are expected to approach zero. For example, if $\sigma_1, \ldots, \sigma_m$ are chosen to be the eigenvalues of the trailing $m \times m$ principal submatrix of the current iterate, then – under some mild extra assumptions – the $(n - m + 1, n - m)$ subdiagonal entry of $A_i$ converges quadratically to zero [30]. The classical deflation criterion is to consider a subdiagonal element $a_{l,l+1}^{(i)}$ negligible if it satisfies

$$\left| a_{l,l+1}^{(i)} \right| \leq \mathbf{u} \left( \left| a_{l,l}^{(i)} \right| + \left| a_{l+1,l+1}^{(i)} \right| \right), \tag{2}$$

where $a_{kl}^{(i)}$ denotes the $(k, l)$ entry of $A_i$ and $\mathbf{u}$ the unit roundoff. A negligible subdiagonal entry is set to zero, effectively bringing $A_i$ to block upper triangular form:

$$A_i = \begin{bmatrix} A_{11}^{(i)} & A_{12}^{(i)} \\ 0 & A_{22}^{(i)} \end{bmatrix}, \quad A_{11}^{(i)} \in \mathbb{C}^{l \times l}, \quad A_{22}^{(i)} \in \mathbb{C}^{(n-l) \times (n-l)}.$$

This allows one to apply all subsequent QR iterations to the diagonal blocks $A_{11}^{(i)}$ and $A_{22}^{(i)}$ separately and therefore deflates the problem of computing the Schur decomposition of an $n \times n$ matrix into two smaller problems. The QR algorithm is said to have converged when all deflated diagonal blocks are $1 \times 1$.

The state-of-the-art LAPACK implementation of the QR algorithm attains high performance by making use of level 3 BLAS operations [6, 20] and employing additional deflation criteria going far beyond the classical criterion (2). Specifically, the *aggressive early deflation* strategy developed by Braman, Byers, and Mathias [7] often detects converged eigenvalues much earlier than (2) and therefore significantly decreases the overall number of QR iterations needed until convergence. In this paper, we approach this deflation technique from a rather different direction, based on Krylov subspace relations implicitly maintained during the QR algorithm. It turns out that aggressive early deflation amounts to finding and extracting converged Ritz pairs from a Krylov subspace $\mathcal{K}_w(A^{\mathsf{H}}, u_n)$, where $u_n$ denotes the last column of the accumulated unitary transformation matrix. This not only complements the analyses in [7, 32], partially explaining the remarkable success of aggressive early deflation, but also allows for improved convergence bounds. In particular, we can combine the classical convergence theory of the QR algorithm [30] with the convergence of Krylov subspaces to an invariant subspace [4, 5, 24]. The obtained convergence bounds clearly exhibit the benefits of aggressive early deflation.

The rest of this paper is organized as follows. In Section 2, we explore different Krylov subspaces associated with the QR algorithm. The Krylov-Schur algorithm, a reliable means to extract and lock converged Ritz pairs from these Krylov subspaces, is recalled in Section 3. Reinterpreting this algorithm in terms of unitary transformations on the QR iterate $A_i$ in Section 4 reveals its equivalence to aggressive early deflation. Finally, in Section 5 we use this relationship to derive convergence bounds for the QR algorithm with aggressive early deflation. In the following, $\| \cdot \|$ denotes the 2-norm of a vector or matrix while $\| \cdot \|_F$ denotes the Frobenius norm of a matrix.

## 2    Krylov subspace relations

The QR algorithm can be viewed as a nested subspace iteration [3, 8, 23, 28, 31]. Well suited for theoretical purposes, this approach forms the basis of the elegant convergence theory developed by Watkins and Elsner [30]. Starting with a linear subspace $\mathcal{S}_0 \subseteq \mathbb{C}^n$, it can be shown that the QR iteration (1) effects a subspace iteration of the form

$$\mathcal{S}_i = p_i(A)\mathcal{S}_{i-1}, \qquad i = 1, 2, \ldots. \tag{3}$$

If we let $\hat{p}_i = p_i p_{i-1} \cdots p_1$ denote the product of the shift polynomials then

$$\mathcal{S}_i = \hat{p}_i(A)\mathcal{S}_0, \qquad i = 1, 2, \ldots. \tag{4}$$

Setting $A_0 \equiv A$, we can define

$$\mathcal{S}_0 = \operatorname{span}\{e_1, e_2, \ldots, e_k\},$$

where $k$ satisfies $1 \le k \le n$ and $e_j$ denotes the $j$th unit vector of appropriate length. It is important to note that (3) and (4) hold for all $k$ *simultaneously* (giving rise to a *nested subspace iteration*).

To avoid technical difficulties, we assume for the rest of this section that each $p_i(A)$ is nonsingular (a singular $p_i(A)$ results in sudden convergence), which is equivalent to requiring that none of the zeros of $p_i$ coincides with an eigenvalue of $A$. Then the concrete relation of (3) to the QR iteration (1) is revealed by

$$\mathcal{S}_i = \operatorname{span}\{\hat{Q}_i e_1, \hat{Q}_i e_2, \ldots, \hat{Q}_i e_k\} \tag{5}$$

for $i \ge 0$ with $\hat{Q}_i = Q_1 \cdots Q_i$. Here, $Q_1, \ldots, Q_i$ are the unitary matrices computed during the QR iteration while $\hat{Q}_0$ is defined to be the identity $I_n$. A simple way to show (5) is to note that (1) implies a QR decomposition

$$\hat{p}_i(A) = \hat{Q}_i(R_i R_{i-1} \cdots R_1) \tag{6}$$

with nonsingular, upper triangular $R_i R_{i-1} \cdots R_1$.

The circumstance that the implicit shifted QR algorithm operates on Hessenberg matrices, links it intimately to Krylov subspace methods, see, e.g., [29] for a recent discussion. Additionally assuming that $A$ is in unreduced Hessenberg form, it is well known that

$$\mathcal{S}_i = \mathcal{K}_k(A, u_1) = \operatorname{span}\{u_1, Au_1, \ldots, A^{k-1}u_1\}, \tag{7}$$

where $u_1$ denotes the first column of $\hat{Q}_i$. In fact, if we let $u_j$ denote the $j$th column of $\hat{Q}_i$ and partition

$$A_i = \hat{Q}_i^{\mathsf{H}} A \hat{Q}_i = \begin{array}{c} \\ k \\ 1 \\ n-k-1 \end{array} \overset{\begin{array}{ccc} k & 1 & n-k-1 \end{array}}{\begin{bmatrix} H_{11} & H_{12} & H_{13} \\ H_{21} & H_{22} & H_{23} \\ 0 & H_{32} & H_{33} \end{bmatrix}} \tag{8}$$

then trivially

$$A[u_1, u_2, \ldots, u_k] = [u_1, u_2, \ldots, u_k]H_{11} + u_{k+1}H_{21}. \tag{9}$$

Under the given assumptions, $A_i$ is in unreduced Hessenberg form and hence the relation (9) happens to be an unreduced Arnoldi decomposition, which implies (7) by [25, Thm. 5.1.1].

Although occasionally mentioned in the literature [10, 28], it is less well known that $\mathcal{S}_i^\perp$, the orthogonal complement of $\mathcal{S}_i$, is also a Krylov subspace. To show this, we employ the "flip transpose" of a matrix. Given an $l \times m$ matrix $B$ let $B^\mathsf{F} = F_m B^\mathsf{H} F_l$, where $F_j$ denotes the $j \times j$ flip matrix, having ones on the anti-diagonal and zeros everywhere else. It can be directly seen that a square matrix $B^\mathsf{F}$ is in upper Hessenberg (triangular) form if and only if $B$ is in upper Hessenberg (triangular) form. Setting $w = n - k - 1$, the partitioning (8) immediately implies

$$A^\mathsf{H}[u_{n-w+1}, \dots, u_{n-1}, u_n] = [u_{n-w+1}, \dots, u_{n-1}, u_n]H_{33}^\mathsf{H} + u_{n-w}H_{32}^\mathsf{H}$$

and, after applying the flip matrix $F \equiv F_w$,

$$A^\mathsf{H}[u_n, u_{n-1}, \dots, u_{n-w+1}] = [u_n, u_{n-1}, \dots, u_{n-w+1}]H_{33}^\mathsf{F} + u_{n-w}H_{32}^\mathsf{H}F. \tag{10}$$

Again, (10) is an unreduced Arnoldi decomposition and therefore

$$\mathrm{span}\{u_n, u_{n-1}, \dots, u_{n-w+1}\} = \mathcal{K}_w(A^\mathsf{H}, u_n)$$

holds for every $w = 1, \dots, n$. Adjusting $w$ to $w = n - k$, this proves that $\mathcal{S}_i^\perp$ indeed coincides with the Krylov subspace $\mathcal{K}_{n-k}(A^\mathsf{H}, u_n)$. Since $\hat{p}_i(A)$ is invertible, (6) shows that $u_n$ is parallel to $(\hat{p}_i(A)^{-1})^\mathsf{H}e_n =: \hat{p}_i(A)^{-\mathsf{H}}e_n$. Therefore, using the fact that $A$ and $\hat{p}_i(A)^{-1}$ commute,

$$\mathcal{K}_w(A^\mathsf{H}, u_n) = \mathcal{K}_w(A^\mathsf{H}, \hat{p}_i(A)^{-\mathsf{H}}e_n) = \hat{p}_i(A)^{-\mathsf{H}}\mathcal{K}_w(A^\mathsf{H}, e_n).$$

## 3 Extracting Ritz pairs from Krylov subspaces

Given an Arnoldi decomposition of the form (10), a conventional way to extract approximations to eigenvectors from the corresponding Krylov subspace is to compute Ritz pairs and check their residuals. A *Ritz value* $\lambda$ from the subspace $\mathcal{K}_w(A^\mathsf{H}, u_n)$ is defined as an eigenvalue of the $w \times w$ matrix $H_{33}^\mathsf{F}$. The corresponding *Ritz vector* is given by $x = U_w z$, where $z$ is an eigenvector of $H_{33}^\mathsf{F}$ belonging to $\lambda$ and $U_w = [u_{n-w+1}, \dots, u_{n-1}, u_n]$. Taken together, $(\lambda, x)$ form a so-called *Ritz pair*.

