

# Inexact methods for the low rank solution to large scale Lyapunov equations

Patrick Kürschner\*      Melina A. Freitag †

September 20, 2018

## Abstract

The rational Krylov subspace method (RKSM) and the low-rank alternating directions implicit (LR-ADI) iteration are established numerical tools for computing low-rank solution factors of large-scale Lyapunov equations. In order to generate the basis vectors for the RKSM, or extend the low-rank factors within the LR-ADI method the repeated solution to a shifted linear system is necessary. For very large systems this solve is usually implemented using iterative methods, leading to inexact solves within this inner iteration. We derive theory for a relaxation strategy within these inexact solves, both for the RKSM and the LR-ADI method. Practical choices for relaxing the solution tolerance within the inner linear system are then provided. The theory is supported by several numerical examples.

## 1 Introduction

We consider the numerical solution of large scale Lyapunov equations of the form

$$AX + XA^T = -BB^T, \quad (1)$$

where  $A \in \mathbb{R}^{n \times n}$ ,  $B \in \mathbb{R}^{n \times r}$ . Lyapunov equations play a fundamental role in many areas of applications, such as control theory, model reduction and signal processing, see, e.g. [2, 10]. Here we assume that  $\Lambda(A) \in \mathbb{C}_-$  and that the right hand side is of low rank, e.g.  $r \ll n$ , which is often the case in practise. A large and growing amount of literature considers the solution to (1), see [36] and references therein for an overview of developments and methods.

If  $A$  is large and sparse, the solution matrix  $X$  is dense, and possibly cannot be stored. For a low-rank right hand side, however, the solution  $X$  often has a very small numerical rank [28, 32] and many algorithms have been developed that approximate  $X$  by a low-rank matrix  $X \approx ZZ^T$ , where  $Z \in \mathbb{R}^{n \times s}$ ,  $s \ll n$ . Several methods belong to such low-rank

---

\*Max Planck Institute for Dynamics of Complex Technical Systems, Computational Methods in Systems and Control Theory, Magdeburg, Germany, [kuerschner@mpi-magdeburg.mpg.de](mailto:kuerschner@mpi-magdeburg.mpg.de)

†Department of Mathematical Sciences, University of Bath, Claverton Down, BA2 7AY, United Kingdom, [m.a.freitag@bath.ac.uk](mailto:m.a.freitag@bath.ac.uk)

algorithms, namely projection type methods based on Krylov subspaces [21, 34, 16, 15, 35] and low-rank alternating directions implicit (ADI) methods [28, 25, 10, 9, 22].

This paper considers both the rational Krylov subspace method (from the family of projection type methods) and the low-rank ADI method. One of the computationally most expensive parts in both methods is that, in each iteration, shifted linear systems of the form

$$(A - \sigma I)y = z, \quad z \in \mathbb{R}^n,$$

need to be solved, where the shift  $\sigma$  is usually variable and both the shift  $\sigma$  and the right hand side  $z$  depend on the particular algorithm used. Normally these linear systems are solved by sparse-direct or iterative methods. When iterative methods, such as preconditioned Krylov subspace methods, are used to solve the linear systems, then these solves are implemented inexactly and we obtain a so-called inner-outer iterative method. The outer method is (in our case) a rational Krylov method or a low-rank ADI iteration. The inner problem is the iterative solution to the linear systems. The inner solves are often carried out at least as accurately as the required solution accuracy for the Lyapunov equation (cf., e.g., the numerical experiments in [15]), usually regarding the associated Lyapunov residual norms. It turns out that this is not necessary, and the underlying theory is the main contribution of this paper.

Inexact methods have been considered for the solution to linear systems and eigenvalue problems (see [37, 42, 33, 18] and references therein). One focus has been on inexact inverse iteration and similar methods, where, in general, the accuracy of the inexact solve has to be increased to obtain convergence [17]. For subspace expansion methods, such as the Krylov methods considered in [37, 42, 33], it has been observed numerically and shown theoretically that it is possible to relax the solve tolerance as the outer iteration proceeds.

In this paper we prove that for the linear solve both for the rational Krylov subspace method and the low-rank ADI algorithm we can relax the solve tolerance within the inner solve, a similar behavior as observed in [33, 18] for Krylov methods applied to eigenvalue methods and in [37, 42] for inexact matrix-vector products within linear system solves. To this end, we provide practical relaxation strategies for both methods and give numerical examples.

The paper is organised as follows. In Section 2 we review results about rational Krylov subspace methods. Those are used to show important properties about Galerkin projections. A new inexact rational Arnoldi decomposition is derived, extending the theory of [29]. We show in Theorem 2, Corollary 3 and Corollary 5 that the entries of the solution to the projected Lyapunov equation have a decreasing pattern. This crucial novel result is then used to demonstrate that, in the inexact rational Krylov subspace method, the linear system solve can be relaxed, in a way inversely proportional to this pattern. Section 3 is devoted to low-rank ADI methods. We show that the low-rank factors arising within the inexact ADI iteration satisfy an Arnoldi like decomposition in Theorem 8. This theory is again new and significant for deriving a straightforward relaxation strategy for inexact low-rank ADI methods. Finally, in Section 4 we test several practical relaxation strategies and provide numerical evidence for our findings. Our examples show that, in particular for very large problems, we can save up to half the number of inner iterations within the both methods for solving Lyapunov equations.

Notation: We use  $\|\cdot\|$  to denote the 2-norm of a matrix or a vector, and  $(\cdot)^*$  for the complex conjugate transpose of a matrix or a vector.

## 2 Rational Krylov subspace methods and inexact solves

### 2.1 Introduction and rational Arnoldi decompositions

Projection methods for (1) follow Galerkin principles similar to, e.g., the Arnoldi method for eigenvalue computations or linear systems (in this case called full orthogonalization method). Let  $\mathcal{Q} = \text{span}\{Q\} \subset \mathbb{C}^n$  be a subspace of  $\mathbb{C}^n$  with an orthogonal basis matrix  $Q \in \mathbb{C}^{n \times k}$ ,  $Q^*Q = I_k$ ,  $k \ll n$ . We look for low-rank approximate solutions to (1) in the form  $Q\tilde{X}Q^*$  with  $\tilde{X} = \tilde{X}^* \in \mathbb{R}^{k \times k}$ , i.e.

$$X \in \mathcal{Z} := \{Q\tilde{X}Q^* \in \mathbb{C}^{n \times n}, \tilde{X} = \tilde{X}^*, \text{span}\{Q\} = \mathcal{Q}\}.$$

Imposing a Galerkin condition [21, 30] onto the Lyapunov residual of solutions of such form  $\mathcal{R}(Q\tilde{X}Q^*) = A(Q\tilde{X}Q^*) + (Q\tilde{X}Q^*)A^* + BB^* \perp \mathcal{Z}$  leads to the projected Lyapunov equation

$$Q^*AQ\tilde{X} + \tilde{X}Q^*A^*Q + Q^*BB^*Q = 0,$$

i.e.,  $\tilde{X}$  is the solution of a small,  $k$ -dimensional Lyapunov equation which can be solved by algorithms employing dense numerical linear algebra, e.g., the Bartels-Stewart method [3]. Note that  $A + A^* < 0$  has to hold in order to ensure  $\Lambda(Q^*AQ) \subset \mathbb{C}_-$ , but in practice this approach often also works when this condition is violated.

Usually, one produces sequences of subspaces of increasing dimensions in an iterative manner, e.g.  $\mathcal{Q}_1 \subseteq \mathcal{Q}_2 \subseteq \dots \subseteq \mathcal{Q}_k$ . For practical problems, using the standard (block) Krylov subspace

$$\mathcal{Q}_k = \mathcal{K}_k(A, B) = \text{span}\{B, AB, \dots, A^{k-1}B\} \quad (2)$$

is not sufficient and leads to a slowly converging process and, hence, large low-rank solution factors. A better performance can in most cases be achieved by using rational Krylov subspaces which we use here in the form

$$\mathcal{Q}_k = \mathcal{RK}_k(A, B, \boldsymbol{\xi}) := \text{span}\left\{B, (A - \xi_2 I)^{-1}B, \dots, \prod_{i=2}^k (A - \xi_i I)^{-1}B\right\}, \quad (3)$$

with shifts (poles)  $\xi_i \in \mathbb{C}_+$  (hence  $\xi_i \notin \Lambda(A)$ ),  $i = 2, \dots, k$ . An orthonormal basis for  $\mathcal{RK}_k$  can be computed by the (block) rational Arnoldi process [29]. The resulting rational Krylov subspace method (RKSM) for computing approximate solutions to (1) is given in Algorithm 1. The shift parameters are crucial for a rapid reduction of the Lyapunov residual norm  $\|\mathcal{R}_j\| := \|\mathcal{R}(Q_j Y_j Q_j^*)\|$  and can be generated a-priori or adaptively in the course of the iteration [16]. The extension of the orthonormal basis  $\text{range}(Q_j)$  by  $w$  in Line 6 should be done by a robust orthogonalization process, e.g., using an repeated (block) Gram-Schmidt process. The orthogonalization coefficients are accumulated in the block Hessenberg matrix  $H_j = [h_{i,k}] \in \mathbb{R}^{jr \times jr}$ ,  $h_{i,k} \in \mathbb{R}^{r \times r}$ ,  $i, k = 1, \dots, j$ ,  $i < k + 2$ . If the

---

**Algorithm 1:** Rational Krylov subspace method for solving (1)

---

**Input** :  $A, B$  as in (1), shifts  $\{\xi_1, \dots, \xi_{j_{\max}}\} \subset \mathbb{C}_+$ , tolerance  $0 < \tau \ll 1$ .  
**Output:**  $Q \in \mathbb{C}^{n \times jr}$ ,  $\tilde{X} \in \mathbb{C}^{jr \times jr}$  such that  $Q\tilde{X}Q^* \approx X$  with  $jr \ll n$ .

- 1 Compute thin QR decomposition of  $B$ :  $Q_1\beta = B$ ,  $j = 1$ .
- 2 **while**  $\|\mathcal{R}(Q_j Y_j Q_j^*)\| > \tau \|BB^*\|$  **do**
- 3     Solve  $(A - \xi_{j+1}I)w_{j+1} = q_j$  for  $w_{j+1}$ .
- 4     Orthogonally extend basis matrix  $Q_j$ :
- 5      $h_{1:j,j} = Q_j^* w_{j+1}$ ,  $\hat{w}_{j+1} = w_{j+1} - Q_j h_{1:j,j}$ ,
- 6      $q_{j+1} h_{j+1,j} = \hat{w}_{j+1}$ ,  $Q_{j+1} = [Q_j, q_{j+1}]$ .
- 7     Projected equation:  $T_j = Q_j^*(A Q_j)$ ,  $\tilde{B}_j = Q_j^* B = [\beta^*, 0]^*$ .
- 8     Solve  $T_j Y_j + Y_j T_j^* + \tilde{B}_j \tilde{B}_j^* = 0$  for  $Y_j$ .
- 9     Estimate residual norm  $\mathcal{R}(Q_j Y_j Q_j^*)$ .
- 10     $j = j + 1$ .
- 11  $\tilde{X} = Y_j$ ,  $Q = Q_j$

---

new basis blocks are normalized using a thin QR factorization, then the  $h_{i+1,i}$  matrices are upper triangular.

In the following we summarise properties of the RKSM method, in particular known results about the rational Arnoldi decomposition and the Lyapunov residual, which will be crucial later on. We simplify this discussion of the rational Arnoldi process and Algorithm 1 to the case  $r = 1$ . From there, one can generalize to the situation  $r > 1$ .

Assuming at first that the linear systems in Line 3 of Algorithm 1 are solved *exactly*, then, after  $j$  steps of Algorithm 1, the generated quantities satisfy a rational Arnoldi decomposition [29, 20, 16]

$$A Q_j = Q_j T_j + g_j h_{j+1,j} e_j^* H_j^{-1}, \quad g_j := q_{j+1} \xi_{j+1} - (I - Q_j Q_j^*) A q_{j+1}, \quad (4)$$

from which it follows that the restriction of  $A$  onto  $\mathcal{RK}_k$  can be expressed as

$$T_j = Q_j^*(A Q_j) = (I + H_j D_j) H_j^{-1} - Q_j^* A q_{j+1} h_{j+1,j} e_j^* H_j^{-1}, \quad (5)$$

with  $D_j = \text{diag}(\xi_2, \dots, \xi_{j+1})$ .  $T_j$ ,  $H_j$  and  $D_j$  are square matrices of size  $j$ . The relation (5) can be used to compute the restriction  $T_j$  efficiently in terms of memory usage. The rational Arnoldi decomposition can also be expressed as

$$A Q_{j+1} H_{j+1,j} = Q_{j+1} M_{j+1,j}, \quad M_{j+1,j} := \begin{bmatrix} I + H_j D_j \\ \xi_{j+1} h_{j+1,j} e_j^* \end{bmatrix}, \quad (6)$$

see, e.g., [29, 12] and  $T_j$  can be associated with the Hessenberg-Hessenberg pair  $(M_{j,j}, H_{j,j})$ . Here we generally use  $H_j = H_{j,j}$  and  $M_j = M_{j,j}$  for simplicity.  $H_{j+1,j}$  and  $M_{j+1,j}$  are  $j + 1 \times j$  matrices.