With the normalization $\|x\| = \|z\| = 1$, a Ritz pair is usually regarded as converged towards an eigenpair of $A^\mathsf{H}$ if the norm of the residual $r = A^\mathsf{H}x - \lambda x$ is sufficiently small. A practical criterion for deciding upon smallness can be found in the ARPACK manual [22, Sec. 4.6]:

$$\|r\| \le \max\{\mathbf{u}\|H_{33}^\mathsf{F}\|_F, \mathsf{tol} \times |\lambda|\}, \tag{11}$$

where $\mathbf{u}$ denotes the unit roundoff and $\mathsf{tol}$ is a tolerance chosen by the user. Note that (10) implies $r = (H_{32}^\mathsf{H}Fz)u_{n-w}$ and hence $\|r\| = |H_{32}^\mathsf{H}Fz|$. For $\mathsf{tol} = 0$, the criterion (11) yields normwise backward stability: $(\lambda, x)$ is the exact eigenpair of the perturbed matrix $A^\mathsf{H} + (\triangle A)^\mathsf{H}$, where

$$\triangle A = -xr^\mathsf{H} = -(z^\mathsf{H}FH_{32})(U_w z)u_{n-w}^\mathsf{H} \tag{12}$$

satisfies $\|\triangle A\|_F = \|r\| \le \mathbf{u}\|H_{33}^\mathsf{F}\|_F \le \mathbf{u}\|A\|_F$.

### 3.1 Locking converged Ritz values

Stewart's *Krylov-Schur algorithm* [26] provides a numerically reliable means to detect, extract and lock converged Ritz pairs. For some Ritz value $\lambda$, an ordered Schur decomposition [14] of $H_{33}^{\mathsf{F}}$ is computed such that $\lambda$ appears in the top left corner of the triangular factor:

$$V^{\mathsf{H}} H_{33}^{\mathsf{F}} V = \begin{bmatrix} \lambda & T_{12} \\ 0 & T_{22} \end{bmatrix}, \tag{13}$$

where $V \in \mathbb{C}^{w \times w}$ is unitary and $T_{22} \in \mathbb{C}^{(w-1) \times (w-1)}$. A corresponding Ritz vector is given by $x = U_w z$ with $z = V e_1$. If (11) is satisfied, we partition $U_w V = [x, \widehat{U}_{w-1}]$ and

$$H_{32}^{\mathsf{H}} F V = [\bar{s}_1, s_2^{\mathsf{H}}], \tag{14}$$

where $s_1 = z^{\mathsf{H}} F H_{32}$. Together with (10), this implies

$$A^{\mathsf{H}}[x, \widehat{U}_{w-1}] = [x, \widehat{U}_{w-1}] \begin{bmatrix} \lambda & T_{12} \\ 0 & T_{22} \end{bmatrix} + u_{n-w}[\bar{s}_1, s_2^{\mathsf{H}}]. \tag{15}$$

If (11) is not satisfied, we can test any other Ritz pair by reordering the eigenvalues in the Schur form (13).

Continuing the described process for the Ritz values contained in $T_{22}$ eventually leads to a decomposition of the form

$$A^{\mathsf{H}}[X, \widehat{U}_{w-d}] = [X, \widehat{U}_{w-d}] \begin{bmatrix} \widehat{T}_{11} & \widehat{T}_{12} \\ 0 & \widehat{T}_{22} \end{bmatrix} + u_{n-w}[\hat{s}_1^{\mathsf{H}}, \hat{s}_2^{\mathsf{H}}], \tag{16}$$

where the upper triangular matrix $\widehat{T}_{11}$ contains all $d$ converged Ritz values on its diagonal and

$$\|\hat{s}_1\| \leq \sqrt{d} \max\{\mathbf{u}, \mathsf{tol}\} \|H_{33}^{\mathsf{F}}\|_F. \tag{17}$$

The matrices $U_w$ and $[X, \widehat{U}_{w-d}]$ span the same space, i.e., there is a unitary matrix $\widehat{V}$ such that $[X, \widehat{U}_{w-d}] = U_w \widehat{V}$.

### 3.2 Restoring the Arnoldi decomposition

Note that (16) is not a standard Arnoldi decomposition, which would require $[\hat{s}_1^{\mathsf{H}}, \hat{s}_2^{\mathsf{H}}]$ to be a multiple of $e_w^{\mathsf{H}}$. The vector $\hat{s}_1^{\mathsf{H}}$ is negligible but $\hat{s}_2^{\mathsf{H}}$ is not. Therefore, the last (optional) step of the Krylov-Schur algorithm consists of transforming $\hat{s}_2^{\mathsf{H}}$ back to a multiple of $e_{w-d}^{\mathsf{H}}$ while maintaining $\begin{bmatrix} \widehat{T}_{11} & \widehat{T}_{12} \\ 0 & \widehat{T}_{22} \end{bmatrix}$ in upper Hessenberg form.

To achieve this goal, first a unitary matrix $V_1$ is computed such that $V_1^{\mathsf{H}} \hat{s}_2 = \beta e_{w-d}$ with $|\beta| = \|\hat{s}_2\|$. By a row-oriented version of the usual Hessenberg reduction algorithm [25, Pg. 312], a unitary matrix $V_2 = \widetilde{V}_2 \oplus 1$ can be computed such that $\widetilde{T}_{22} = V_2^{\mathsf{H}}(V_1^{\mathsf{H}} \widehat{T}_{22} V_1) V_2$ is again in Hessenberg form. Setting $\widetilde{T}_{12} = \widehat{T}_{12} V_1 V_2$ and $\widetilde{U}_{w-d} = \widehat{U}_{w-d} V_1 V_2$ finally yields

$$A^{\mathsf{H}}[X, \widetilde{U}_{w-d}] = [X, \widetilde{U}_{w-d}] \begin{bmatrix} \widehat{T}_{11} & \widetilde{T}_{12} \\ 0 & \widetilde{T}_{22} \end{bmatrix} + u_{n-w}[\hat{s}_1^{\mathsf{H}}, \bar{\beta} e_{w-d}^{\mathsf{H}}], \tag{18}$$

which becomes an Arnoldi decomposition after setting $\hat{s}_1$ to zero.

## 4   Aggressive early deflation

In this section, we reinterpret the extraction of Ritz pairs from the Arnoldi decomposition (10) in terms of transformations operating on the $n \times n$ Hessenberg matrix $A_i$. Let us recall the partitioning (8) with $w = n - k - 1$:

$$
A_i = \hat{Q}_i^{\mathsf{H}} A \hat{Q}_i =
\begin{matrix}
 & & n-w-1 & 1 & w \\
 w & \\
 1 & \\
 n-w-1 &
\end{matrix}
\begin{bmatrix}
H_{11} & H_{12} & H_{13} \\
H_{21} & H_{22} & H_{23} \\
0 & H_{32} & H_{33}
\end{bmatrix}.
$$

In the context of the QR algorithm, we refer to $w$ as the *deflation window size* and assume $w \ll n$ to make sure that the cost of the extraction process remains modest in comparison to the cost of a QR iteration.

### 4.1   Locking converged Ritz values

Note that all variables used in this section refer precisely to the same quantities introduced in Section 3. A Ritz value $\lambda$ has been defined to be an eigenvalue of $H_{33}^{\mathsf{F}}$ and $z$ denotes a corresponding normalized eigenvector. It directly follows that $\bar{\lambda}$ is an eigenvalue of $H_{33}$ having $F_w z$ as a corresponding *left* eigenvector. Setting $V_{\mathsf{F}} = F_w V F_w$, the ordered Schur decomposition (13) implies

$$
V_{\mathsf{F}}^{\mathsf{H}} H_{33} V_{\mathsf{F}} = F_w \left(V^{\mathsf{H}} H_{33}^{\mathsf{F}} V\right)^{\mathsf{H}} F_w = F_w
\begin{bmatrix}
\bar{\lambda} & 0 \\
T_{12}^{\mathsf{H}} & T_{22}^{\mathsf{H}}
\end{bmatrix}
F_w =
\begin{bmatrix}
T_{22}^{\mathsf{F}} & F_{w-1} T_{12}^{\mathsf{H}} \\
0 & \bar{\lambda}
\end{bmatrix},
\tag{19}
$$

which is an ordered Schur decomposition for $H_{33}$. Moreover, it follows from (14) that

$$
V_{\mathsf{F}}^{\mathsf{H}} H_{32} = F_w V^{\mathsf{H}} (F_w H_{32}) = F_w
\begin{bmatrix}
s_1 \\
s_2
\end{bmatrix}
=
\begin{bmatrix}
F_{w-1} s_2 \\
s_1
\end{bmatrix}.
$$

These two relations yield

$$
(I \oplus V_{\mathsf{F}})^{\mathsf{H}} A_i (I \oplus V_{\mathsf{F}}) =
\begin{bmatrix}
H_{11} & H_{12} & \check{H}_{13} & \check{H}_{14} \\
H_{21} & H_{22} & \check{H}_{23} & \check{H}_{24} \\
0 & F_{w-1} s_2 & T_{22}^{\mathsf{F}} & T_{12}^{\mathsf{F}} \\
0 & s_1 & 0 & \bar{\lambda}
\end{bmatrix},
$$

where we define $\begin{bmatrix} \check{H}_{13} & \check{H}_{14} \\ \check{H}_{23} & \check{H}_{24} \end{bmatrix} := \begin{bmatrix} H_{13} \\ H_{23} \end{bmatrix} V_{\mathsf{F}}$. Note that the vector $\begin{bmatrix} F_{w-1} s_2 \\ s_1 \end{bmatrix}$ is called the *spike* in [7] .