The final part of this section discusses how the Lyapunov residual in Line 9 of Algorithm 1 can be computed efficiently. If the projected Lyapunov equation in Line 8

$$T_j Y_j + Y_j T_j^* + \tilde{B}_j \tilde{B}_j^* = 0 \quad (7)$$

is solved for  $Y_j$ , then for the Lyapunov residual  $\mathcal{R}_j$  after  $j$  steps of the RKSM, we have

$$\mathcal{R}_j = \mathcal{R}(Q_j Y_j Q_j^*) = F_j + F_j^*, \quad (8a)$$

with the rank-1 matrix

$$F_j := L_j Q_j^* \in \mathbb{C}^{n \times n}, \quad L_j := g_j h_{j+1,j} e_j^* H_j^{-1} Y_j \in \mathbb{C}^{n \times j}, \quad (8b)$$

as was shown in [16, 35]. The term  $L_j$  is sometimes referred to as ‘‘semi-residual’’.

Note that relation (8a) is common for projection methods for (1), but the special structure of  $L_j$  in (8b) arises from the rational Arnoldi process. Because  $g_j^* Q_j = 0$  we have that  $F_j^2 = 0$  and  $\|\mathcal{R}_j\| = \|F_j\| = \|L_j\|$ , which enables an efficient way for computing the norm of the Lyapunov residual for the approximate solution  $Q_j Y_j Q_j^*$  via

$$\|\mathcal{R}_j\| = \|g_j\| \|h_{j+1,j} e_j^* H_j^{-1} Y_j\|. \quad (9)$$

## 2.2 The inexact rational Arnoldi method

The solution of the linear systems for  $w_{j+1}$  in each step of RKSM (3 in Algorithm 1) is one of the computationally most expensive stages of the rational Arnoldi method. In this section we investigate inexact solves of this linear systems, e.g., by iterative Krylov subspace methods, but we assume that this is the only source of inaccuracy in the algorithm. Clearly, some of the above properties do not hold anymore if the linear systems are solved inaccurately. Let

$$s_j := q_j - (A - \xi_{j+1} I) \tilde{w}_{j+1}, \quad \|s_j\| \leq \tau_j^{\text{ls}}$$

be the residual vectors with respect to the linear systems and inexact solution vectors  $\tilde{w}_{j+1}$  with  $\tau_j^{\text{ls}} < 1$  being the solve tolerance of the linear system at step  $j$  of the rational Arnoldi method.

The derivations in [29],[16, Proof of Prop. 42.] can be modified with

$$q_j = (A - \xi_{j+1} I) Q_{j+1} h_{1:j+1,j} + s_j,$$

in order to obtain

$$A Q_{j+1} H_{j+1,j} = Q_{j+1} M_{j+1,j} - S_j, \quad S_j := [s_1, \dots, s_j], \quad M_{j+1,j} := \begin{bmatrix} I + H_j D_j \\ \xi_{j+1} h_{j+1,j} e_j^* \end{bmatrix},$$

instead of (6). Hence  $A Q_j H_j = [Q_j (I + H_j D_j) + (\xi_{j+1} q_{j+1} - A q_{j+1}) e_j^* h_{j+1,j}] - S_j$  leading to the perturbed rational Arnoldi relation

$$A Q_j = Q_j T_j^{\text{impl.}} + g_j e_j^* h_{j+1,j} H_j^{-1} - S_j H_j^{-1}, \quad (10)$$

where  $T_j^{\text{impl.}} := [(I + H_j D_j) - Q_j^* A q_{j+1}] e_j^* h_{j+1,j} H_j^{-1}$  marks the restriction of  $A$  by the implicit formula (5). However, the explicit restriction of  $A$  can be written as

$$T_j^{\text{expl.}} := Q_j^* (A Q_j) = T_j^{\text{impl.}} - Q_j^* S_j H_j^{-1}, \quad (11)$$

which highlights the problem that the implicit (computed) restriction  $T_j^{\text{impl.}}$  from (5) is not the true restriction of  $A$  onto  $\text{range}(Q_j)$ . In fact, the above derivations reveal that

$$T_j^{\text{impl.}} = Q_j^*(A + E_j)Q_j, \quad E_j := S_j H_j^{-1} Q_j^*,$$

i.e., the implicit restriction (5) is the exact restriction of a perturbation of  $A$  [19, 23].

Similar to [19] we prefer to use  $T_j^{\text{expl.}}$  for defining the projected problem as this keeps the whole process slightly closer to the original matrix  $A$ . Unlike  $T_j^{\text{impl.}}$ , the explicit restriction  $T_j^{\text{expl.}}$  is not connected to a Hessenberg-Hessenberg pair, since  $T_j^{\text{expl.}} H_j$  does by (11) not have upper Hessenberg structure. To compute  $T_j^{\text{expl.}}$  one can either use (11) or explicitly generate  $T_j^{\text{expl.}} = Q_j^*(A Q_j)$  by adding new columns and rows to  $T_{j-1}^{\text{expl.}}$ . In both approaches an additional  $n \times j$  array,  $S_j \in \mathbb{C}^{n \times j}$  or  $W_j := A Q_j$ , has to be stored. Using  $T_j^{\text{expl.}}$  leads to the *inexact rational Arnoldi relation*

$$A Q_j = Q_j T_j^{\text{expl.}} + \tilde{g}_j H_j^{-1}, \quad \tilde{g}_j := g_j h_{j+1,j} e_j^* - (I - Q_j Q_j^*) S_j, \quad (12)$$

which we employ in the subsequent investigations.

**Lemma 1.** Consider the approximate solution  $Q_j Y_j Q_j^*$  of (1) after  $j$  iterations of inexact RKSM, where  $Y_j$  solves the projected Lyapunov equation (7) defined by either  $T_j^{\text{expl.}}$  or  $T_j^{\text{impl.}}$ . Then the true Lyapunov residual matrix  $\mathcal{R}_j^{\text{true}} = \mathcal{R}(Q_j Y_j Q_j^*)$  can be written in the following forms.

(a) If  $T_j^{\text{expl.}}$  is used it holds

$$\mathcal{R}_j^{\text{true}} = \tilde{F}_j + \tilde{F}_j^*, \quad \tilde{F}_j := \tilde{g}_j H_j^{-1} Y_j Q_j^* = F_j - (I - Q_j Q_j^*) S_j H_j^{-1} Y_j Q_j^*, \quad (13a)$$

(b) and, if otherwise  $T_j^{\text{impl.}}$  is used, it holds

$$\mathcal{R}_j^{\text{true}} = \hat{F}_j + \hat{F}_j^*, \quad \hat{F}_j := F_j - S_j H_j^{-1} Y_j Q_j^*. \quad (13b)$$

*Proof.* For case (a), using (12) immediately yields

$$\begin{aligned} \mathcal{R}_j^{\text{true}} &= \mathcal{R}(Q_j Y_j Q_j^*) = A Q_j Y_j Q_j^* + Q_j Y_j Q_j^* A^* + B B^* \\ &= [Q_j T_j^{\text{expl.}} + \tilde{g}_j H_j^{-1}] Y_j Q_j^* + Q_j Y_j [Q_j T_j^{\text{expl.}} + \tilde{g}_j H_j^{-1}]^* + Q_j Q_j^* B B^* Q_j Q_j^* \\ &= Q_j \left[ T_j^{\text{expl.}} Y_j + Y_j (T_j^{\text{expl.}})^* + Q_j^* B B^* Q_j \right] Q_j^* + \tilde{g}_j H_j^{-1} Y_j Q_j^* + Q_j Y_j H_j^{-*} \tilde{g}_j^* = \tilde{F}_j + \tilde{F}_j^*. \end{aligned}$$

Case (b) follows similarly using (10).  $\square$

Hence in both cases the true residual  $\mathcal{R}_j^{\text{true}}$  is a perturbation of the computed residual given by (8a)-(8b) which we denote in the remainder by  $\mathcal{R}_j^{\text{comp.}}$ . In case  $T_j^{\text{expl.}}$  is used, since  $Q_j \perp \tilde{g}_j$  one can easily see that  $\|\mathcal{R}_j^{\text{true}}\| = \|\tilde{F}_j\| = \|\tilde{g}_j H_j^{-1} Y_j\|$ , a property not shared when using  $T_j^{\text{impl.}}$ .

Our next step is to analyze the difference between true and computed residual. In the following we use  $T_j^{\text{expl.}}$  to define the projected problem.

**Definition 1.** The *residual gap* after  $j$  steps of inexact RKSM for Lyapunov equations is given by

$$\Delta\mathcal{R}_j := \mathcal{R}_j^{\text{true}} - \mathcal{R}_j^{\text{comp.}} = F_j + F_j^* - (\tilde{F}_j + \tilde{F}_j^*) =: \eta_j + \eta_j^*,$$

where  $\eta_j := F_j - \tilde{F}_j = (I - Q_j Q_j^*) S_j H_j^{-1} Y_j Q_j^*$ .

We have

$$\|\Delta\mathcal{R}_j\| = \|\eta_j + \eta_j^*\| = \|\eta_j\|,$$

due to the left and rightmost factors of  $\eta_j$  being orthogonal when  $T_j^{\text{expl.}}$  is used. Suppose the desired accuracy is so that  $\|\mathcal{R}_j^{\text{true}}\| \leq \varepsilon$ , where  $\varepsilon > 0$  is a given threshold. In practice the computed residual norms often show a decreasing behavior very similar to the exact method. However, the norm of the residual gap  $\|\eta_j\|$  indicates the attainable accuracy of the inexact rational Arnoldi method because  $\|\mathcal{R}_j^{\text{true}}\| \leq \|\mathcal{R}_j^{\text{comp.}}\| + \|\eta_j\|$  and the true residual norm is bounded by  $\|\eta_j\|$  even if  $\|\mathcal{R}_j^{\text{comp.}}\| \leq \varepsilon$ , which would indicate convergence of the computed residuals. This motivates to enforce  $\|\eta_j\| < \varepsilon$ , such that small true residual norms  $\|\mathcal{R}_j^{\text{true}}\| \leq 2\varepsilon$  are obtained overall. Since

$$\|\eta_j\| \leq \|S_j H_j^{-1} Y_j\| = \left\| \sum_{k=1}^j s_k e_k^* H_j^{-1} Y_j \right\| \leq \sum_{k=1}^j \|s_k\| \|e_k^* H_j^{-1} Y_j\|, \quad (14)$$

it is sufficient that only one of the factors in each addend in the sum is small and the other one is bounded by, say,  $\mathcal{O}(1)$  in order to achieve  $\|\eta_j\| \leq \varepsilon$ . In particular, if the  $\|e_k^* H_j^{-1} Y_j\|$  terms decrease with  $k$ , the linear residual norms  $\|s_k\|$  are allowed to increase to some extent. Hence, the solve tolerance  $\tau_j^{\text{ls}}$  of the linear solves can be relaxed in the course of the outer iteration which has coined the term *relaxation*. For this to happen, however, we first need to investigate if there is a decay of  $\|e_k^* H_j^{-1} Y_j\|$  as  $k$  increases.

## 2.3 Properties of the solution of the Galerkin system

By (14), the norm of the residual gap  $\eta_j$  strongly depends on the solution  $Y_j$  of the projected Lyapunov equation (7). We will see in Theorem 2 and Corollary 3 that the entries of  $Y_j$  decrease away from the diagonal in a manner proportional to the Lyapunov residual norm.

In the second part of this section we consider the rows of  $H_j^{-1} Y_j$ , since the residual formula (8a-8b) and the residual gap (14) depend on this quantity. It turns out that the norm of those rows also decay with the Lyapunov residual norms. Both decay bounds will later be used to develop practical relaxation criteria for achieving  $\|\eta_j\| \leq \varepsilon$ .

Consider the solution to the projected Lyapunov equation (7). We are interested in the transition from step  $k$  to  $j$ , where  $k < j$ . At first, we investigate this transition for a general Galerkin method including RKSM as a special case.

**Theorem 2.** Assume  $A + A^* < 0$  and consider a Galerkin projection method for (1) with  $T_j = Q_j^*(A Q_j)$ ,  $Q_j^* Q_j = I_j$ , and the first basis vector given by  $B = q_1 \beta$ . Let the  $k \times k$  matrix  $Y_k$  and the  $j \times j$  matrix  $Y_j$  be the solution to the projected Lyapunov equation (7)

after  $k$  and  $j$  steps of this method (e.g., Algorithm 1 with  $r = 1$ ), respectively, where  $k < j$ . Consider the  $j \times j$  difference matrix  $\Delta Y_{j,k} := Y_j - \begin{bmatrix} Y_k & 0 \\ 0 & 0 \end{bmatrix}$ , where the zero blocks are of appropriate size. Then

$$\|\Delta Y_{j,k}\| \leq c_A \|\mathcal{R}_k^{\text{true}}\|, \quad c_A := \frac{(1 + \sqrt{2})^2}{2\alpha_A}, \quad (15)$$

where  $\alpha_A := \frac{1}{2}|\lambda_{\min}(A + A^*)|$ , and  $\mathcal{R}_k^{\text{true}}$  is the Lyapunov residual matrix after  $k$  steps of Algorithm 1.