Let us recall that the residual $r$ of the Ritz value $\lambda$ satisfies $r = (H_{32}^{\mathsf{H}} F z) u_{n-w} = \bar{s}_1 u_{n-w}$ and hence $\|r\| = |s_1|$. Thus if (11) is satisfied for $\mathsf{tol} = 0$, the trailing spike element can be safely set to zero and $\bar{\lambda}$ can be regarded as a computed eigenvalue of $A$ without spoiling the numerical backward stability of the QR algorithm. If $|s_1| > \mathsf{u}\|H_{33}\|_F$, we can test any other eigenvalue of $H_{33}$ by considering a differently ordered Schur decomposition. This is equivalent to the search for converged Ritz values, since the ordered Schur decompositions of $H_{33}$ and $H_{33}^{\mathsf{F}}$ are connected to each other in the one-to-one relationship (19).

Extracting and locking further converged Ritz values of $H_{33}^{\mathsf{F}}$ as described in Section 3 eventually yields a unitary matrix $\widehat{V}_{\mathsf{F}} \in \mathbb{C}^{w \times w}$ such that

$$\widehat{A}_i = (I \oplus \widehat{V}_{\mathsf{F}})^{\mathsf{H}} A_i (I \oplus \widehat{V}_{\mathsf{F}}) = \begin{bmatrix} H_{11} & H_{12} & \widehat{H}_{13} & \widehat{H}_{14} \\ H_{21} & H_{22} & \widehat{H}_{23} & \widehat{H}_{24} \\ 0 & F_{w-d}\hat{s}_2 & \widehat{T}_{22}^{\mathsf{F}} & \widehat{T}_{12}^{\mathsf{F}} \\ 0 & F_d\hat{s}_1 & 0 & \widehat{T}_{11}^{\mathsf{F}} \end{bmatrix}, \tag{20}$$

where the diagonal of the upper triangular matrix $\widehat{T}_{11}^{\mathsf{F}}$ contains the complex conjugates of all converged Ritz values. The $d$ trailing spike elements in $F_d\hat{s}_1$ are all negligible, since (17) amounts to $\|F_d\hat{s}_1\| \leq \sqrt{d}\mathbf{u}\|H_{33}\|_F \leq \sqrt{d}\mathbf{u}\|A\|_F$ for $\mathsf{tol} = 0$.

## 4.2  Restoring the Hessenberg form

To continue the QR algorithm, the matrix $\widehat{A}_i$ in (20) needs to be restored to Hessenberg form. This is exactly what is achieved by the unitary transformations defined in Section 3.2. Setting $V_{1,\mathsf{F}} = F_{w-d}V_1 F_{w-d}$ and $V_{2,\mathsf{F}} = F_{w-d}V_2 F_{w-d}$ yields

$$V_{1,\mathsf{F}}^{\mathsf{H}}(F_{w-d}\hat{s}_2) = \beta e_1, \quad V_{2,\mathsf{F}}^{\mathsf{H}}(V_{1,\mathsf{F}}^{\mathsf{H}}\widehat{T}_{22}^{\mathsf{F}}V_{1,\mathsf{F}})V_{2,\mathsf{F}} = \widetilde{T}_{22}^{\mathsf{F}},$$

where $\widetilde{T}_{22}^{\mathsf{F}}$ is in upper Hessenberg form. Hence,

$$\widetilde{A}_i = (I \oplus V_{1,\mathsf{F}}V_{2,\mathsf{F}} \oplus I_d)^{\mathsf{H}} \widehat{A}_i (I \oplus V_{1,\mathsf{F}}V_{2,\mathsf{F}} \oplus I_d) = \begin{bmatrix} H_{11} & H_{12} & \widetilde{H}_{13} & \widehat{H}_{14} \\ H_{21} & H_{22} & \widetilde{H}_{23} & \widehat{H}_{24} \\ 0 & \beta e_1 & \widetilde{T}_{22}^{\mathsf{F}} & \widetilde{T}_{12}^{\mathsf{F}} \\ 0 & F_d\hat{s}_1 & 0 & \widehat{T}_{11}^{\mathsf{F}} \end{bmatrix}$$

with $\begin{bmatrix} \widetilde{H}_{13} \\ \widetilde{H}_{23} \end{bmatrix} = \begin{bmatrix} \widetilde{H}_{13} \\ \widetilde{H}_{23} \end{bmatrix} V_{1,\mathsf{F}}V_{2,\mathsf{F}}$. Setting $\hat{s}_1$ to zero returns $\widetilde{A}_i$ to Hessenberg form and the QR algorithm can be continued on the leading $(n - d) \times (n - d)$ principal submatrix of $\widetilde{A}_i$.

## 4.3  Comparison to classical deflation

The reader is invited to check that the presented deflation algorithm obtained via the extraction of Ritz pairs from the Krylov subspace $\mathcal{K}_w(A^{\mathsf{H}}, u_n)$ is precisely the same as the aggressive early deflation algorithm originally described in [7]. Let us contemplate what the different deflation criteria for the QR algorithm mean in terms of Krylov subspaces.

The classical deflation criterion (2) tests the smallness of $|a_{l,l+1}^{(i)}|$ for each $l = 1, \ldots, n-1$. From the Arnoldi decomposition (10) it is immediately clear that $|a_{l,l+1}^{(i)}|$ is the minimal norm of a backward error matrix $\triangle A$ such that $\mathcal{K}_{n-l}(A^{\mathsf{H}}, u_n)$ becomes an invariant subspace of $(A + \triangle A)^{\mathsf{H}}$. In other words, classical deflation considers the nested sequence of Krylov subspaces

$$\mathcal{K}_1(A^{\mathsf{H}}, u_n), \ \mathcal{K}_2(A^{\mathsf{H}}, u_n), \ \ldots, \ \mathcal{K}_{n-1}(A^{\mathsf{H}}, u_n)$$

and checks whether any of them is a good approximation to an invariant subspace of $A^{\mathsf{H}}$ as a whole. If the shifts in the QR algorithm are chosen as the eigenvalues of the $m \times m$ trailing submatrix, a choice sometimes called *Francis shifts*, then this is most likely to happen for $\mathcal{K}_{n-l}(A^{\mathsf{H}}, u_n)$ with $n-l \leq m$, as the convergence theory [30] states that $\mathcal{K}_m(A^{\mathsf{H}}, u_n)$ converges

locally quadratically to an invariant subspace of $A^{\mathsf{H}}$. Roughly speaking (neglecting the effects of ill-conditioning), the approximation quality of $\mathcal{K}_{n-l}(A^{\mathsf{H}}, u_n)$ as a whole is determined by the poorest Ritz vector approximation that can be extracted from $\mathcal{K}_{n-l}(A^{\mathsf{H}}, u_n)$. Hence, slowly converging Ritz vectors hinder the deflation of other, quickly converging Ritz vectors. Aggressive early deflation is fundamentally different and avoids this effect. Only one Krylov subspace $\mathcal{K}_w(A^{\mathsf{H}}, u_n)$ for some fixed value of $w$ is considered. Moreover, not the convergence of $\mathcal{K}_w(A^{\mathsf{H}}, u_n)$ as a whole but the convergence of each individual Ritz vector to an eigenvector of $A^{\mathsf{H}}$ is checked. Provided that $A$ has distinct eigenvalues, all Ritz vectors from the Krylov subspace $\mathcal{K}_m(A^{\mathsf{H}}, u_n)$ converge locally quadratically, but some may converge at a significantly faster rate and can be deflated much earlier. This is demonstrated by the following example.