*Proof.* The residual matrix  $N_{j,k}$  of (7) w.r.t.  $T_j$  and  $\begin{bmatrix} Y_k & 0 \\ 0 & 0 \end{bmatrix}$  is given by

$$N_{j,k} := T_j \begin{bmatrix} Y_k & 0 \\ 0 & 0 \end{bmatrix} + \begin{bmatrix} Y_k & 0 \\ 0 & 0 \end{bmatrix} T_j^* + \begin{bmatrix} [\beta\beta^* & 0] & 0 \\ 0 & 0 \end{bmatrix}.$$

as  $T_j$  is built cumulatively. Since  $Q_j \begin{bmatrix} Y_k & 0 \\ 0 & 0 \end{bmatrix} Q_j^* = Q_k Y_k Q_k^*$ , it follows that  $\|N_{j,k}\| \leq \|\mathcal{R}_k^{\text{true}}\|_2$ . The difference matrix  $\Delta Y_{j,k}$  satisfies the Lyapunov equation  $T_j \Delta Y_{j,k} + \Delta Y_{j,k} T_j^* = -N_{j,k}$ , and since  $\Lambda(T_j) \subset \mathbb{C}_-$  it can be expressed via the integral

$$\Delta Y_{j,k} = \int_0^\infty e^{T_j t} N_{j,k} e^{T_j^* t} dt. \quad (16)$$

Moreover, using results from [14], we have

$$\|e^{T_j t}\| \leq (1 + \sqrt{2}) \max_{z \in \mathcal{W}(T_j)} |e^{zt}| = (1 + \sqrt{2}) e^{z \max_{z \in \mathcal{W}(T_j)} \text{Re}(z)t},$$

where  $\mathcal{W}(\cdot)$  denotes the field of values. Using the assumption  $A + A^* < 0$  it holds that  $\mathcal{W}(A) \in \mathbb{C}_-$  and, consequently,  $\mathcal{W}(T_j) \subseteq \mathcal{W}(A)$ ,  $\forall j > 0$ . Hence,

$$e^{z \max_{z \in \mathcal{W}(T_j)} \text{Re}(z)t} \leq e^{-\alpha_A t} \quad \text{with} \quad \alpha_A := \frac{1}{2}|\lambda_{\min}(A + A^*)|.$$

From (16) we obtain  $\|\Delta Y_{j,k}\| \leq \int_0^\infty \|e^{T_j t}\|^2 dt \|N_{j,k}\| \leq \frac{(1 + \sqrt{2})^2 \|\mathcal{R}_k^{\text{true}}\|}{2\alpha_A}$ . □

Theorem 2 shows that the difference matrix  $\Delta Y_{j,k}$  decays at a similar rate as the Lyapunov residual norms. The above theorem can be used to obtain results about the entries of the solution  $Y_j$  of the projected Lyapunov equation (7).

**Corollary 3.** Let the assumptions of Theorem 2 hold. For the  $\ell, i$ -th entry of  $Y_j$  we have

$$|e_\ell^* Y_j e_i| \leq c_A \|\mathcal{R}_k^{\text{true}}\|, \quad \ell, i = 1, \dots, j, \quad (17)$$

where  $k < \max(\ell, i)$  and  $c_A$  is as in Theorem 2.

*Proof.* We have

$$e_\ell^* Y_j e_i = e_\ell^* \begin{bmatrix} Y_k & 0 \\ 0 & 0 \end{bmatrix} e_i + e_\ell^* \Delta Y_{j,k} e_i. \quad (18)$$

For  $k < \ell$  the first summand vanishes and hence, using Theorem 2

$$|e_\ell^* Y_j e_i| = |e_\ell^* \Delta Y_{j,k} e_i| \leq \frac{(1 + \sqrt{2})^2}{2\alpha_A} \|\mathcal{R}_k^{\text{true}}\|, \quad k = 1, \dots, \ell - 1. \quad (19)$$

Since,  $Y_j = Y_j^*$  the indices  $\ell, i$  can be interchanged, s.t.  $k < \max(\ell, i)$  and we end up with the final bound (17).  $\square$

**Remark 1.** Note that the sequence  $\{\|\mathcal{R}_k^{\text{true}}\|\}$  is not necessarily monotonically decreasing and we can further extend the bound in (17) and obtain

$$|e_\ell^* Y_j e_i| \leq c_A \min_{p=1, \dots, \max(\ell, i)} \|\mathcal{R}_{p-1}^{\text{true}}\|, \quad \ell, i = 1, \dots, j.$$

Corollary 3 shows that the  $\ell, i$ -th entry of  $Y_j$  can be bounded from above using the Lyapunov residual norm at step  $k < \max(\ell, i)$ . This indicates that the further away from the diagonal, the smaller the entries of  $Y_j$  will become in the course of the iteration, provided the true residual norms exhibit a sufficiently decreasing behavior. We can observe this characteristic in Figure 1 in Section 4. The decay of Lyapunov solutions has also been investigated by different approaches. Especially when  $T_j$  is banded, which is, e.g., the case when a Lanczos process is applied to  $A = A^*$ , more refined decay bounds for the entries of  $Y_j$  can be established, see, e.g., [13, 27].

In exact and inexact RKSM, the formula for residuals (8), (13) and residual gaps (14) depend not only on  $Y_j$  but rather on  $H_j^{-1} Y_j$ . In particular, the rows of  $H_j^{-1} Y_j$  appear in (14) and will later be substantial for defining relaxation strategies. Hence, we shall investigate if the norms  $\|e_k^* H_j^{-1} Y_j\|$  also exhibit a decay for  $1 \leq k \leq j$ . For the last row, i.e.  $k = j$ , using (9) readily reveals

$$\|e_j^* H_j^{-1} Y_j\| \leq \frac{\|\mathcal{R}_j^{\text{comp.}}\|}{|h_{j+1,j}| \|g_j\|}, \quad \text{assuming } g_j \neq 0, h_{j+1,j} \neq 0.$$

In the spirit of Theorem 2, Corollary 3, we would like to bound the  $\ell$ th row of  $H_j^{-1} Y_j$  by the  $\ell - 1$ th computed Lyapunov residual norm  $\|\mathcal{R}_{\ell-1}^{\text{comp.}}\|$ . For this we require the following lemma showing that, motivated by similar results in [42], the first  $k$  ( $1 < k \leq j$ ) entries of  $e_k^* H_j^{-1}$  are essentially determined by the left null space of  $\underline{H}_j := H_{j+1,j}$  and  $e_{k-1}^* H_{k-1}^{-1}$  modulo scaling.

**Lemma 4.** Let  $\underline{H}_j = \begin{bmatrix} H_j \\ h_{j+1,j} e_j^* \end{bmatrix} \in \mathbb{C}^{j+1,j}$  with  $H_j$  an unreduced upper Hessenberg matrix,  $\text{rank}(\underline{H}_j) = j$ ,  $H_k := H_{1:k,1:k}$  nonsingular  $\forall 1 \leq k \leq j$ , and let  $\omega \in \mathbb{C}^{1 \times j+1}$  satisfy  $\omega \underline{H}_j = 0$ . Define the vectors  $f_m^{(k)} := e_k^* H_m^{-1}$ ,  $1 \leq k, m \leq j$ . Then it holds for  $1 \leq k \leq j$

$$f_j^{(j)} = -\frac{\omega_{1:j}}{\omega_{j+1} h_{j+1,j}}, \quad f_j^{(k)} = \frac{v_j^{(k)}}{\phi_j^{(k)}}, \quad \text{where} \quad (20a)$$

$$v_j^{(k)} := \omega_{1:j} + [0_{1,k}, [0_{1,j-k-1}, h_{j+1,j} \omega_{j+1}] H_{k+1:j,k+1:j}^{-1}], \quad \phi_j^{(k)} := v_j^{(k)} H_j e_k. \quad (20b)$$

Moreover, for  $k > 1$ , the first  $k$  entries of  $v_j^{(k)}$  can be expressed by

$$v_j^{(k)}(1:k) = [-h_{k,k-1} f_{k-1}^{(k-1)}, 1] \omega(k). \quad (20c)$$

*Proof.* See Appendix A. □

**Corollary 5.** Let the assumptions of Theorem 2, Corollary 3, and Lemma 4 hold and consider the RKSM as a special Galerkin projection method. Let  $\omega \in \mathbb{C}^{1 \times j+1}$  be a unit vector with  $\omega \underline{H}_j = 0$ , where  $\underline{H}_j = \begin{bmatrix} H_j \\ e_j^* h_{j+1,j} \end{bmatrix}$  is the upper Hessenberg matrix obtained after  $j$  steps of the rational Arnoldi decomposition. Then

$$\|e_\ell^* H_j^{-1} Y_j\| \leq \begin{cases} c_A \|e_1^* H_j^{-1}\| \|\mathcal{R}_0^{\text{true}}\|, & \ell = 1, \\ \frac{1}{\phi_j^{(\ell)} \|g_{\ell-1}\|} \|\mathcal{R}_{\ell-1}^{\text{comp.}}\| + c_A \|e_\ell^* H_j^{-1}\| \|\mathcal{R}_{\ell-1}^{\text{true}}\|, & \ell = 2, \dots, j \end{cases} \quad (21)$$

with  $g_k$  from (4) and  $\phi_j^{(\ell)}$  from Lemma 4.

*Proof.* Using Lemma 4 for  $f_j^{(\ell)} := e_\ell^* H_j^{-1}$ ,  $1 < \ell \leq j$ , and the structure (8b) of the computed Lyapunov residuals yields

$$\begin{aligned} e_\ell^* H_j^{-1} Y_j &= f_j^{(\ell)} \left( \begin{bmatrix} Y_{\ell-1} & 0 \\ 0 & 0 \end{bmatrix} + \Delta Y_{j,\ell-1} \right) = [-h_{\ell,\ell-1} f_{\ell-1}^{(\ell-1)} \omega(\ell)] / \phi_j^{(\ell)} Y_{\ell-1} + f_j^{(\ell)} \Delta Y_{j,\ell-1} \\ &= -\frac{g_{\ell-1}^* L_{\ell-1}}{\|g_{\ell-1}\|^2 \phi_j^{(\ell)}} \omega(\ell) + f_j^{(\ell)} \Delta Y_{j,\ell-1}. \end{aligned}$$

Taking norms, using that  $\omega$  is a unit norm vector, noticing for  $\ell = 1$  only the second term exists, and applying Theorem 2 gives the result. □

Corollary 5 shows that, similar to the entries of  $Y_j$ , the rows of  $H_j^{-1} Y_j$  can be bounded using the previous Lyapunov residual norm. However, due to the influence of  $H_j^{-1}$ , the occurring constants in front of the Lyapunov residual norms can potentially be very large.

## 2.4 Relaxation strategies and stopping criteria for the inner iteration

In order to achieve accurate results, the difference between the true and computed residual, the residual gap, needs to be small.

Corollary 5 indicates that  $\|e_k^* H_j^{-1} Y_j\|$  decreases with the computed Lyapunov residual, and hence  $\|s_k\|$  may be relaxed during the RKSM iteration in a manner inverse proportional to the norm of the Lyapunov residual (assuming the Lyapunov residual norm decreases).

**Theorem 6** (Theoretical relaxation strategy in RKSM). Let the assumptions of Theorem 2 and Corollary 5 hold. Assume we carry out  $j$  steps of Algorithm 1 always using the explicit projection  $T_j^{\text{expl}}$ . If we choose the solve tolerances  $\|s_k\|$ ,  $1 \leq k \leq j$  for the inexact solves within RKSM such that

$$\|s_k\| \leq \tau_k^{\text{LS}} = \begin{cases} \frac{\varepsilon}{j c_A \|e_1^* H_j^{-1}\| \|\mathcal{R}_0^{\text{true}}\|}, & k = 1, \\ \frac{\varepsilon}{\phi_j^{(k)} \|g_{k-1}\| \|\mathcal{R}_{k-1}^{\text{comp.}}\| + j c_A \|e_k^* H_j^{-1}\| \|\mathcal{R}_{k-1}^{\text{true}}\|}, & k > 1, \end{cases} \quad (22)$$

with the same notation as before, then, for the residual gap  $\|\eta_j\| < \varepsilon$  holds.

*Proof.* Consider a single addend in the sum expression (14) for the norm of the residual gap

$$\|s_k\| \|e_k^* H_j^{-1} Y_j\| \leq \|s_k\| \left( \frac{1}{\phi_j^{(k)} \|g_{k-1}\|} \|\mathcal{R}_{k-1}^{\text{comp.}}\| + c_A \|e_k^* H_j^{-1}\| \|\mathcal{R}_{\ell-1}^{\text{true}}\| \right), k > 1.$$

Choosing  $s_k$  such that (22) is satisfied for  $1 \leq k \leq j$  then gives  $\|\eta_j\| \leq \sum_{k=1}^j \frac{\varepsilon}{j} = \varepsilon$  where we have used (14), Theorem 2, Corollaries 3 and 5.  $\square$

The true norms  $\|\mathcal{R}_{k-1}^{\text{true}}\|$  can be estimated by  $\|\mathcal{R}_{k-1}^{\text{true}}\| \leq \|\mathcal{R}_{k-1}^{\text{comp.}}\| + \|\eta_{k-1}\|$ . For this, we might either use some estimation for  $\|\eta_{k-1}\|$  or simply assume that all previous residual gaps were sufficiently small, i.e.,  $\|\eta_{k-1}\| \leq \varepsilon$ .