**Example 1** *We applied the QR algorithm with $m = 4$ Francis shifts to a $250 \times 250$ random matrix generated by the following* MATLAB *commands:*

```
randn('state',0); A = hess(randn(250)+1i*randn(250));
```

| | Classical deflation: $\lvert a_{l+1,l}^{(i)} \rvert$, $l =$ | | | | Aggressive early deflation: min $\lVert r \rVert$, $w =$ | | | |
|---|---|---|---|---|---|---|---|---|
| | $n-4$ | $n-3$ | $n-2$ | $n-1$ | 2 | 4 | 6 | 8 |
| $i=0$ | $2.3\times10^{+0}$ | $1.7\times10^{+0}$ | $1.9\times10^{+0}$ | $2.1\times10^{+0}$ | $1.3\times10^{+0}$ | $2.2\times10^{-1}$ | $9.1\times10^{-2}$ | $7.3\times10^{-2}$ |
| $i=1$ | $9.1\times10^{-1}$ | $1.2\times10^{+0}$ | $1.8\times10^{+0}$ | $2.7\times10^{+0}$ | $1.2\times10^{+0}$ | $6.0\times10^{-2}$ | $1.5\times10^{-2}$ | $6.4\times10^{-3}$ |
| $i=2$ | $4.9\times10^{-1}$ | $5.8\times10^{-1}$ | $1.3\times10^{+0}$ | $6.0\times10^{-1}$ | $2.1\times10^{-1}$ | $2.0\times10^{-3}$ | $8.6\times10^{-5}$ | $5.0\times10^{-5}$ |
| $i=3$ | $4.6\times10^{-2}$ | $9.2\times10^{-2}$ | $1.7\times10^{+0}$ | $1.8\times10^{-2}$ | $8.5\times10^{-3}$ | $1.2\times10^{-6}$ | $9.6\times10^{-9}$ | $7.7\times10^{-9}$ |

Table 1: Magnitudes of trailing subdiagonal elements (columns 2–5) and minimal residuals for the Ritz pairs extracted from $\mathcal{K}_w(A^{\mathsf{H}}, u_n)$ (columns 6–9) after $i$ QR iterations with 4 Francis shifts applied to the matrix from Example 1.

*Table 1 compares classical with aggressive early deflation and clearly exhibits the advantages of the latter. For example, consider the Krylov subspace $\mathcal{K}_4(A^{\mathsf{H}}, u_n)$ after $i = 3$ QR iterations. The norms of the residuals for the four corresponding Ritz pairs are $3.8\times10^{-2}$, $1.3\times10^{-3}$, $6.0\times10^{-4}$, and $1.2\times10^{-6}$. The magnitude of the corresponding subdiagonal element, $\lvert a_{n-4,n-3}^{(3)} \rvert = 4.6\times10^{-2}$, is nearly the maximum of these numbers, demonstrating that the convergence of $\mathcal{K}_4(A^{\mathsf{H}}, u_n)$ as a whole is determined by the poorest Ritz pair approximation.*

In Table 1, aggressive early deflation shows dramatic improvements already for $w = m$, i.e., when the size of the deflation window coincides with the number of Francis shifts. Choosing $w > m$ enlarges the Krylov subspace and adds a Krylov subspace acceleration to the quadratically converging Ritz vectors. This is explored in more detail in the next section.

**Remark 2** *If each QR iteration is based on a single shift that is defined to be the $(n, n)$ entry of the current iterate then it is well known that $u_n$ undergoes a Rayleigh-quotient iteration, see, e.g., [25, Sec. 2.1]. To be more specific, let $u_n^{(i-1)}$ denote the last column of the accumulated unitary transformation matrix $U_{i-1}$ satisfying $A_{i-1} = U_{i-1}^{\mathsf{H}} A U_{i-1}$. Then*

$$u_n^{(i)} = \frac{\left(A^{\mathsf{H}} - \sigma_1^{(i)}\right)^{-1} u_n^{(i-1)}}{\left\lVert \left(A^{\mathsf{H}} - \sigma_1^{(i)}\right)^{-1} u_n^{(i-1)} \right\rVert},$$

*provided that $\sigma_1^{(i)} = u_n^{(i-1)\mathsf{H}} A u_n^{(i-1)}$ is not an eigenvalue of $A^{\mathsf{H}}$. Using aggressive early deflation, we deflate converged Ritz pairs from the Krylov subspace $\mathcal{K}_w(A^{\mathsf{H}}, u_n^{(i)})$. This shows*

*that the single-shifted QR iteration equipped with such a deflation strategy is in fact a Krylov subspace method where the starting vector undergoes a quadratically convergent iteration.*

There are situations for which the classical criterion detects deflations that go undetected if only aggressive early deflation is used. Probably the most practically relevant one is that linear convergence phenomena occur if the shifts vary little in the course of several QR iterations. The typical consequence is that one or more of the leading subdiagonal elements of $A_i$ approach zero. As $w \ll n$, the Krylov subspace $\mathcal{K}_w(A^\mathsf{H}, u_n)$ will not benefit from this effect. As a remedy, one could additionally use a variant of aggressive early deflation that considers Ritz pairs from $\mathcal{K}_w(A, u_1)$. However, numerical experiments reported in [17] reveal that these linear convergence phenomena seem to be too rare to justify the extra computational effort.

A quite different situation is pointed out in [7, Sec. 2.5] and it is interesting to see this example in the light of Krylov subspaces. Let

$$A = \begin{bmatrix} 2 & 3 & 4 & 5 & 6 \\ 1 & 2 & 0 & 0 & 1 \\ 0 & 1 & 2 & 0 & 0 \\ 0 & 0 & \varepsilon & 2 & 0 \\ 0 & 0 & 0 & 1 & 2 \end{bmatrix}, \tag{21}$$

where $0 < \varepsilon \ll 1$. The norm of the residual for every Ritz pair from $\mathcal{K}_4(A^\mathsf{H}, e_5)$ is approximately given by $\sqrt{\varepsilon}$, even though a perturbation of norm $\varepsilon$ *within* the deflation window deflates two eigenvalues. Caused by the high nonnormality of $A(2:5, 2:5)$, standard Ritz pair extraction fails to detect the converged left eigenvector contained in $\mathcal{K}_4(A^\mathsf{H}, e_5)$. Refined Ritz pairs [15, 16] aim to avoid this effect. For each Ritz value $\lambda$, the refined Ritz vector is chosen to minimize the norm of the corresponding residual. Unfortunately, for our example each refined Ritz pair has residual norm of about $\sqrt{\varepsilon/2}$. In [18], a computational procedure based on the distance to uncontrollability was described that optimizes both the refined Ritz vector *and* value. For our example, this leads to a refined Ritz pair whose residual norm nearly attains $\varepsilon$. So, in principle, refined Ritz pairs can be used to improve aggressive early deflation even further. However, preliminary numerical experiments with several matrices from the Matrix Market collection [2] indicate that a situation like (21) occurs rarely in practice, at least too seldom to justify the extra computational effort needed for extracting refined Ritz pairs.

## 5 Convergence bounds

To discuss the impact of aggressive early deflation on the convergence of the QR algorithm, let us return to the notation of Section 2. Watkins and Elsner [30, Theorem 6.3] have shown that $\mathcal{S}_i = \mathcal{K}_{n-m}(A, u_1)$ converges locally quadratically to an $(n-m)$-dimensional invariant subspace $\mathcal{X}$ of $A$, provided that $A$ has distinct eigenvalues and that in each iteration $m$ Francis shifts $\sigma_1, \ldots, \sigma_m$ are chosen. To formalize this statement, let us employ the *(containment) gap* between two subspaces $\mathcal{S}$ and $\mathcal{T}$ (not necessarily of the same dimension) defined as

$$d(\mathcal{S}, \mathcal{T}) := \sup_{\substack{s \in \mathcal{S} \\ \|s\|=1}} \inf_{t \in \mathcal{T}} \|s - t\|.$$

For convenience, we write $d(y, \mathcal{T})$ if $\mathcal{S}$ is spanned by a single vector $y$. With this notation, the convergence statement reads as follows:

> *If $d(\mathcal{S}_i, \mathcal{X}) \to 0$ converges to zero as $i \to \infty$ then this convergence is quadratic.*

At the time of writing, the task of finding conditions for global convergence, which seems to hold almost always in practice, remains open and the results in this paper have nothing to contribute to this question.

Relevant properties of the gap $d(\cdot, \cdot)$ are summarized, e.g., in [4, 30]. We have $d(\mathcal{S}, \mathcal{T}) = 0$ if and only if $\mathcal{S} \subseteq \mathcal{T}$. If $\Pi_{\mathcal{S}}$ and $\Pi_{\mathcal{T}}$ denote orthogonal projections onto $\mathcal{S}$ and $\mathcal{T}$, respectively, then $d(\mathcal{S}, \mathcal{T}) = \|(I - \Pi_{\mathcal{T}})\Pi_{\mathcal{S}}\|$. This readily implies $d(\mathcal{T}^{\perp}, \mathcal{S}^{\perp}) = d(\mathcal{S}, \mathcal{T})$. The last property allows us to replace $d(\mathcal{S}_i, \mathcal{X})$ by

$$d(\mathcal{Y}, \mathcal{S}_i^{\perp}) = d(\mathcal{Y}, \mathcal{K}_m(A^{\mathsf{H}}, u_n)) = d(\mathcal{Y}, \hat{p}_i(A)^{-\mathsf{H}}\mathcal{K}(A^{\mathsf{H}}, e_n)),$$

with $\mathcal{Y} = \mathcal{X}^{\perp}$, in the convergence discussion of the QR algorithm. Note that $\mathcal{Y}$ is an invariant subspace of $A^{\mathsf{H}}$, or, equivalently, a *left* invariant subspace of $A$.

## 5.1 Convergence of eigenvectors

As discussed in Section 4.3, aggressive early deflation has two advantages: (1) approximation by Ritz vectors instead of the whole Krylov subspace, and (2) additional Krylov subspace acceleration. The convergence theory by Watkins and Elsner must be modified in order to accommodate both advantages; the following theorem addresses the first one.