**Practical relaxation strategies for inexact RKSM** The relaxation strategy in Theorem 6 is far from practical. First, the established bounds for the entries and rows of  $Y_j$  and  $H_j^{-1} Y_j$ , respectively, can be a vast overestimation of the true norms. Hence, the potentially large denominators in (22) result in very small solve tolerances  $\tau_k^{\text{LS}}$  and, therefore, prevent a substantial relaxation of the inner solve accuracies. Second, several quantities in the used bounds are unknown at step  $k < j$ , e.g.  $H_j^{-1}$ ,  $Y_j$ , and the constants  $\phi_j^{(k)}$ . If we knew  $H_j^{-1}$ ,  $Y_j$ , we could employ a relaxation strategy of the form  $\|s_k\| \leq \frac{\varepsilon}{j \|e_k^* H_j^{-1} Y_j\|}$  and only use Corollary 5 as theoretical indication that  $\|e_k^* H_j^{-1} Y_j\|$  decreases as the outer method converges.

In the following we therefore aim to develop a more practical relaxation strategy by trying to estimate the relevant quantity  $\|e_k^* H_j^{-1} Y_j\|$  differently using the most recent available data. Suppose an approximate solution with residual norm  $\|\mathcal{R}_k^{\text{true}}\| \leq \varepsilon$ ,  $0 < \varepsilon \ll 1$  is desired which should be found after at most  $j_{\text{max}}$  rational Arnoldi steps. This goal is achieved if  $\|\eta_{j_{\text{max}}}\| \leq \frac{\varepsilon}{2}$  and if  $\|\mathcal{R}_{j_{\text{max}}}^{\text{comp.}}\| \leq \frac{\varepsilon}{2}$  is obtained by the inexact RKSM.

Consider the left null vectors of the augmented Hessenberg matrices,  $\omega_k \underline{H}_k = 0$ ,  $k \leq j_{\text{max}}$ . It is easy to show that  $\omega_k$  can be updated recursively, in particular it is possible to compute  $\omega_m = \omega_k(1 : m + 1)$ ,  $m \leq k \leq j_{\text{max}}$ . Consequently, at the beginning of step  $k$  we already have  $\underline{H}_{k-1}$  and hence, also the first  $k$  entries of  $\omega_{j_{\text{max}}}$  without knowing the full matrix  $\underline{H}_{j_{\text{max}}}$ . Using Lemma 4 and proceeding similar as in the proof of Corollary 5

results in

$$\begin{aligned}
\|e_k^* H_{j_{\max}}^{-1} Y_{j_{\max}}\| &= \|e_k^* H_{j_{\max}}^{-1} ([Y_{k-1} \ 0] + \Delta Y_{j_{\max}, \ell-1})\| \\
&= \left\| -\frac{\omega_{j_{\max}}(k)}{\phi_{j_{\max}}^{(k)}} [h_{k,k-1} e_{k-1}^* H_{k-1}^{-1}, *] [Y_{k-1} \ 0] + e_k^* H_{j_{\max}}^{-1} \Delta Y_{j_{\max}, \ell-1} \right\| \\
&= \left\| -\frac{h_{k,k-1} \omega_k(k)}{\phi_{j_{\max}}^{(k)}} e_{k-1}^* H_{k-1}^{-1} Y_{k-1} + e_k^* H_{j_{\max}}^{-1} \Delta Y_{j_{\max}, k-1} \right\| \\
&\leq \left| \frac{h_{k,k-1} \omega_k(k)}{\phi_{j_{\max}}^{(k)}} \right| \|e_{k-1}^* H_{k-1}^{-1} Y_{k-1}\| + \|e_k^* H_{j_{\max}}^{-1} \|c_A\| \mathcal{R}_{k-1}^{\text{true}}\| \\
&\approx \left| \frac{h_{k,k-1} \omega_k(k)}{\phi_{j_{\max}}^{(k)}} \right| \|e_{k-1}^* H_{k-1}^{-1} Y_{k-1}\|,
\end{aligned}$$

if  $\|e_k^* H_{j_{\max}}^{-1} \|c_A\| \mathcal{R}_{k-1}^{\text{true}}\|$  is small. Only the scaling parameter  $\phi_{j_{\max}}^{(k)}$  contains missing data at the beginning of step  $k$ , because  $H_{k-1}$ ,  $Y_{k-1}$ ,  $\omega_k$  are known from the previous step. We suggest to omit the unknown data and use the following *practical relaxation strategy*

$$\|s_k\| \leq \tau_k^{\text{LS}} = \begin{cases} \delta \frac{\varepsilon}{j_{\max}}, & k = 1, \\ \delta \frac{\varepsilon}{j_{\max} \|h_{k,k-1} e_{k-1}^* H_{k-1}^{-1} Y_{k-1}\|}, & k > 1, \end{cases} \quad (23a)$$

where  $0 < \delta \leq 1$  is a safeguard constant to accommodate for the estimation error resulting from approximating  $\|e_k^* H_{j_{\max}}^{-1} Y_{j_{\max}}\|$  and omitting unknown quantities (e.g.,  $\phi_{j_{\max}}^{(k)}$ ).

The reader might notice by following Algorithm 1 closely, that the built up subspace at the beginning of iteration step  $k \leq j_{\max}$  is already  $k$ -dimensional and since we are using the explicit projection  $T_k^{\text{expl.}}$  to define the Galerkin systems, we can already compute  $Y_k$  directly after building  $Q_k$ . This amounts to a simple rearrangement of Algorithm 1 by moving Lines 7, 8 before Line 3. Hence, the slight variation

$$\|e_k^* H_{j_{\max}}^{-1} Y_{j_{\max}}\| \approx \left| \frac{\omega_k(k)}{\phi_{j_{\max}}^{(k)}} \right| \|[-h_{k,k-1} e_{k-1}^* H_{k-1}^{-1}, 1] Y_k\|$$

of the above estimate is obtained, which suggests the use of the (slightly different) *practical relaxation strategy*

$$\|s_k\| \leq \tau_k^{\text{LS}} = \begin{cases} \delta \frac{\varepsilon}{j_{\max}}, & k = 1, \\ \delta \frac{\varepsilon}{j_{\max} \|[-h_{k,k-1} e_{k-1}^* H_{k-1}^{-1}, 1] Y_k\|}, & k > 1. \end{cases} \quad (23b)$$

For both dynamic stopping criteria, in order to prevent too inaccurate and too accurate linear solves, it is reasonable to enforce  $\tau_k^{\text{LS}} \in [\tau_{\min}^{\text{LS}}, \tau_{\max}^{\text{LS}}]$ , where  $0 < \tau_{\min}^{\text{LS}} < \tau_{\max}^{\text{LS}} \leq 1$  indicate minimal and maximal linear solve thresholds.

The numerical examples in Section 4 show that these practical relaxation strategies are effective and can reduce the number of inner iterations for RKSM by up to 50 per cent.

## 2.5 Implementation issues and generalizations

This section contains several remarks on the implementation of the inexact RKSM, in particular the case when the right hand side of the Lyapunov equation has rank greater than one, as well as considerations of the inner iterative solver and preconditioning. We also briefly comment on extensions to generalized Lyapunov equations and algebraic Riccati equations.

**The case  $r > 1$**  The previous analysis was restricted to the case  $r = 1$  but the block situation,  $r > 1$ , can be handled similarly by using appropriate block versions of the relevant quantities, e.g.,  $q_k, w_k, s_k \in \mathbb{C}^{n \times r}$ ,  $h_{ij} \in \mathbb{C}^{r \times r}$ ,  $\omega \in \mathbb{C}^{r \times (j+1)r}$ , and  $e_k$  by  $e_k \otimes I_r$ , as well as replacing absolute values by spectral norms in the right places. When solving the linear system with  $r$  right hand sides  $q_k$ , such that,  $\|s_k\| \leq \tau_k^{\text{LS}}$  is achieved, one can either use block variants of the iterative methods (see, e.g., [39, 40]), or simply apply a single vector method and sequentially consider every column  $q_k(:, \ell)$ ,  $\ell = 1, \dots, r$  and demand that  $s_k(:, \ell) \leq \tau_k^{\text{LS}}/r$ .

**Choice of inner solver** One purpose of low-rank solvers for large matrix equations is to work in a memory efficient way. Using a long-recurrence method such as GMRES to solve unsymmetric inner linear systems defies this purpose in some sense because it requires storing the full Krylov basis. Unless a very good preconditioner is available, this can lead to significant additional storage requirements within the inexact low-rank method, where the Krylov method consumes more memory than the actual low-rank Lyapunov solution factor of interest. Therefore we exclusively used short-recurrence Krylov methods (e.g., BiCGstab) for the upcoming numerical examples defined by an unsymmetric matrix  $A$ .

**Preconditioning** The preceding relaxation strategies relate to the residuals  $s_k$  of the underlying linear systems. For enhancing the performance of the Krylov subspace methods, using preconditioners is typically inevitable. Then the inner iteration itself inherently only works with the preconditioned residuals  $s_k^{\text{Prec.}}$  which, if left or two-sided preconditioning is used, are different from the residuals  $s_k$  of the original linear systems. Since  $\|s_k^{\text{Prec.}}\| \leq \tau_k^{\text{LS}}$  does not imply  $\|s_k\| \leq \tau_k^{\text{LS}}$  this can result in linear systems solved not accurately enough to ensure small enough Lyapunov residual gaps. Hence, some additional care is needed to respond to these effects from preconditioning. The obvious approach is to use right preconditioning which gives  $\|s_k^{\text{Prec.}}\| = \|s_k\|$ .

**Complex shifts** In practice  $A, B$  are usually real but some of the shift parameters can occur in complex conjugate pairs. Then it is advised to reduce the amount of complex operations by working with a real rational Arnoldi method [29] that constructs a real rational Arnoldi decomposition and slightly modified formulae for the computed Lyapunov residuals, in particular for  $F_j$ . The actual derivations are tedious and are omitted here for the sake of brevity, but our implementation for the numerical experiments works exclusively with the real Arnoldi process.

**Generalized Lyapunov equations** Often, generalized Lyapunov equations of the form

$$AXM^* + MXA^* = -BB^*, \quad (24)$$

with an additional, nonsingular *mass matrix*  $M \in \mathbb{R}^{n \times n}$  have to be solved. Projection methods tackle (24) by implicitly running on equivalent Lyapunov equations defined by  $A_M := L_M^{-1}AL_M^{-1}$ ,  $B_M := L_M^{-1}B$  using a factorization  $M = L_M U_M$ , which could be a LU-factorization or, if  $M$  is positive definite, a Cholesky factorization ( $U_M = L_M^*$ ). Other possibilities are  $L_M = M$ ,  $U_M = I$  and  $L_M = I$ ,  $U_M = M$ . Basis vectors for the projection subspace are obtained by

$$Q_1 \beta = L_M^{-1}B, \quad (A - \xi_{j+1}M)\hat{w}_{j+1} = L_M q_j, \quad w_{j+1} = U_M \hat{w}_{j+1}. \quad (25)$$

After convergence, the low-rank approximate solution of the original problem (24) is given by  $X \approx (U_M^{-1}Q_j)Y_j(U_M^{-1}Q_j)^*$ , where  $Y_j$  solves the reduced Lyapunov equation defined by the restrictions of  $A_M$  and  $B_M$ . This requires solving extra linear systems defined by (factors of)  $M$  in certain stages of Algorithm 1: setting up the first basis vector, building the restriction  $T_j$  of  $A_M$  (either explicitly or implicitly using (5)), and recovering the approximate solution after termination. Since the coefficients of these linear system do not change throughout the iteration, often sparse factorizations of  $M$  are computed once at the beginning and reused every time they are needed. In this case the above analysis can be applied with minor changes of the form

$$s_j := L_M q_j - (A - \xi_{j+1}M)\hat{w}_{j+1}, \quad \|s_j\| \leq \tau_j^{\text{ls}}, \quad w_{j+1} = U_M \hat{w}_{j+1}.$$

for the inexact linear solves. We obtain an inexact rational Arnoldi decomposition with respect to  $A_M$  of the form

$$\begin{aligned} A_M Q_j &= Q_j \hat{T}_j^{\text{expl.}} + \hat{g}_j e_j^* h_{j+1,j} H_j^{-1} - (I - Q_j Q_j^*) L_M^{-1} S_j H_j^{-1}, \\ \text{with } \hat{T}_j^{\text{expl.}} &= Q_j^* A_M Q_j, \quad \hat{g}_j = q_{j+1} \xi_{j+1} - (I - Q_j Q_j^*) A_M q_{j+1}. \end{aligned}$$

If  $Q_j Y_j Q_j^*$  is an approximate solution of the equivalent Lyapunov equation defined by  $A_M$ ,  $B_M$ , and  $Y_j$  solves the reduced Lyapunov equation defined by  $\hat{T}_j^{\text{expl.}}$  and  $Q_j^* B_M$ , then the associated residual is

$$\begin{aligned} A_M Q_j Y_j Q_j^* + Q_j Y_j Q_j^* A_M^* + B_M B_M^* &= \hat{F}_j + \Delta \hat{F}_j + (\hat{F}_j + \Delta \hat{F}_j)^*, \\ \hat{F}_j &:= \hat{g}_j^* h_{j+1,j} H_j^{-1} Y_j Q_j^*, \quad \Delta \hat{F}_j := (I - Q_j Q_j^*) L_M^{-1} S_j H_j^{-1} Y_j Q_j^*. \end{aligned}$$

Hence, the generalized residual gap is  $\hat{\eta}_j = \Delta \hat{F}_j$ . If  $L_M = I$ ,  $U_M = M$ , bounding  $\|\hat{\eta}_j\|$  works in the same way as in the case  $M = I$ , otherwise an additional constant  $1/\sigma_{\min}(L_M)$  (or an estimation thereof) has to be multiplied to the established bounds. Allowing inexact solves of the linear systems defined by (factors of)  $M$  substantially complicates the analysis. In particular, the transformation to a standard Lyapunov equation is essentially not given exactly since in that case only a perturbed version of  $A_M$  and its restriction are available. This situation is, hence, similar to the case when no exact matrix vector products with  $A$  are available. If  $L_M \neq I$ , also  $B_M$  is not available exactly leading to further perturbations in the basis generation. For these reasons, we will not further pursue inexact solves with  $M$  or its factors. This is also motivated from practical situations, where solving with  $M$ , or computing a sparse factorization thereof, is usually much less demanding compared to factorising  $A - \xi_{j+1}M$ .