**Theorem 3** *Let $A \in \mathbb{C}^{n \times n}$ have a block Schur decomposition*

$$Q^{\mathsf{H}}AQ = \begin{bmatrix} A_{11} & A_{12} \\ 0 & A_{22} \end{bmatrix}, \quad A_{11} \in \mathbb{C}^{(n-m) \times (n-m)}, \quad A_{22} \in \mathbb{C}^{m \times m},$$

*for some unitary matrix $Q \in \mathbb{C}^{n \times n}$, such that $\lambda(A_{11}) \cap \lambda(A_{22}) = \emptyset$. Choose a left eigenvector $y$ belonging to a simple eigenvalue $\lambda_2 \in \lambda(A_{22})$. Moreover, let $\mathcal{Y}$ denote the left invariant subspace of $A$ belonging to $\lambda(A_{22})$.*

*Consider a function $f$ that is analytic and nonzero on the spectrum of $A$. Then*

$$d(y, f(A)^{-\mathsf{H}}\mathcal{V}) \leq \frac{C_T}{1 - C_T\, d(\mathcal{Y}, \mathcal{V})}\, \|f(A_{11})^{-1}\|\, |f(\lambda_2)|\, d(y, \mathcal{V}), \tag{22}$$

*holds for any $m$-dimensional subspace $\mathcal{V}$ satisfying $C_T\, d(\mathcal{Y}, \mathcal{V}) < 1$, where*

$$C_T = \left(\|P_{A_{11}}\| + \sqrt{1 + \|P_{A_{11}}\|^2}\right)\left(\|P_{\lambda_2}\| + \sqrt{1 + \|P_{\lambda_2}\|^2}\right).$$

*Here, $P_{A_{11}} \in \mathbb{C}^{n \times n}$ denotes the spectral projector for $A$ associated with $\lambda(A_{11})$, and $P_{\lambda_2} \in \mathbb{C}^{m \times m}$ denotes the spectral projector for $A_{22}$ associated with $\lambda_2$.*

*Proof.* In the following, $\kappa(\cdot)$ denotes the 2-norm condition number of a matrix. By [11], there are invertible lower block triangular matrices $T_1, T_2$ with

$$\kappa(T_1) = \|P_{A_{11}}\| + \sqrt{1 + \|P_{A_{11}}\|^2}, \quad \kappa(T_2) = \|P_{\lambda_2}\| + \sqrt{1 + \|P_{\lambda_2}\|^2},$$

such that

$$T_1^{-1}\begin{bmatrix} A_{11} & A_{12} \\ 0 & A_{22} \end{bmatrix}T_1 = \begin{bmatrix} A_{11} & 0 \\ 0 & A_{22} \end{bmatrix}, \quad T_2^{-1}A_{22}T_2 = \begin{bmatrix} \widetilde{A}_{22} & 0 \\ 0 & \lambda_2 \end{bmatrix}.$$

The matrix $T := QT_1(I \oplus T_2)$ satisfies $\kappa(T) \leq \kappa(T_1)\kappa(T_2) = C_T$ and block diagonalizes $A$:

$$D := T^{-1}AT = A_{11} \oplus \tilde{A}_{22} \oplus \lambda_2. \tag{23}$$

In the following, we first derive a bound for the accordingly transformed subspaces $\widetilde{\mathcal{V}} = T^H\mathcal{V}$ and $f(D)^{-H}\widetilde{\mathcal{V}}$, which will then be turned into a bound for the original subspaces.

Since $\tilde{y} := T^H y$ is a left eigenvector of the block diagonal matrix $D = T^{-1}AT$ belonging to the simple eigenvalue $\lambda_2$, we must have $\tilde{y} = \beta e_n$ for some $\beta \in \mathbb{C}$. To simplify the description, we assume without loss of generality that $\tilde{y} = e_n$. Analogously, we have

$$\widetilde{\mathcal{Y}} := T^H\mathcal{Y} = \text{span} \begin{bmatrix} 0 \\ I_m \end{bmatrix}. \tag{24}$$

By [30, Lemma 4.1], the condition $C_T\, d(\mathcal{Y}, \mathcal{V}) < 1$ implies $d(\widetilde{\mathcal{Y}}, \widetilde{\mathcal{V}}) < 1$, Hence, no vector in $\widetilde{\mathcal{V}}$ is perpendicular to $\widetilde{\mathcal{Y}}$. Taking (24) into account, this means that no nonzero vector contained in $\widetilde{\mathcal{V}}$ can have $m$ trailing zero entries. Therefore we can choose a basis of the form $V = \begin{bmatrix} F \\ I_m \end{bmatrix}$ for $\widetilde{\mathcal{V}}$. The norm of $F$ determines the distance between $\widetilde{\mathcal{V}}$ and $\widetilde{\mathcal{Y}}$. Specifically, we have

$$\|F\| = d(\widetilde{\mathcal{Y}}, \widetilde{\mathcal{V}})/\sqrt{1 - d(\widetilde{\mathcal{Y}}, \widetilde{\mathcal{V}})^2}, \tag{25}$$

see, e.g., [27].

To obtain a bound on $d\big(e_n, f(D)^{-H}\widetilde{\mathcal{V}}\big)$, we select the particularly convenient vector

$$v_1 = Ve_m = \begin{bmatrix} Fe_m \\ e_m \end{bmatrix} \in \widetilde{\mathcal{V}}. \tag{26}$$

We have

$$f(D)^{-H}v_1 = \begin{bmatrix} f(A_{11})^{-H}Fe_m \\ 1/\overline{f(\lambda_2)}e_m \end{bmatrix},$$

which implies

$$d\big(e_n, f(D)^{-H}v_1\big) \leq \|f(A_{11})^{-H}Fe_m\||f(\lambda_2)| \leq \|f(A_{11})^{-1}\||f(\lambda_2)|\,\|e_n - v_1\|. \tag{27}$$

Note that $v_1$ is in general *not* the vector in $\widetilde{\mathcal{V}}$ that is closest to $e_n$, so we cannot simply replace $\|e_n - v_1\|$ by $d\big(e_n, \widetilde{\mathcal{V}}\big)$ in (27). In fact, the closest vector is the solution of the minimization problem $d(e_n, \widetilde{\mathcal{V}}) = \inf_{v \in \widetilde{\mathcal{V}}} \|e_n - v\|$ and takes the form $v_0 = V(V^H V)^{-1}V^H e_n = V(V^H V)^{-1}e_m$. Lemma 8 from Appendix A reveals the relationship

$$\|e_n - v_1\| \leq \frac{\sqrt{1 + \|F\|^2}}{\sqrt{1 + \|F\|^2} - \|F\|}\|e_n - v_0\| = \frac{\|e_n - v_0\|}{1 - d(\widetilde{\mathcal{Y}}, \widetilde{\mathcal{V}})}, \tag{28}$$

where we used (25) for the latter equality.

Combining (27) with (28) shows

$$d\big(e_n, f(D)^{-H}\widetilde{\mathcal{V}}\big) \leq d\big(e_n, f(D)^{-H}v_1\big) \leq \frac{\|f(A_{11})^{-1}\||f(\lambda_2)|}{1 - d(\widetilde{\mathcal{Y}}, \widetilde{\mathcal{V}})}d(\tilde{y}, \widetilde{\mathcal{V}}).$$

Recall that $y = T^{-H}e_n$, $\mathcal{V} = T^{-H}\widetilde{\mathcal{V}}$, and $f(A)^{-H}\mathcal{V} = T^{-H}f(D)^{-H}\widetilde{\mathcal{V}}$. Using Lemma 9 from Appendix A, we obtain

$$d\big(y, f(A)^{-H}\mathcal{V}\big) \leq C_T \frac{\|f(A_{11})^{-1}\|\,|f(\lambda_2)|}{1 - d(\widetilde{\mathcal{Y}}, \widetilde{\mathcal{V}})} d(y, \mathcal{V}),$$

using $\kappa(T) \leq C_T$. The proof is completed after applying the estimate $d(\widetilde{\mathcal{Y}}, \widetilde{\mathcal{V}}) \leq C_T\, d(\mathcal{Y}, \mathcal{V})$.
□

**Remark 4** *The condition $C_T\, d(\mathcal{Y}, \mathcal{K}_m(A^H, e_n)) < 1$ in Theorem 3 is an artifact of the proof technique and can be removed. Let $\mathcal{U}$ denote the right invariant subspace of $A$ belonging to $\lambda(A_{11})$. Then $\mathcal{U} \cap \mathcal{V} = \{0\}$ is sufficient to guarantee that the condition $d(\widetilde{\mathcal{Y}}, \widetilde{\mathcal{V}}) < 1$ needed in the proof of Theorem 3 is satisfied. Hence, there is a constant $C$ independent of $f$ such that*

$$d\big(y, f(A)^{-H}\mathcal{V}\big) \leq C\,\|f(A_{11})^{-1}\|\,|f(\lambda_2)|\,d(y, \mathcal{V}),$$

*even under this rather mild assumption. Note that $\mathcal{U} \cap \mathcal{V} = \{0\}$ is always satisfied if $A$ is in unreduced Hessenberg form and $\mathcal{V} = \mathcal{K}_m(A^H, e_n)$ [30].*

Theorem 3 is applied to the convergence analysis of the QR algorithm by setting $f = \hat{p}_i$ and $\mathcal{V} = \mathcal{K}_m(A^H, e_n)$. The classical convergence analysis [30, Lemma 4.4] provides bounds of the form

$$d\big(\mathcal{Y}, \hat{p}_i(A)^{-H}\mathcal{K}_m(A^H, e_n)\big) \leq C\,\|\hat{p}_i(A_{11})^{-1}\|\,\|\hat{p}_i(A_{22})\|\,d(\mathcal{Y}, \mathcal{K}_m(A^H, e_n)) \qquad (29)$$

for some constant $C$. *All* shifts (i.e., the roots of $p_i$) need to converge simultaneously to eigenvalues of $A_{22}$ in order to attain superlinear converge. In contrast, the bound (22) predicts superlinear convergence even if only one shift converges to a simple eigenvalue of $A_{22}$.