**Extended Krylov subspace methods** A special case of the rational Krylov subspace appears when only the shifts zero and infinity are used, leading to the extended Krylov subspace  $\mathcal{EK}_k(A_M, B_M) = \mathcal{K}_k(A_M, B_M) \cap \mathcal{K}_k(A_M^{-1}, A_M^{-1}B_M)$  (using the notation from the previous subsection). Usually, in the resulting extended Arnoldi process the basis is expanded by vectors from  $\mathcal{K}_k(A_M, B_M)$  and  $\mathcal{K}_k(A_M^{-1}, A_M^{-1}B_M)$  in an alternating fashion, starting with  $\mathcal{K}_k(A_M, B_M)$ . The extended Krylov subspace method (EKSM) [34] for (1) and (24) uses a Galerkin projection onto  $\mathcal{EK}_k(A_M, B_M)$ . In each step the basis is orthogonally expanded by  $w_{j+1} = [A_M q_j(:, 1:r), A_M^{-1} q_j(:, r+1:2r)]$ , where  $q_j$  contains the last  $2r$  basis vectors. This translates to the following linear systems and matrix vector products

$$U_M z = q_j(:, 1:r), \quad L_M w_{j+1}(:, 1:r) = Az$$

and  $Az = L_M q_j(:, r+1:2r) \quad w_{j+1}(:, r+1:2r) = U_M z,$

that have to be dealt with. Similar formula for the implicit restriction of  $A$  and the Lyapunov residual as in RKSM can be found in [34]. Since these coefficient matrices do not change during the iteration, a very efficient strategy is to compute, if possible, sparse factorizations of  $A, M$  once before the algorithm and reuse them in every step. Incorporating inexact linear solves by using inexact sparse factorizations  $A \approx \tilde{L}_A \tilde{U}_A$ ,  $M \approx \tilde{L}_M \tilde{U}_M$  would make it difficult to incorporate relaxed solve tolerances, since there is little reason to compute a less accurate factorization once a highly accurate one has been constructed. For the same reasons stated in the paragraph above, we restrict ourselves to the iterative solution of the linear systems defined by  $A$ . These linear systems affect only the columns in  $w_{j+1}(:, r+1:2r)$ . In particular, by proceeding as for inexact RKSM, one can show that

$$A_M Q_j = Q_{j+1} T_{j+1,j}^{\text{expl.}} - (I - Q_j Q_j^*) S_j^{\text{EK}} \quad \text{with}$$

$$T_{j+1,j}^{\text{expl.}} = Q_{j+1}^* A_M Q_j, \quad S_j^{\text{EK}} := [s_1^{\text{EK}}, \dots, s_j^{\text{EK}}], \quad s_i^{\text{EK}} := [0, U_M^{-1} s_i] \in \mathbb{C}^{n \times 2r},$$

where  $s_i := L_M q_i(:, r+1:2r) - Az$ ,  $1 \leq i \leq j$ . Note that allowing inexact solves w.r.t.  $M$  or even inexact matrix vector products with  $A, M$  would destroy the zero block columns in  $s_i^{\text{EK}}$ .

**Algebraic Riccati equations** The RKSM in Algorithm 1 can be generalized to compute low-rank approximate solutions of generalized algebraic Riccati equations (AREs)

$$AXM^* + MXA^* - MXCC^*XM + BB^* = 0, \quad C \in \mathbb{C}^{n \times p}, \quad p \ll n, \quad (26)$$

see, e.g., [38, 35]. The majority of steps in Algorithm 1 remain unchanged, the major difference is that the Galerkin system is now a reduced ARE

$$T_j Y_j + Y_j T_j^* - Y_j (Q_j^* C_M) (C_M^* Q_j) Y_j + Q_j^* B_M B_M^* Q_j = 0, \quad C_M := U_M^{-*} C,$$

to be solved. Since we do not alter the underlying rational Arnoldi process, most of the properties of the RKSM hold again, especially the residual formulas in both exact and inexact case, and a residual gap is defined again as in the Lyapunov case. Differences occur in the bounds for the entries of  $Y_j$  and rows of  $H_j^{-1} Y_j$ , since Theorem 2 cannot be

formulated in the same way. Under some additional assumptions, a bound of the form  $\|\Delta Y_{j,j-1}\| \leq c_{\text{ARE}} \|\mathcal{R}_{j-1}\|$  can be established [35], where the constant  $c_{\text{ARE}}$  is different from  $c_A$ . We leave concrete generalizations of Theorem 2, Corollaries 3, 5 for future research and only show in some experiments that relaxation strategies of the form (23a), (23b) also work for the inexact RKSM for AREs.

### 3 The inexact low-rank ADI iteration

#### 3.1 Derivation, basic properties, and introduction of inexact solves

Using the Cayley transformation  $\mathcal{C}(A, \alpha) := (A + \alpha I)^{-1}(A - \bar{\alpha}I)$ ,  $\alpha \in \mathbb{C}_-$ , (1) can be reformulated as discrete-time Lyapunov equation (symmetric Stein equation)

$$X = \mathcal{C}(A, \alpha)X\mathcal{C}(A, \alpha)^* - 2\operatorname{Re}(\alpha)\mathcal{B}(\alpha)\mathcal{B}(\alpha)^*, \quad \mathcal{B}(\alpha) := (A + \alpha I)^{-1}B.$$

For suitably chosen  $\alpha_j$  this motivates the iteration

$$X_j = \mathcal{C}(A, \alpha_j)X_{j-1}\mathcal{C}(A, \alpha_j)^* - 2\operatorname{Re}(\alpha_j)\mathcal{B}(\alpha_j)\mathcal{B}(\alpha_j)^*, \quad (27)$$

which is known as alternating directions implicit (ADI) iteration [43] for Lyapunov equations. It converges for shift parameters  $\alpha_j \in \mathbb{C}_-$  because  $\rho(\mathcal{C}(A, \alpha_j)) < 1$  and it holds [25, 32, 43, 7]

$$X_j - X = \mathcal{A}_j(X_0 - X)\mathcal{A}_j^*, \quad (28)$$

$$\mathcal{R}_j = AX_j + X_jA^* + BB^* = \mathcal{A}_j\mathcal{R}_0\mathcal{A}_j^*, \quad (29)$$

where  $\mathcal{A}_j := \prod_{i=1}^j \mathcal{C}(A, \alpha_i)$ .

A low-rank version of the ADI iteration is obtained by setting  $X_0 = 0$  in (27), exploiting that  $(A + \alpha_j I)^{-1}$  and  $(A - \bar{\alpha}_i I)$  commute for  $i, j \geq 1$ , and realising that the iterates are given in low-rank factored form  $X_j = Z_j Z_j^*$  with low-rank factors  $Z_j$  constructed by

$$\begin{aligned} Z_j &= [\gamma_1 v_1, \gamma_2 v_2, \dots, \gamma_j v_j] = [Z_{j-1}, \gamma_j v_j], \quad \gamma_j := \sqrt{-2\operatorname{Re}(\alpha_j)}, \\ v_j &= (A - \bar{\alpha}_{j-1} I)(A + \alpha_j I)^{-1} v_{j-1}, \quad j \geq 1, \quad v_1 := (A + \alpha_1 I)^{-1} B, \end{aligned} \quad (30)$$

see [25, 31] for more detailed derivations. Thus, in each step  $Z_{j-1}$  is augmented by  $r$  new columns  $v_j$ . Moreover, from (29) with  $X_0 = 0$  and (30) it is evident that

$$\mathcal{R}_j = w_j w_j^*, \quad w_j := \mathcal{A}_j B = w_{j-1} - \gamma_j^2 (A + \alpha_j I)^{-1} w_{j-1}, \quad w_0 := B, \quad (31)$$

see also [7, 46, 44]. Hence, the residual matrix has at most rank  $r$  and its norm can be cheaply computed as  $\|R_j\|_2 = \|w_j^* w_j\|_2$  which coined the name *residual factors* for the  $w_j$ . The low-rank ADI (LR-ADI) iteration using these residual factors is

$$v_j = (A + \alpha_j I)^{-1} w_{j-1}, \quad w_j = w_{j-1} + \gamma_j^2 v_j, \quad w_0 := B. \quad (32a)$$

---

**Algorithm 2:** Inexact LR-ADI iteration.

---

**Input** : Matrices  $A$ ,  $M$ ,  $B$  defining (24), set of shift parameters  $\{\alpha_1, \dots, \alpha_{j_{\max}}\} \subset \mathbb{C}_-$ , tolerance  $0 < \varepsilon \ll 1$ .

**Output:**  $Z_j \in \mathbb{C}^{n \times rj}$ , such that  $ZZ^* \approx X$ .

- 1  $w_0 = B$ ,  $Z_0 = []$ ,  $\mathcal{R}_0 = \|w_0^* w_0\|$ ,  $j = 1$ .
- 2 **while**  $\|\mathcal{R}_{j-1}\| \geq \varepsilon$  **do**
- 3     Get  $v_j$  s.t.  $s_j = w_{j-1} - (A + \alpha_j M)v_j$ ,  $\|s_j\| \leq \delta_j$ .
- 4      $w_j = w_{j-1} - 2 \operatorname{Re}(\alpha_j) Mv_j$ .
- 5      $Z_j = [Z_{j-1}, \sqrt{-2 \operatorname{Re}(\alpha_j)} v_j]$ .
- 6      $\mathcal{R}_j = \|w_j^* w_j\|_2$ .
- 7      $j = j + 1$ .

---

For generalized Lyapunov equations (24), this iteration has to be modified to (see, e.g. [7, 22])

$$v_j = (A + \alpha_j M)^{-1} w_{j-1}, \quad w_j = w_{j-1} + \gamma_j^2 Mv_j, \quad w_0 := B. \quad (32b)$$

The main computational effort in each step is obviously the solution of the shifted linear system with  $(A + \alpha_j M)$  and  $r$  right hand sides for  $v_j$  in (32). Allowing inexact linear solves but keeping the other steps in (32) unchanged results in the *inexact low-rank ADI iteration* illustrated in Algorithm 2. We point out that a different notion of an inexact ADI iteration can be found in [26] in the context of operator Lyapunov equations, where inexactness refers to the approximation of infinite dimensional operators by finite dimensional ones.

Of course, allowing  $s_j \neq 0$  will violate some of the properties that were used to derive the LR-ADI iteration. Inexact solves within the closely related Smith methods have been investigated in, e.g., [32, 41], from the viewpoint of an inexact instationary iteration (27) that led to rather conservative results on the allowed sized of the norm of the linear system residual  $s_j$ . The analysis we present here follows a different path by exploiting the well-known connection of the LR-ADI iteration to rational Krylov subspaces [24, 25, 15].

**Theorem 7** ([24, 25, 44]). The low-rank solution factors  $Z_j$  after  $j$  steps of the exact LR-ADI iteration ( $\|s_j\| = 0$ ,  $\forall j \geq 1$ ) span a (block) rational Krylov subspace:

$$\operatorname{range}(Z_j) \subseteq \left\{ (A + \alpha_1 M)^{-1} B, \dots, \left[ \prod_{i=2}^j (A + \alpha_i M)^{-1} \right] (A + \alpha_1 M)^{-1} B \right\} \quad (33)$$

Although for LR-ADI, the  $Z_j$  do not have orthonormal columns and there is no rational Arnoldi process in Algorithm 2, it is still possible to find decompositions similar to (4) and (10).

**Theorem 8.** The low-rank solution factors  $Z_j$  after  $j$  steps of the inexact LR-ADI iteration (Algorithm 2) satisfy a rational Arnoldi like decomposition

$$AZ_j = MZ_j T_j + w_j g_j^* - S_j \Gamma_j, \quad S_j := [s_1, \dots, s_j], \quad \text{where} \quad (34)$$

$$T_j = - \begin{bmatrix} \alpha_1 & 0 & \dots & 0 \\ \gamma_1 \gamma_2 & \alpha_2 & \ddots & \vdots \\ \vdots & & \ddots & 0 \\ \gamma_1 \gamma_j & \dots & & \alpha_j \end{bmatrix} \otimes I_r, \quad g_j := \begin{bmatrix} \gamma_1 \\ \vdots \\ \gamma_j \end{bmatrix} \otimes I_r, \quad \Gamma_j := \operatorname{diag}(g_j) \quad (35)$$

Moreover, the Lyapunov residual matrix associated with the approximation  $X_j = Z_j Z_j^* \approx X$  is given by

$$\mathcal{R}_j^{\text{true}} = AZ_j Z_j^* M^* + MZ_j Z_j^* A^* + BB^* = -S_j \Gamma_j Z_j^* M^* - MZ_j \Gamma_j S_j^* + w_j w_j^*. \quad (36)$$

Note that here,  $T_j$  and  $g_j$  denote different quantities than in Section 2.