**Corollary 5** *Under the notation of Theorem 3, let $p_1, p_2, \ldots$ be a sequence of monomials of bounded degree such that $\hat{p}_i(A)$ is invertible for all $\hat{p}_i = p_1 p_2 \cdots p_i$. Assume that one root of $p_j$ converges to a simple eigenvalue $\lambda_2$ as $j \to \infty$, while the other roots and $p_j(A_{11})^{-1}$ remain bounded. Then for every $\rho > 0$ there is a constant $C$ so that*

$$d\big(y, \hat{p}_i(A)^{-H}\mathcal{K}_m(A^H, e_n)\big) \leq C\rho^i,$$

*provided $\mathcal{U} \cap \mathcal{K}_m(A^H, e_n) = \{0\}$, where $\mathcal{U}$ denotes the right invariant subspace of $A$ belonging to $\lambda(A_{11})$.*

*Proof.* Let $d$ be the maximal degree of $p_j$ and let $C_2 = \max\{1, \widetilde{C}_2 + |\lambda_2|\}$, where $\widetilde{C}_2$ is a uniform upper bound for the magnitude of all roots of $p_j$ for $j = 1, 2, \ldots$. Then $|p_j(\lambda_2)| \leq C_2^{d-1}|\lambda_2 - \sigma^{(j)}|$, where $\sigma^{(j)}$ denotes the root that converges to $\lambda_2$. Moreover, there is a constant $C_1$ such that $\|p_j(A_{11})^{-1}\| \leq C_1$ for all $j$. Choosing $k$ sufficiently large guarantees $C_1 C_2^{d-1}|\lambda_2 - \sigma^{(j)}| \leq \rho$ for all $j \geq k$. Hence, there is a constant $\widetilde{C}$ such that

$$\|\hat{p}_i(A_{11})^{-1}\|\,|\hat{p}_i(\lambda_2)| \leq \big(C_1 C_2^{d-1}\big)^i \prod_{j=1}^{i} |\lambda_2 - \sigma^{(j)}| \leq \widetilde{C}\rho^i.$$

Combined with Theorem 3 and Remark 4, this concludes the proof. □

The following theorem incorporates the Krylov subspace acceleration benefited from choosing the deflation window size $w$ larger than $m$.

**Theorem 6** *Under the notation of Theorem 3, assume that $C_T\, d(\mathcal{Y}, \mathcal{K}_m(A^{\mathsf{H}}, e_n)) < 1$. Then for any function $f$ that is analytic and nonzero on the spectrum of $A$,*

$$d\big(y, f(A)^{-\mathsf{H}}\mathcal{K}_w(A^{\mathsf{H}}, e_n)\big) \leq \frac{C_T \|f(A_{11})^{-1}\| \, |f(\lambda_2)| d(y, \mathcal{K}_m(A^{\mathsf{H}}, e_n))}{1 - C_T\, d(\mathcal{Y}, \mathcal{K}_m(A^{\mathsf{H}}, e_n))} \inf_{\phi \in \mathcal{P}_{w-m}} \frac{\|\phi(A_{11})\|}{|\phi(\lambda_2)|}, \quad (30)$$

*where $\mathcal{P}_{w-m}$ denotes the set of all polynomials of degree at most $w - m$.*

*Proof.* As in the proof of Theorem 3, we first block diagonalize $A$ and obtain bounds for the accordingly transformed subspaces. Let $T^{-1}AT = D$ with $D$ as in (23) and set $\tilde{y} = T^{\mathsf{H}}y$, $\widetilde{\mathcal{Y}} = T^{\mathsf{H}}\mathcal{Y}$. Without loss of generality, we may again assume $\tilde{y} = e_n$. Moreover, we have

$$T^{\mathsf{H}}f(A)^{-\mathsf{H}}\mathcal{K}_w(A^{\mathsf{H}}, e_n) = \mathcal{K}_w(D^{\mathsf{H}}, t)$$

with $t = f(D)^{-\mathsf{H}}T^{\mathsf{H}}e_n$. It follows that

$$
\begin{aligned}
d(e_n, \mathcal{K}_w(D^{\mathsf{H}}, t)) &= \inf_{v \in \mathcal{K}_w(D^{\mathsf{H}}, t)} \|e_n - v\| = \inf_{p \in \mathcal{P}_{w-1}} \|e_n - p(D^{\mathsf{H}})t\| \\
&= \inf_{\phi \in \mathcal{P}_{w-m}} \inf_{q \in \mathcal{P}_{m-1}} \|e_n - \phi(D^{\mathsf{H}})q(D^{\mathsf{H}})t\| \quad (31) \\
&= \inf_{\phi \in \mathcal{P}_{w-m}} \inf_{q \in \mathcal{P}_{m-1}} \frac{1}{|\phi(\bar{\lambda}_2)|} \|\phi(D^{\mathsf{H}})(e_n - q(D^{\mathsf{H}})t)\| \\
&= \inf_{\phi \in \mathcal{P}_{w-m}} \inf_{v \in \mathcal{K}_m(D^{\mathsf{H}}, t)} \frac{1}{|\phi(\bar{\lambda}_2)|} \|\phi(D^{\mathsf{H}})(e_n - v)\|. \quad (32)
\end{aligned}
$$

The implicit assumption that no root of $\phi$ coincides with $\bar{\lambda}_2$ can be justified by the following simple argument. If $\phi(\bar{\lambda}_2)$ is zero then $\|e_n - \phi(D^{\mathsf{H}})q(D^{\mathsf{H}})t\| \geq 1$. But since already the choice $\phi \equiv q \equiv 0$ gives $\|e_n - \phi(D^{\mathsf{H}})q(D^{\mathsf{H}})t\| = 1$, the infimum in (31) is not increased if $\phi(\bar{\lambda}_2) \neq 0$ is imposed.

The vector $\hat{v} = -\overline{f(\lambda_2)}f(D)^{-\mathsf{H}}v_1$ with $v_1$ defined as in (26) satisfies $\hat{v} \in \mathcal{K}_m(D^{\mathsf{H}}, t)$ and

$$e_n - \hat{v} = \left[ \begin{array}{c} \overline{f(\lambda_2)}f(A_{11})^{-\mathsf{H}}Fe_m \\ 0 \end{array} \right].$$

Choosing $v = \hat{v}$ in (32) and using (28) yields

$$d(e_n, \mathcal{K}_w(D^{\mathsf{H}}, t)) \leq \frac{\|f(A_{11})^{-\mathsf{H}}\| |f(\lambda_2)|}{1 - d(\widetilde{\mathcal{Y}}, \mathcal{K}_m(D^{\mathsf{H}}, t))} \, d(e_n, \widetilde{\mathcal{V}}) \inf_{\phi \in \mathcal{P}_{w-m}} \frac{\|\phi(A_{11})\|}{|\phi(\lambda_2)|}.$$

Together with Lemma 9 from Appendix A, this completes the proof. $\square$

Remark 4 applies analogously to Theorem 6. Comparing (30) with (22), the bound gains the factor $\inf_{\phi \in \mathcal{P}_{w-m}} \|\phi(A_{11})\|/|\phi(\lambda_2)|$. Consider a compact set $\Omega_1 \subset \mathbb{C}$ containing the eigenvalues of $A_{11}$ and define $\kappa(\Omega_1)$ to be the smallest constant for which

$$\|g(A_{11})\| \leq \kappa(\Omega_1) \max_{z \in \Omega_1} |g(z)| \quad (33)$$

holds uniformly for every analytic function $g$ on $\Omega_1$. Then, trivially,

$$\inf_{\phi \in \mathcal{P}_{w-m}} \frac{\|\phi(A_{11})\|}{|\phi(\lambda_2)|} \leq \kappa(\Omega_1) \inf_{\phi \in \mathcal{P}_{w-m}} \frac{\max\{|\phi(z)| : z \in \Omega_1\}}{|\phi(\lambda_2)|}.$$

Such approximation problems play a prominent role in the convergence analysis of Krylov subspace methods for nonnormal matrices. It is clear that $\inf_{\phi \in \mathcal{P}_{w-m}} \max\{|\phi(z)| \ : \ z \in \Omega_1\}/|\phi(\lambda_2)|$ cannot be larger than 1 (choose $\phi \equiv 1$) and decays as $w - m$ becomes larger. Estimates for this decay can be found using potential theory, see [4, 12], but this is beyond the scope of this paper.

In LAPACK 3.1 [1], the default values for the number of shifts and the size of the deflation window for $3000 \leq n < 6000$ are $m = 128$ and $w = 192$, respectively. As explained in more detail in [9], this choice was based on computational experiments with pseudo-random matrices and matrices from [2]. A good choice of these parameters is crucial for the performance of the QR algorithm. For example, choosing a larger $w$ increases the cost of the deflation procedure (which requires $O(w^2 n)$ operations) but also results in earlier deflations and therefore in fewer multishift QR iterations (each requires $O(mn^2)$ operations). While Theorem 6 does not provide any new insight into the choice of $m$, it does shed some light on a beneficial choice of $w$. Unfortunately, the exact computation of the quantity $\inf_{\phi \in \mathcal{P}_{w-m}} \frac{\|\phi(A_{11})\|}{|\phi(\lambda_2)|}$ in (30) is infeasible, simply because $A_{11}$ and $\lambda_2$ are available only after the Schur form of $A$ has been computed. The results in this paper should thus be understood only as a first step towards developing a strategy that chooses $w$ nearly optimal in each iteration. For this purpose, cheap estimates on the decay of $\inf_{\phi \in \mathcal{P}_{w-m}} \frac{\|\phi(A_{11})\|}{|\phi(\lambda_2)|}$ as $w$ increases need to be developed, possibly using heuristics on the distribution of the eigenvalues of $A_{11}$.

## 5.2  Convergence of invariant subspaces

Not only the convergence of individual eigenvectors but also the convergence of the left invariant subspace $\mathcal{Y}$ as a whole is improved using Krylov subspaces of larger dimension. To quantify this effect, we can combine the bound (29) with results by Beattie, Embree, and Rossi [4] on the convergence of invariant subspaces in Krylov subspaces.

**Theorem 7** *Let $A \in \mathbb{C}^{n \times n}$ have a block Schur decomposition*

$$Q^{\mathsf{H}} A Q = \left[ \begin{array}{cc} A_{11} & A_{12} \\ 0 & A_{22} \end{array} \right], \quad A_{11} \in \mathbb{C}^{(n-m) \times (n-m)}, \quad A_{22} \in \mathbb{C}^{m \times m},$$

*for some unitary matrix $Q \in \mathbb{C}^{n \times n}$, such that $\lambda(A_{11}) \cap \lambda(A_{22}) = \emptyset$. Let $\mathcal{Y}$ denote the left invariant subspace of $A$ belonging to $\lambda(A_{22})$. Assume that $C_T d\big(\mathcal{Y}, \mathcal{K}_m(A^{\mathsf{H}}, e_n)\big) < 1$, where $C_T = \|P_{A_{11}}\| + \sqrt{1 + \|P_{A_{11}}\|^2}$ with $P_{A_{11}}$ being the spectral projector belonging to $A_{11}$.*
*Then for any function $f$ that is analytic and nonzero on the spectrum of $A$,*

$$d\big(\mathcal{Y}, f(A)^{-\mathsf{H}} \mathcal{K}_w(A^{\mathsf{H}}, e_n)\big) \leq C \, \|f(A_{11})^{-1}\| \|f(A_{22})\| \inf_{\phi \in \mathcal{P}_{w-m}} \|\phi(A_{11})\| \|\phi(A_{22})^{-1}\|, \qquad (34)$$

*where*

$$C = \frac{C_T^2 \, d(\mathcal{Y}, \mathcal{K}_m(A^{\mathsf{H}}, e_n))}{\sqrt{1 - C_T^2 \, d(\mathcal{Y}, \mathcal{K}_m(A^{\mathsf{H}}, e_n))^2}}.$$

*Proof.* There is an invertible lower block triangular matrix $T_1$ with $\kappa(T_1) = \|P_{A_{11}}\| + \sqrt{1 + \|P_{A_{11}}\|^2}$ such that

$$(QT_1)^{-1} A (QT_1) = \left[ \begin{array}{cc} A_{11} & 0 \\ 0 & A_{22} \end{array} \right] =: D.$$

Let $T = QT$ and $\widetilde{\mathcal{Y}} = T^{\mathsf{H}}\mathcal{Y}$. Then

$$T^{\mathsf{H}} f(A)^{-\mathsf{H}} \mathcal{K}_w(A^{\mathsf{H}}, e_n) = \mathcal{K}_w(D^{\mathsf{H}}, f(D)^{-\mathsf{H}} u)$$

with $u = T^{\mathsf{H}} e_n$. Partition $u = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}$ with $u_2 \in \mathbb{C}^m$. Then Theorem 3.4 in [4] states

$$d\big(\widetilde{\mathcal{Y}}, \mathcal{K}_w(D^{\mathsf{H}}, f(D)^{-\mathsf{H}} u)\big) \leq C_1 \inf_{\phi \in \mathcal{P}_{w-m}} \|\phi(A_{11})\| \|\phi(A_{22})^{-1}\|,$$

where

$$C_1 = \max_{\psi \in \mathcal{P}_{m-1}} \frac{\|\psi(A_{11})^{\mathsf{H}} f(A_{11})^{-\mathsf{H}} u_1\|}{\|\psi(A_{22})^{\mathsf{H}} f(A_{22})^{-\mathsf{H}} u_2\|}.$$

Using the fact that $\psi(D)^{\mathsf{H}}$ and $f(D)^{-\mathsf{H}}$ commute,

$$
\begin{aligned}
C_1 &\leq \|f(A_{11})^{-1}\| \|f(A_{22})\| \max_{\psi \in \mathcal{P}_{m-1}} \frac{\|\psi(A_{11})^{\mathsf{H}} u_1\|}{\|\psi(A_{22})^{\mathsf{H}} u_2\|} \\
&= \|f(A_{11})^{-1}\| \|f(A_{22})\| \max_{\left[\begin{smallmatrix} v_1 \\ v_2 \end{smallmatrix}\right] \in \mathcal{K}_m(D^{\mathsf{H}}, u)} \frac{\|v_1\|}{\|v_2\|} \\
&\leq \|f(A_{11})^{-1}\| \|f(A_{22})\| \frac{d\big(\widetilde{\mathcal{Y}}, \mathcal{K}_m(D^{\mathsf{H}}, u)\big)}{\sqrt{1 - d\big(\widetilde{\mathcal{Y}}, \mathcal{K}_m(D^{\mathsf{H}}, u)\big)^2}},
\end{aligned}
$$

where the latter inequality follows as in the proof of Lemma 4.4 in [30]. The proof is concluded by applying Lemma 4.2 in [30] twice, relating the distances between the transformed subspaces to the distances of the original subspaces. □

Once again, an analogue of Remark 4 applies, showing that $C_T \, d\big(\mathcal{Y}, \mathcal{K}_m(A^{\mathsf{H}}, e_n)\big) < 1$ can be replaced by the weaker assumption $\mathcal{U} \cap \mathcal{K}_m(A^{\mathsf{H}}, e_n) = \{0\}$. Comparing (34) with the classical bound (29), the extra factor $\inf_{\phi \in \mathcal{P}_{w-m}} \|\phi(A_{11})\| \|\phi(A_{22})^{-1}\|$ is gained. If $\Omega_1, \Omega_2 \subset \mathbb{C}$ are compact sets containing $\lambda(A_{11}), \lambda(A_{22})$, respectively, then

$$\inf_{\phi \in \mathcal{P}_{w-m}} \|\phi(A_{11})\| \|\phi(A_{22})^{-1}\| \leq \kappa(\Omega_1)\kappa(\Omega_2) \min_{\phi \in \mathcal{P}_{w-m}} \frac{\max\{|\phi(z)| : z \in \Omega_1\}}{\min\{|\phi(z)| : z \in \Omega_2\}},$$

where $\kappa(\Omega_1)\kappa(\Omega_2)$ are defined as in (33). Note that the factor $\kappa(\Omega_2)$ can actually be removed, see [5, Theorem 3.3], at the expense of having a different polynomial approximation problem. In any case, the bound can be expected to decay as $w - m$ increases, see [4, Sec. 4] and [5, Sec. 3] for estimates of this decay. When using Francis shifts, this means that besides the quadratically vanishing $\hat{p}_i(A_{22})$ we have an additional factor that may become very small for larger $w$, resulting in nearly superquadratic convergence.

Finally, let us remark that there is always a polynomial $q$, depending on $A$ and $f$, such that $q(A) = f^{-1}(A)$. This implies an equivalence between the QR algorithm and the Arnoldi method with polynomial restarts, as discussed by Lehoucq [21]. In principle, this connection could be used to apply the bounds in [4, 5] verbatim to the QR algorithm. In contrast, Theorem 7 treats the convergence obtained from the shifts and from the Krylov subspace separately, having the advantage that it allows more insight into the benefits gained from choosing a larger deflation window.

## 6    Conclusions

This paper contributes to the understanding of why aggressive early deflation works so well in practice. A very intuitive explanation can be drawn from practical experiences with Krylov subspace methods for computing eigenvalues. Extracting Ritz pairs (= aggressive early deflation) is a much more effective strategy for detecting converged eigenvalues than only testing the subdiagonal entries of the Hessenberg factor in an Arnoldi decomposition (= classical deflation). The convergence bounds from Section 5 provide a mathematical explanation, stating individual bounds for each converging eigenvalue and showing that the bounds are multiplied by a factor that approaches zero as the deflation window size $w$ increases.

This paper should be seen as a first step towards developing a strategy for choosing $w$ optimally in each QR iteration. In a serial computing environment, the current default value for $w$ implemented in LAPACK already delivers good performance across a wide range of examples and more sophisticated strategies might not lead to significant speedup. However, in a parallel computing environment, where aggressive early deflation will constitute a bottleneck, we expect the performance to become more sensitive with respect to the choice of $w$.

Finally, it is tempting to ask whether other Ritz pair extraction techniques, such as refined Ritz vectors, could be used to enhance the convergence of the QR algorithm even further. Although preliminary numerical experiments have indicated no obvious beneficial effect on the average performance of the QR algorithm, the use of these techniques in avoiding exceptional global convergence failures remains to be studied.

## 7    Acknowledgments

## A    Appendix

This section collects two elementary facts needed in the proofs of Section 5.