*Proof.* For  $S_j = 0$  the decomposition (34) was established in [22, 45] and can be entirely derived from the relations (32). Inexact solves in the sense  $w_{j-1} - (A + \alpha_j M)v_j = s_j$  can be inserted in a straightforward way leading to (34). By construction, it holds  $B = w_0 = w_j - MZ_j g_j$  which, together with (34), gives

$$\mathcal{R}_j^{\text{true}} = MZ_j(T_j + T_j^* + g_j g_j^*)Z_j^* M^* + w_j w_j^* - S_j \Gamma_j Z_j^* M^* - MZ_j \Gamma_j S_j^*$$

and (36) follows by verifying that  $T_j + T_j^* = -g_j g_j^*$ .  $\square$

**Remark 2.** The LR-ADI iteration is in general not a typical Galerkin projection method using orthonormal bases of the search spaces and solutions of reduced matrix equations. In [44, 45] it is shown that the exact LR-ADI iteration can be seen as an implicit Petrov-Galerkin type method with a hidden left projection subspace. In particular, as we exploited in the above proof, the relation  $T_j + T_j^* + g_j g_j^* = 0$  can be interpreted as reduced matrix equation solved by the identity matrix  $I_j$ . It is also possible to state a decomposition similar to (34) which incorporates  $w_0 = B$  instead of  $w_j$  [15, 45]. This can be used to state conditions which indicate when the LR-ADI approximation satisfies a Galerkin condition [15, Theorem 3.4]. Similar to inexact RKSM, if  $s_j \neq 0$  these result do in general not hold anymore.

### 3.2 Computed Lyapunov residuals, residual gap, and relaxation strategies in inexact LR-ADI

Similar to inexact RKSM when inexact solves are allowed in LR-ADI, the computed Lyapunov residuals  $\mathcal{R}_j^{\text{comp.}} = w_j w_j^*$  are different from the true residuals  $\mathcal{R}_j^{\text{true}}$  and, thus,  $\|w_j\|^2$  ceases to be a safe way to assess the quality of the current approximation  $Z_j Z_j^*$ . In the exact case,  $s_j = 0$ , it follows from (31) and  $\rho_j := \rho(\mathcal{C}_j) < 1$  that the Lyapunov residual norms decrease in the form  $\|\mathcal{R}_j\| \leq c \rho_j^2 \|\mathcal{R}_{j-1}\|$  for some  $c \geq 1$ . For the factors  $w_j$  of the computed Lyapunov residuals  $\mathcal{R}_j^{\text{comp.}}$  in the inexact method we have the following result.

**Lemma 9.** Let  $w_j^{\text{exact}}$  be the factors of the Lyapunov residuals of the exact LR-ADI method, i.e.,  $s_i = 0, 1 \leq i \leq j$ . The factors  $w_j$  of the computed Lyapunov residuals of the inexact LR-ADI are given by

$$w_j = w_j^{\text{exact}} + \sum_{i=1}^j \left[ \prod_{k=i+1}^j C_k \right] (C_i - I) s_i.$$

*Proof.* For simplicity, let  $r = 1$ ,  $M = I_n$ , and exploit  $\mathcal{C}_j := \mathcal{C}(A, \alpha_j) = A_{-\bar{\alpha}_j} A_{\alpha_j}^{-1} = I + \gamma_j^2 A_{\alpha_j}^{-1}$ ,  $A_\alpha = A + \alpha I$ . Denoting the errors by  $\delta v_j = A_{\alpha_j}^{-1} w_{j-1} - v_j = A_{\alpha_j}^{-1} s_j$  for  $j \geq 1$ ,

we have

$$\begin{aligned} w_j &= w_{j-1} + \gamma_j^2 v_j = w_{j-1} + \gamma_j^2 A_{\alpha_j}^{-1} w_{j-1} - \gamma_j^2 \delta v_j = \mathcal{C}_j w_{j-1} - \gamma_j^2 \delta v_j \\ &= \dots = \prod_{k=1}^j \mathcal{C}_k w_0 - \sum_{i=1}^j \gamma_i^2 \left[ \prod_{k=i+1}^j \mathcal{C}_k \right] \delta v_i = \prod_{k=1}^j \mathcal{C}_k w_0 - \sum_{i=1}^j \left[ \prod_{k=i+1}^j \mathcal{C}_k \right] (\mathcal{C}_i - I) s_i \end{aligned}$$

and we notice that the first term is exactly the exact Lyapunov residual factor from (31).  $\square$

Constructing  $w_j w_j^*$  from the above formula and taking norms indicates that, by the contraction property of  $\mathcal{C}_i$ , the linear residual  $s_i$  get damped. In fact, similar to the inexact projection methods, in practice we often observe that the computed Lyapunov residual norms  $\|\mathcal{R}_j^{\text{comp.}}\| = \|w_j w_j^*\|$  also show a decreasing behavior.

In the next step we analyze the difference between the computed and true residuals, in a similar manner as we did for the RKSM in Section 2. Theorem 8 motivates the definition of a *residual gap* analogue to the inexact RKSM.

**Definition 2.** The *residual gap* after  $j$  steps of the inexact LR-ADI iteration is given by

$$\Delta \mathcal{R}_j^{\text{ADI}} := \mathcal{R}_j^{\text{true}} - \mathcal{R}_j^{\text{comp.}} = \mathcal{R}_j^{\text{true}} - w_j w_j^* = \eta_j^{\text{ADI}} + (\eta_j^{\text{ADI}})^*, \quad \eta_j^{\text{ADI}} := -S_j \Gamma_j Z_j^* M^*. \quad (37)$$

Assuming we have  $\|\mathcal{R}_j^{\text{comp.}}\| = \|w_j w_j^*\| \leq \varepsilon$  and are able to bound the residual gap, e.g.,  $\|\Delta \mathcal{R}_j^{\text{ADI}}\| \leq \varepsilon$ , then we achieve small true residual norms  $\|\mathcal{R}_j^{\text{true}}\| \leq 2\varepsilon$ . A theoretical approach for bounding  $\|\Delta \mathcal{R}_j^{\text{ADI}}\|$  is given in the next Theorem.

**Theorem 10** (Theoretical relaxation for inexact LR-ADI). Let the residual gap be given by Definition 2 with  $w_j, \gamma_j$  as in (30)-(32b). Let  $j_{\max}$  be the maximum number of steps of Algorithm 2,  $\sigma_{\min,k} := \sigma_{\min}(A + \alpha_k M)$ , and  $0 < \varepsilon < 1$  the desired residual tolerance.

(a) If, for  $1 \leq k \leq j_{\max}$ , the linear system residual satisfies

$$\|s_k\| \leq \frac{1}{2} \left( \sqrt{\|w_{k-1}\|^2 + \frac{2\varepsilon \sigma_{\min,k}}{\|M\| \gamma_k^2 j_{\max}}} - \|w_{k-1}\| \right), \quad (38a)$$

then  $\|\Delta \mathcal{R}_{j_{\max}}\| \leq \varepsilon$ .

(b) Let  $\|S_{k-1} \Gamma_{k-1} Z_{k-1}^* M^*\| \leq u_{k-1}$  with  $u_0 = 0$ . If, for  $1 \leq k \leq j_{\max}$ , the linear system residual satisfies

$$\|s_k\| \leq \frac{1}{2} \left( \sqrt{\|w_{k-1}\|^2 + 2 \left( \frac{k\varepsilon}{j_{\max}} - 2u_{k-1} \right) \frac{\sigma_{\min,k}}{\|M\| \gamma_k^2}} - \|w_{k-1}\| \right), \quad (38b)$$

then  $\|\Delta \mathcal{R}_{j_{\max}}\| \leq \varepsilon$ .

*Proof.* Consider the following estimate

$$\begin{aligned} \|\Delta \mathcal{R}_{j_{\max}}^{\text{ADI}}\| &\leq 2\|\eta_{j_{\max}}^{\text{ADI}}\| = 2\|S_{j_{\max}}\Gamma_{j_{\max}}Z_{j_{\max}}^*M^*\| \\ &\leq 2\|S_{j_{\max}-1}\Gamma_{j_{\max}-1}Z_{j_{\max}-1}^*M^*\| + 2\gamma_{j_{\max}}^2\|s_{j_{\max}}\|\|Mv_{j_{\max}}\| \end{aligned} \quad (39)$$

$$\leq 2\sum_{k=1}^{j_{\max}}\gamma_k^2\|s_k\|\|Mv_k\|. \quad (40)$$

Moreover,

$$\|Mv_k\| \leq \|M\|\|(A + \alpha_k M)^{-1}(w_{k-1} - s_k)\| \leq \|M\|\frac{(\|w_{k-1}\| + \|s_k\|)}{\sigma_{\min}(A + \alpha_k M)}. \quad (41)$$

If the linear residual norms  $\|s_k\|$  are so that each addend in the sum (40) is smaller than  $\frac{\varepsilon}{j_{\max}}$  we achieve  $\|\eta_{j_{\max}}^{\text{ADI}}\| \leq \varepsilon/2$ . With  $\phi_k := \frac{2\gamma_k^2\|M\|}{\sigma_{\min}(A + \alpha_k M)}$ ,  $\omega_k := \|w_{k-1}\| = \sqrt{\|\mathcal{R}_{k-1}\|}$  this means we require  $\phi_k(\omega_k\|s_k\| + \|s_k\|^2) \leq \frac{\varepsilon}{j_{\max}}$ . Hence, the desired largest allowed value  $\|s_k\|$  is given by the positive root of the inherent quadratic equation  $\phi_k(\omega_k\varsigma + \varsigma^2) - \frac{\varepsilon}{j_{\max}} = 0$  such that

$$\|s_k\| \leq \frac{1}{2}\left(\sqrt{\omega_k^2 + \frac{4\varepsilon}{\phi_k j_{\max}}} - \omega_k\right) \quad (42)$$

leading to the desired result (a). The second strategy can be similarly shown by using (39) and

$$\|\eta_k^{\text{ADI}}\| \leq 2u_{k-1} + \phi_k(\|w_{k-1}\|\|s_k\| + \|s_k\|^2)$$

and finding  $\|s_k\|$  such that the right hand side in the above inequality is pushed below  $\frac{k\varepsilon}{j_{\max}}$ .  $\square$

The motivation behind the second stopping strategy (38b) is that we can take the previous  $\|\eta_{k-1}^{\text{ADI}}\|$  into account. This can be helpful if at steps  $i \leq k-1$  the used iterative method produced smaller linear residual norms than demanded, such that the linear residual norms  $\|s_i\|$  are allowed to grow slightly larger in later iteration steps  $i > k-1$ .

**Practical relaxation strategies for inexact LR-ADI** The proposed stopping criteria (38) are not very practical, since the employed bound (41) will often give substantial overestimation of  $\|Mv_k\|$  by several orders of magnitude which, in turn, will result in smaller inner tolerances  $\tau^{\text{LS}}$  than actually needed. Furthermore, computing or estimating the smallest singular value of the large matrix  $A + \alpha_k M$  is possible, e.g. by Lanczos-type approaches, but the extra effort for this does not pay off. Here we proposed some variations of the above approaches that are better applicable in an actual implementation. From the algorithmic description (Algorithm 2) of the LR-ADI iteration it holds

$$\|Mv_k\| = \frac{1}{\gamma_k^2}\|w_k - w_{k-1}\| \leq \frac{1}{\gamma_k^2}(\|w_k\| + \|w_{k-1}\|). \quad (43)$$

In practice, the sequence  $\{\|w_k\|\} = \{\sqrt{\|\mathcal{R}_k^{\text{comp.}}\|}\}$  is often monotonically decreasing. Assuming  $\|w_k\| \leq \|w_{k-1}\|$  suggests to use  $\|Mv_k\| \leq \frac{2\|w_{k-1}\|}{\gamma_k^2}$  in (40) leading to the relaxation criterion

$$\|s_k\| \leq \tau_k^{\text{LS}} = \frac{\varepsilon}{4j_{\max}\sqrt{\|\mathcal{R}_{k-1}^{\text{comp.}}\|}}. \quad (44a)$$

Starting from (39), assuming  $\|S_{k-1}\Gamma_{k-1}Z_{k-1}^*M^*\| \leq u_{k-1} < \frac{(k-1)\varepsilon}{2j_{\max}}$ ,  $u_0 = 0$ , and enforcing  $\|\Delta\mathcal{R}_k^{\text{ADI}}\| < \frac{k\varepsilon}{j_{\max}}$  we obtain

$$\|s_k\| \leq \tau_k^{\text{LS}} = \frac{\frac{k\varepsilon}{j_{\max}} - 2u_{k-1}}{4\sqrt{\|\mathcal{R}_{k-1}^{\text{comp.}}\|}}. \quad (44b)$$

The second relaxation strategy (44b) requires  $\|S_{k-1}\Gamma_{k-1}Z_{k-1}^*M^*\|$  or an approximation thereof. A basic approximation can be computed via

$$\|S_{k-1}\Gamma_{k-1}Z_{k-1}^*M^*\| \leq u_{k-1} \leq u_{k-2} + \gamma_{k-1}^2\|Mv_{k-1}\|\|s_{k-1}\|. \quad (45)$$

Note that  $Mv_{k-1}$  is required anyway to continue Algorithm 2. The linear residual norms  $\|s_i\|$ ,  $1 \leq i \leq k-1$  are sometimes available as byproducts of Krylov subspace solvers for linear systems if right preconditioning is used. In case of other forms of preconditioning, the  $s_i$  might need to be computed explicitly, which requires extra matrix vector products with  $A$  (and  $M$ ), or the norms  $\|s_i\|$  have to be estimated in some other way. Similar to RKSM, for problems defined by real but unsymmetric coefficients, pairs of complex conjugated shifts can occur in LR-ADI. These can be dealt with efficiently using the machinery developed in [8, 7, 22] to ensure that the majority of operations remains in real arithmetic. By following these results it is easily shown that if steps  $k-2, k-1$  used a complex conjugated pairs of shifts, then in the formula (45) the real and imaginary parts of both  $v_{k-2}, s_{k-2}$  enter the update.