**Lemma 8** *Let* $V = \begin{bmatrix} F \\ I_m \end{bmatrix} \in \mathbb{C}^{n \times m}$. *Then*

$$\|e_n - V e_m\| \leq \frac{\sqrt{1 + \|F\|^2}}{\sqrt{1 + \|F\|^2} - \|F\|} \, \|e_n - V(V^{\mathsf{H}}V)^{-1} e_m\|.$$

*Proof.* Set $r_0 = e_n - V(V^{\mathsf{H}}V)^{-1} e_m$ and $r_1 = e_n - V e_m$. Then

$$
\begin{aligned}
r_0 &= \begin{bmatrix} -F(I + F^{\mathsf{H}}F)^{-1} \\ I - (I + F^{\mathsf{H}}F)^{-1} \end{bmatrix} e_m = \begin{bmatrix} -(I + FF^{\mathsf{H}})^{-1}F \\ (I + F^{\mathsf{H}}F)^{-1}F^{\mathsf{H}}F \end{bmatrix} e_m \\
&= \begin{bmatrix} (I + FF^{\mathsf{H}})^{-1} & F(I + F^{\mathsf{H}}F)^{-1} \\ -(I + F^{\mathsf{H}}F)^{-1}F^{\mathsf{H}} & (I + F^{\mathsf{H}}F)^{-1} \end{bmatrix} r_1 = (I + E)r_1,
\end{aligned}
$$

where

$$E = \begin{bmatrix} -F(I + F^{\mathsf{H}}F)^{-1}F^{\mathsf{H}} & F(I + F^{\mathsf{H}}F)^{-1} \\ -(I + F^{\mathsf{H}}F)^{-1}F^{\mathsf{H}} & -F^{\mathsf{H}}(I + FF^{\mathsf{H}})^{-1}F \end{bmatrix}.$$

A singular value decomposition of $F$ verifies $\|E\| = \|F\|/\sqrt{1 + \|F\|^2}$. Thus $\|E\| < 1$ and $(I + E)$ is invertible. Finally,

$$\|r_1\| \leq \|(I + E)^{-1}\|\|r_0\| \leq \frac{1}{1 - \|E\|}\|r_0\| = \frac{\sqrt{1 + \|F\|^2}}{\sqrt{1 + \|F\|^2} - \|F\|}\|r_0\|$$

concludes the proof. □

**Lemma 9** *Let $\mathcal{T}, \mathcal{U}$ be subspaces and let $P$ be an invertible matrix. Then for any $s \in \mathbb{C}^n$, $d(s, \mathcal{T}) \leq \alpha \, d(s, \mathcal{U})$ implies $d(Ps, P\mathcal{T}) \leq \alpha \|P\|\|P^{-1}\| \, d(Ps, P\mathcal{U})$.*

*Proof.* Without loss of generality, we may assume $\|s\| = 1$. Then

$$\begin{aligned}
d(Ps, P\mathcal{T}) &= \inf_{t \in \mathcal{T}} \frac{\|Ps - Pt\|}{\|Ps\|} \leq \frac{\|P\|}{\|Ps\|} \inf_{t \in \mathcal{T}} \|s - t\| \\
&\leq \alpha \frac{\|P\|}{\|Ps\|} \inf_{u \in \mathcal{U}} \|s - u\| \leq \alpha \|P\|\|P^{-1}\| \inf_{u \in \mathcal{U}} \frac{\|Ps - Pu\|}{\|Ps\|} \\
&= \alpha \|P\|\|P^{-1}\| \, d(Ps, P\mathcal{U}).
\end{aligned}$$

□

# References

[1] E. Anderson, Z. Bai, C. H. Bischof, S. Blackford, J. W. Demmel, J. J. Dongarra, J. Du Croz, A. Greenbaum, S. Hammarling, A. McKenney, and D. C. Sorensen. *LAPACK Users' Guide.* SIAM, Philadelphia, PA, third edition, 1999.

[2] Z. Bai, D. Day, J. W. Demmel, and J. J. Dongarra. A test matrix collection for non-Hermitian eigenvalue problems (release 1.0). Technical Report CS-97-355, Department of Computer Science, University of Tennessee, Knoxville, TN, USA, March 1997. Also available online from `http://math.nist.gov/MatrixMarket`.

[3] F. L. Bauer. Das Verfahren der Treppeniteration und verwandte Verfahren zur Lösung algebraischer Eigenwertprobleme. *Z. Angew. Math. Phys.*, 8:214–235, 1957.

[4] C. A. Beattie, M. Embree, and J. Rossi. Convergence of restarted Krylov subspaces to invariant subspaces. *SIAM J. Matrix Anal. Appl.*, 25(4):1074–1109, 2004.

[5] C. A. Beattie, M. Embree, and D. C. Sorensen. Convergence of polynomial restart Krylov methods for eigenvalue computations. *SIAM Rev.*, 47(3):492–515, 2005.

[6] K. Braman, R. Byers, and R. Mathias. The multishift $QR$ algorithm. I. Maintaining well-focused shifts and level 3 performance. *SIAM J. Matrix Anal. Appl.*, 23(4):929–947, 2002.

[7] K. Braman, R. Byers, and R. Mathias. The multishift $QR$ algorithm. II. Aggressive early deflation. *SIAM J. Matrix Anal. Appl.*, 23(4):948–973, 2002.

[8] H. J. Buurema. *A geometric proof of convergence for the QR method.* Rijksuniversiteit te Groningen, Groningen, 1970. Doctoral dissertation, University of Groningen.

[9] R. Byers. LAPACK 3.1 xHSEQR: Tuning and implementation notes on the small bulge multi-shift QR algorithm with aggressive early deflation. LAPACK Working Note 187, 2007.

[10] G. Chen and Z. Jia. A reverse order implicit $Q$-theorem and the Arnoldi process. *J. Comput. Math.*, 20(5):519–524, 2002.

[11] J. W. Demmel. Computing stable eigendecompositions of matrices. *Linear Algebra Appl.*, 79:163–193, 1986.

[12] T. A. Driscoll, K.-C. Toh, and L. N. Trefethen. From potential theory to matrix iterations in six steps. *SIAM Rev.*, 40(3):547–578, 1998.

[13] J. G. F. Francis. The QR transformation, parts I and II. *Computer Journal*, 4:265–271, 332–345, 1961, 1962.

[14] G. H. Golub and C. F. Van Loan. *Matrix Computations.* Johns Hopkins University Press, Baltimore, MD, third edition, 1996.

[15] Z. Jia. Refined iterative algorithms based on Arnoldi's process for large unsymmetric eigenproblems. *Linear Algebra Appl.*, 259:1–23, 1997.

[16] Z. Jia and G. W. Stewart. An analysis of the Rayleigh-Ritz method for approximating eigenspaces. *Math. Comput.*, 70(234):637–647, 2001.

[17] D. Kressner. *Numerical Methods and Software for General and Structured Eigenvalue Problems.* PhD thesis, TU Berlin, Institut für Mathematik, Berlin, Germany, 2004.

[18] D. Kressner. Deflation in Krylov subspace methods and distance to uncontrollability. In *Proceedings of the Conference on Applied Mathematics and Scientific Computing*, 2005. To appear in Annali dell'Universita' di Ferrara.

[19] V. N. Kublanovskaya. On some algorithms for the solution of the complete eigenvalue problem. *Zhurnal Vychislitelnoi Matematiki i Matematicheskoi Fiziki*, 1:555–570, 1961.

[20] B. Lang. Effiziente Orthogonaltransformationen bei der Eigen- und Singulärwertzerlegung. Habilitationsschrift, 1997.

[21] R. B. Lehoucq. Implicitly restarted Arnoldi methods and subspace iteration. *SIAM J. Matrix Anal. Appl.*, 23(2):551–562, 2001.

[22] R. B. Lehoucq, D. C. Sorensen, and C. Yang. *ARPACK users' guide.* SIAM, Philadelphia, PA, 1998. Solution of large-scale eigenvalue problems with implicitly restarted Arnoldi methods.

[23] B. N. Parlett and W. G. Poole, Jr. A geometric theory for the QR, LU and power iterations. *SIAM J. Numer. Anal.*, 10:389–412, 1973.

[24] Y. Saad. Variations on Arnoldi's method for computing eigenelements of large unsymmetric matrices. *Linear Algebra Appl.*, 34:269–295, 1980.

[25] G. W. Stewart. *Matrix Algorithms. Vol. II*. SIAM, Philadelphia, PA, 2001. Eigensystems.

[26] G. W. Stewart. A Krylov-Schur algorithm for large eigenproblems. *SIAM J. Matrix Anal. Appl.*, 23(3):601–614, 2001/02.

[27] G. W. Stewart and J.-G. Sun. *Matrix Perturbation Theory*. Academic Press, New York, 1990.

[28] D. S. Watkins. Understanding the QR algorithm. *SIAM Rev.*, 24(4):427–440, 1982.

[29] D. S. Watkins. Understanding the QR algorithm, part II, 2006. Submitted for publication.

[30] D. S. Watkins and L. Elsner. Convergence of algorithms of decomposition type for the eigenvalue problem. *Linear Algebra Appl.*, 143:19–47, 1991.

[31] J. H. Wilkinson. *The Algebraic Eigenvalue Problem*. Clarendon Press, Oxford, 1965.

[32] J.-P. M. Zemke. Hessenberg eigenvalue-eigenmatrix relations. *Linear Algebra Appl.*, 414(2-3):589–606, 2006.