At the first look, (44) appears to allow somewhat less relaxation compared to RKSM since only the square roots of the computed Lyapunov residual norms appear in the denominator. However, the numerical examples in the next section show that with these relaxation strategies we can reduce the amount of work for inexact low-rank ADI by up to 50 per cent.

## 4 Numerical examples

In this section we consider several numerical examples and apply inexact RKSM and inexact LR-ADI with our practical relaxation strategies. The experiments were carried out in MATLAB<sup>®</sup> 2016a on a Intel<sup>®</sup>Core™2 i7-7500U CPU @ 2.7GHz with 16 GB RAM. We wish to obtain an approximate solution such that the scaled true Lyapunov residual norm satisfies

$$\mathfrak{R} := \|\mathcal{R}^{\text{true}}\|/\|B\|^2 \leq \hat{\varepsilon}, \quad 0 < \hat{\varepsilon} \ll 1,$$

i.e.,  $\varepsilon = \hat{\varepsilon}\|B\|^2$ . In all tests we desire to achieve this accuracy with  $\hat{\varepsilon} = 10^{-8}$  within at most  $j_{\max} = 50$  iteration steps. We exclusively employ dynamic shift generation strategies

Table 1: Properties and setting of used test equations including save guard constant  $\delta$  in (23) and type of preconditioner used. Here,  $\text{iLU}(X, \nu)$  and  $\text{iC}(X, \nu)$  refer to incomplete LU and, respectively, Cholesky factorization of the matrix  $X$  with drop tolerance  $\nu$ .

Example	$n$	$r$	matrices	description	$\delta$	prec.
<b>cd2d</b>	40000	1	$A \neq A^*, M = I$	finite difference discretization of $\Delta u - 100x \frac{\partial u}{\partial x} - 200y \frac{\partial u}{\partial y}$ on $[0, 1]^2$	0.01	$\text{iLU}(A, 10^{-3})$
<b>heat3d</b>	125000	4	$A = A^*, M = I$	finite difference discretization of heat equation on $[0, 1]^3$	1	$\text{iC}(-A, 10^{-2})$
<b>fem3d</b>	24389	1	$A \neq A^*, M = M^* \neq I$	finite element model of 3d convection-diffusion problem from [5]	$\frac{1}{2}$	$\text{iLU}(A, 10^{-2})$
<b>fem3d-are</b>	24389	1	$A \neq A^*, M = M^* \neq I$	extension of <b>fem3d</b> to ARE (26) from [5] including $C \in \mathbb{R}^n$	$\frac{1}{2}$	$\text{iLU}(A, 10^{-2})$

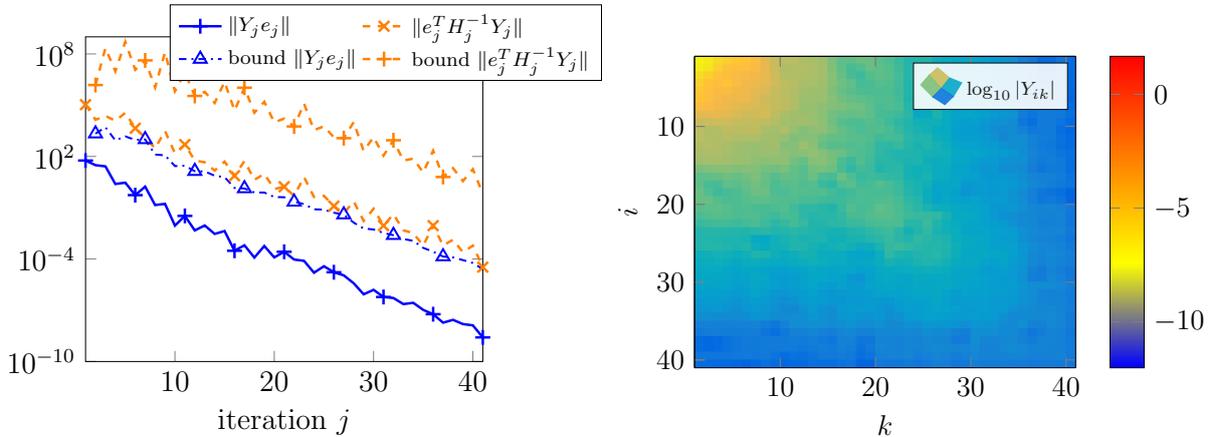
using the approach in [16] for RKSM and shifts based on a projection idea for LR-ADI, see [9, 22] for details. For the latter one, the projection basis was chosen as the last  $\min(\text{coldim}(Z), 10r)$  columns of the low-rank factor  $Z_j$ .

We employ different Krylov subspace solvers and stopping criteria for the arising linear systems of equations, in particular, we use fixed as well as dynamically chosen inner solver tolerances  $\tau^{\text{LS}}$  with the proposed relaxation strategies. The fixed solve tolerances were determined a-priori via trial-and-error such that the residual behavior of inexact method mimicked the residuals of the method using sparse-direct solution approaches for the linear systems. If not stated otherwise, the values  $\tau_{\min}^{\text{LS}} = 10^{-12}$ ,  $\tau_{\max}^{\text{LS}} = 0.1$  are taken as minimal and maximal linear solve tolerances. As Krylov subspace solvers we use BiCGstab and MINRES for problems with unsymmetric and symmetric ( $A = A^*$ ,  $M = M^*$ ) coefficients, respectively. Sparse-direct solves were carried out by the MATLAB *backslash*-routine. We consider four examples, two standard Lyapunov equations (**cd2d** and **heat3d**, where **heat3d** is an example for which  $r > 1$ ), a generalized Lyapunov equation (**fem3d**) and an algebraic Riccati equation (**fem3d-are**), in order to illustrate the theoretical results in this paper. Details and setup on the examples we used are summarized in Table 1. The matrices  $B$  for examples **cd2d** and **heat3d** were generated randomly with a standard Gaussian distribution and the initialization `randn('state', 0)`. Examples **fem3d**, **fem3d-are** provide vectors  $B, C$  [5].

At first we briefly investigate the results in Section 2.3 on the decay of Galerkin solution  $Y_j$  using example **cd2d**. We run the exact RKSM and plot the row norms of  $Y_j$ ,  $H_j^{-1}Y_j$  and the corresponding bounds obtained in Corollaries 3 and 5 against the iteration number  $j$ , as well as the absolute values of the entries of the final Galerkin solution  $Y_j$  in the left plot of Figure 1. Real shift parameters were used for this experiment. The figures are an example to show that our theoretical bounds can indeed be verified, but that they significantly overestimate the true norms. The right plot shows the decay of the entries  $Y_j$  as predicted by Corollary 3. Similar results are obtained for the other examples.

We now experiment with the different practical relaxation strategies (23) and (44) from Sections 2 and 3 for the inner iteration. In Table 2 we report the results for all examples. There, we give, among other relevant information on the performance of the

Figure 1: Illustration of the decay of the Galerkin solution for Example `cd2d`. Left: Row norms of  $Y_j$ ,  $H_j^{-1}Y_j$  and the corresponding bounds (Corollary 3, Corollary 5) against the iteration number  $j$ . Right: Absolute values of the entries of the final  $Y_j$ .



outer method under inexact inner solves, also the final obtained scaled computed residual norms  $\mathfrak{R}_j^{\text{comp}}$  (using the formula (8) for RKSM and (31) for LR-ADI). For assessing the reliability of the value of  $\mathfrak{R}_j^{\text{comp}}$ , the distance  $\delta\mathfrak{R}_j := |\mathfrak{R}_j^{\text{comp}} - \mathfrak{R}_j^{\text{true}}|$  to the true scaled residual norms is also listed, where  $\mathfrak{R}_j^{\text{true}}$  was computed using the Lanczos process on  $\mathcal{R}_j$ .

First, we observe that in all examples the difference between the true and the computed Lyapunov residual norm,  $\delta\mathfrak{R}_j$ , is of the order  $\mathcal{O}(10^{-9})$ , or smaller.

The second observation we make is that both practical relaxation criteria for RKSM, namely (23a) and (23b) are effective and lead to a reduction in overall inner iteration numbers. They both lead to nearly the same results for the number of inner iterations, the gain in using (23b) over (23a) is very minor. For our examples the savings in the total number of inner iterations between the fixed and relaxed stopping criterion within RKSM is between 28 and 50 per cent.

For the LR-ADI method we consider the two relaxation criteria (44a) and (44b). Again, we observe a reduction in the total number of iterations for both relaxation strategies, but we also see that the second relaxation criterion (44b) reduces the iteration numbers even further compared to (44a), so the use of (44b) is generally recommended. Here it pays off that the second strategy takes the previous iterate into account. The total savings in inner iterations between fixed and relaxed tolerance solves for LR-ADI is between 26 and 44 per cent in our examples.

The fairer comparison between the fixed and relaxed tolerance solve is the computation time; in Table 2 we report the computation time for direct solves as well as iterative solves with fixed and relaxed solve tolerances. For all examples we see that the relaxed solve tolerance also leads to (sometimes significant) savings in computation time. For Example `fem3d`, the generalized Lyapunov equation, the absolute time saving is not so significant (both for RKSM and LR-ADI), however, for the other three examples the savings are quite large, in particular for inexact RKSM. The results for the Riccati example `fem3d-are` indicate that the proposed relaxation criteria also work for inexact RKSM for Riccati equations. Note that for Example `cd2d`, the direct solver outperforms the iterative methods with fixed small solve tolerance (both for RKSM and LR-ADI),

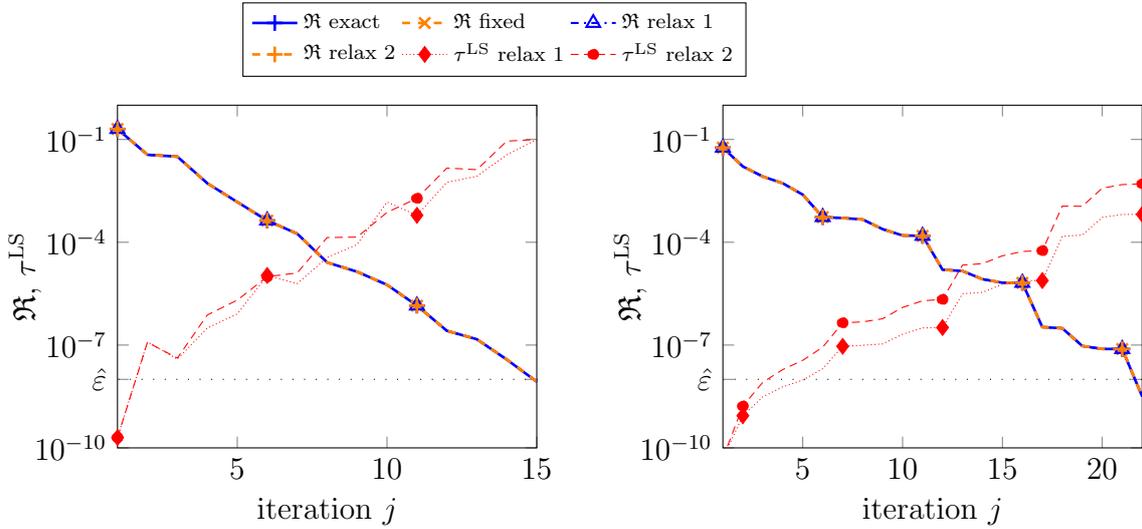
Table 2: Experimental results. The columns denote the outer (i.e., RKSM or LR-ADI) and inner method (which Krylov subspace method), the use inner stopping criterion (fixed or relaxed), the smallest and largest generated inner solve tolerances  $\min \tau^{\text{LS}}$ ,  $\max \tau^{\text{LS}}$ , the number of required outer iterations, the final obtain scaled computed residual norm  $\mathfrak{R}_j^{\text{comp}}$ , the difference  $\delta \mathfrak{R}_j := |\mathfrak{R}_j^{\text{comp}} - \mathfrak{R}_j^{\text{true}}|$  between the true and the computed Lyapunov residual norm, the total number of inner iteration steps, the achieved relative savings regarding the amount of inner iteration steps compared to the run with a fixed inner tolerance, and the computing times in seconds.

Ex.	outer	inner	stop	$\min \tau^{\text{LS}}$	$\max \tau^{\text{LS}}$	it <sup>out</sup>	$\mathfrak{R}_j^{\text{comp}}$	$\delta \mathfrak{R}_j$	it <sup>in</sup>	save	time
cd2d	RKSM	direct	–	–	–	35	5.5e-09	–	–	–	5.6
	RKSM	BICGSTAB	fixed	1.0e-10	–	35	5.5e-09	9.2e-12	789	–	5.4
	RKSM	BICGSTAB	relax (23a)	2.0e-12	1.2e-02	35	6.06e-09	5.8e-10	606	23.2%	4.3
	RKSM	BICGSTAB	relax (23b)	2.0e-12	4.4e-03	35	6.04e-09	4.9e-10	616	21.9%	4.3
	LR-ADI	direct	–	–	–	41	7.5e-09	–	–	–	6.1
	LR-ADI	BICGSTAB	fixed	1.0e-10	–	41	7.5e-09	3.6e-13	1234	–	6.3
	LR-ADI	BICGSTAB	relax (44a)	5.0e-11	9.5e-04	41	7.5e-09	1.3e-11	967	21.6%	4.9
	LR-ADI	BICGSTAB	relax (44b)	5.0e-11	5.9e-03	41	7.3e-09	5.3e-11	932	24.5%	4.7
heat3d	RKSM	direct	–	–	–	15	8.5e-09	–	–	–	88.9
	RKSM	MINRES	fixed	1.0e-09	–	15	8.5e-09	3.6e-14	487	–	26.4
	RKSM	MINRES	relax (23a)	2.0e-10	1.0e-01	15	8.53e-09	2.6e-10	260	46.6%	14.2
	RKSM	MINRES	relax (23b)	2.0e-10	1.0e-01	15	8.58e-09	3.1e-10	251	48.5%	14.2
	LR-ADI	direct	–	–	–	22	3.1e-09	–	–	–	102.8
	LR-ADI	MINRES	fixed	1.0e-09	–	22	3.1e-09	2.0e-13	510	–	20.2
	LR-ADI	MINRES	relax (44a)	5.0e-11	6.5e-04	22	3.1e-09	3.1e-12	332	34.9%	14.7
	LR-ADI	MINRES	relax (44b)	5.0e-11	5.1e-03	22	3.2e-09	6.3e-12	292	42.7%	13.1
fem3d	RKSM	direct	–	–	–	22	8.9e-09	–	–	–	39.8
	RKSM	BICGSTAB	fixed	2.00e-10	–	22	9.1e-09	2.0e-10	444	–	10.1
	RKSM	BICGSTAB	relax (23a)	3.7e-11	1.0e-01	23	6.66e-09	9.6e-10	311	30.0%	8.6
	RKSM	BICGSTAB	relax (23b)	3.9e-11	1.0e-01	23	6.95e-09	8.4e-10	302	32.0%	8.6
	LR-ADI	direct	–	–	–	27	4.2e-09	–	–	–	42.4
	LR-ADI	BICGSTAB	fixed	2.0e-10	–	27	4.2e-09	1.8e-14	359	–	2.0
	LR-ADI	BICGSTAB	relax (44a)	5.0e-11	2.7e-04	27	4.2e-09	1.8e-11	246	31.5%	1.5
	LR-ADI	BICGSTAB	relax (44b)	5.0e-11	2.8e-03	27	4.4e-09	2.5e-10	209	41.8%	1.4
fem3d-are	RKSM	direct	–	–	–	22	5.8e-09	–	–	–	43.6
	RKSM	BICGSTAB	fixed	2.00e-10	–	22	5.8e-09	4.8e-14	436	–	34.8
	RKSM	BICGSTAB	relax (23a)	1.1e-10	1.0e-01	22	6.67e-09	2.3e-09	222	49.1%	12.9
	RKSM	BICGSTAB	relax (23b)	1.1e-10	1.0e-01	22	6.67e-09	2.3e-09	222	49.1%	12.9

however, the relaxed tolerance versions perform similar to the direct methods in terms of computation time. This is to be expected as example `cd2d` represents a two-dimensional problem, where sparse-direct solvers are usually very efficient. For the three-dimensional examples `heat3d`, `fem3d` and `fem3d-are` the iterative solvers significantly outperform the direct method. Using a block-MINRES for the `heat3d` example (with  $r = 4$ ) largely led to very similar results in terms of the numbers of inner iteration steps, but the overall computing times were slightly larger.

As interesting side observation, we point out that each method generated new shift

Figure 2: Results for Example `heat3d`: Scaled computed residual norms  $\mathfrak{R}^{\text{comp}}$  and inner tolerances  $\tau^{\text{LS}}$  vs iteration numbers obtained by (in)exact RKSM (left) and LR-ADI (right)



parameters for each run, which still did not lead to substantial differences regarding the Lyapunov residual behavior, although the adaptive shift generation techniques [16, 9] are based on the eigenvalues of  $T_j$  (rational Ritz values of  $A$ ) and, hence, depend on the built up subspace. One explanation might be that the eigenvalue generation itself can be seen as inexact rational Arnoldi process for the eigenvalue problem and if the inner solve tolerances are chosen intelligently, no large differences regarding the Ritz values should appear [23, 33, 18]. Almost indistinguishable residual curves were also obtained when precomputed shifts were used for both exact and inexact methods.

Figure 2 shows the scaled computed residual norms  $\mathfrak{R}^{\text{comp}}$  and inner tolerances  $\tau^{\text{LS}}$  vs iteration numbers obtained by (in)exact RKSM and LR-ADI for Example `heat3d`.

On the left we plot the convergence history for the exact RKSM (e.g. direct solves within the inner iteration), the inexact RKSM with fixed small solve tolerance within the iterative solve and the inexact RKSM with relaxation criterion (23a) and (23b) within the iterative solution of the inner linear system. All computed residual norms are decreasing and virtually indistinguishable. The solve tolerances when using relaxation criterion (23a) are shown in dotted lines with diamonds and the ones using criterion (23b) are shown in dashed lines with red circles. The relaxation criteria lead to increasing inner solve tolerances, but as already observed in Table 2, both criteria for inexact RKSM give nearly the same results.

The right plot in Figure 2 shows the same results for LR-ADI. Again, the convergence history of the residual norms for inexact LR-ADI using the two relaxation strategies is not visibly distinguishable from the exact LR-ADI. However, we observe that the second relaxation criterion (44b), shown in dashed lines with red circles, gives better results, e.g. more relaxation and hence fewer inner iterations, than the first relaxation criterion (44a), shown in dotted lines with diamonds, a result also observed in Table 2. Similar plots as in Figure 2 can be obtained for other examples.

## 5 Conclusions and future work

The numerical solution of large scale Lyapunov equations is a very important problem in many applications. The rational Krylov subspace method (RKSM) and the low-rank alternating directions implicit (LR-ADI) iteration are well-established methods for computing low-rank solution factors of large-scale Lyapunov equations. The main task in both those methods is to solve a linear system at each step, which is usually carried out iteratively and hence inexactly.

We observed empirically that, when solving the linear system at each iteration step, the accuracy of the solve may be relaxed while maintaining the convergence to the solution of the Lyapunov equation. In this paper we have presented a theoretical foundation for explaining this phenomenon, both for the inexact RKSM method and the inexact low-rank ADI iteration. For both methods we introduced a so-called residual gap, which depends on the accuracy of the linear system solve and on quantities arising from the solution methods for the large scale Lyapunov equation. We analyzed this gap for each method which provided theoretical relaxation criteria for both inexact RKSM and inexact ADI. These criteria are often not applicable in practice as they contain unknown and/or overestimated quantities. Hence, we gave practical relaxation criteria for both methods. Our numerical results indicate that using flexible accuracies gives very good results and can reduce the amount of work for solving large scale Lyapunov equations by up to 50 per cent.

One numerical experiment with inexact RKSM indicated that relaxation strategies might also be fruitful for low-rank methods for algebraic Riccati equations [10, 9, 38, 35, 5, 4], making this an obvious future research topic, together with inexact linear solves in low-rank methods for Sylvester equations [6, 22]. In this work, we restricted ourselves to standard preconditioning techniques. Improved concepts such as *tuned* preconditioners and similar ideas [18, 11] might further enhance the performance of the inner iteration process. Preliminary tests, however, did not yield any performance gain from these techniques worth mentioning, further investigations are necessary in this direction. A further research direction worth pursuing is to reduce the computational effort for solving the sequences of shifted linear systems by storing the Krylov basis obtained from solving one (e.g., the first) linear system and employing subspace recycling techniques as, e.g., discussed for LR-ADI in [24] and for rational Krylov methods in the context of model order reduction in [1]. Allowing inexact matrix vector products, and in case of generalized equations also inexact solves with  $M$ , represents a further, more challenging research perspective.

**Acknowledgements** The authors are grateful to Cost Action EU-MORNET (TD1307) and the Department of Mathematical Sciences at Bath, that provided funding for research visits of PK to the University of Bath, where substantial parts of this work have been conducted. Furthermore, the authors thank Kirk Soodhalter (Trinity College Dublin) for insightful discussion regarding block Krylov subspace methods.

## A Appendix

*Proof of Lemma 4.* From  $\omega \underline{H}_j = 0$  we have for  $k \leq j$

$$0 = \omega \left[ \begin{array}{c|c} H_k & \\ \hline e_k^* h_{k+1,k} & H_{:,k+1:j} \\ 0_{j-k,k} & \end{array} \right],$$

which, for  $k = j$ , immediately leads to  $f_j^{(j)} = -\omega_{1:j}/(h_{j+1,j}\omega_{j+1})$ , the first equality in (20a). Similarly, it is easy to show that for  $k < j$ , a left null space vector  $\hat{\omega} \in \mathbb{C}^{1 \times k+1}$  of  $\underline{H}_k$  is given by the first  $k+1$  entries of the null vector  $\omega$  of  $\underline{H}_j$ . Hence,  $f_k^{(k)} = -\omega_{1:k}/(h_{k+1,k}\omega_{k+1})$  holds for all  $k \leq j$ .

For computing  $f_j^{(k)} = e_k^* H_j^{-1}$ ,  $k < j$ , we use the following partition of  $H_j$  and consider the splitting  $f_j^{(k)} = [u, y]$ ,  $u \in \mathbb{C}^{1 \times k}$ ,  $y \in \mathbb{C}^{1 \times k-j}$ :

$$e_k^* = [0_{1,k-1} \mid 1 \mid 0_{1,j-k}] = f_j^{(k)} H_j = [u, y] \left[ \begin{array}{c|c} H_{k-1} & \\ \hline e_{k-1}^* h_{k,k-1} & H_{:,k:j} \\ 0_{j-k,k-1} & \end{array} \right] = \left[ u \left[ \begin{array}{c} H_{k-1} \\ e_{k-1}^* h_{k,k-1} \end{array} \right] \mid [u, y] H_{:,k:j} \right].$$

This structure enforces conditions on  $[u, y]$ , which we now explore. First,  $u$  has to be a multiple of  $\omega_{1:k}$ . Here we exploited that due to the Hessenberg structure,  $\omega_{1:k} \underline{H}_{k-1} = 0$ . In particular,

$$u_{1:k-1} H_{k-1} + u_k e_{k-1}^* h_{k,k-1} = 0,$$

such that  $u_k h_{k,k-1} e_{k-1}^* H_{k-1}^{-1} = -u_{1:k-1}$ . Since  $f_{k-1}^{(k-1)} = e_{k-1}^* H_{k-1}^{-1}$  we can infer  $u_{1:k-1} = -u_k h_{k,k-1} f_{k-1}^{(k-1)}$  and, consequently, (20c) for  $k > 1$ . Similarly,  $[\omega_{1:k}, y]$  has to satisfy

$$\begin{aligned} 0_{1,j-k} &= [\omega_{1:k}, y] H_{:,k+1:j} = \omega_{1:k} H_{1:k,k+1:j} + y H_{k+1:j,k+1:j} = \omega H_{1:j+1,k+1:j} \\ &= \omega_{1:k} H_{1:k,k+1:j} + \omega_{k+1:j+1} H_{k+1:j+1,k+1:j}, \end{aligned}$$

leading to

$$y = \omega_{k+1:j+1} H_{k+1:j+1,k+1:j} H_{k+1:j,k+1:j}^{-1} = \omega_{k+1:j} + \omega_{j+1} h_{j+1,j} e_{j-k}^* H_{k+1:j,k+1:j}^{-1},$$

where  $e_{j-k}$  is a canonical vector of length  $j-k$ . Hence,

$$v_j^{(k)} = [\omega_{1:k}, y] = \omega_{1:j} + [0_{1,k}, [0_{1,j-k-1}, h_{j+1,j}\omega_{j+1}] H_{k+1:j,k+1:j}^{-1}],$$

leading to (20b). Finally, the normalization constant  $\phi_j^{(k)}$  is obtained by the requirement  $[u, y] H_{1:j,k} = 1$ .  $\square$

## References

- [1] M. I. AHMAD, D. B. SZYLD, AND M. B. VAN GIJZEN, *Preconditioned multi-shift BiCG for  $H_2$ -optimal model reduction*, SIAM J. Matrix Anal. Appl., 38 (2017), pp. 401–424.

- [2] A. C. ANTOULAS, *Approximation of large-scale dynamical systems*, vol. 6 of Advances in Design and Control, SIAM, Philadelphia, PA, 2005.
- [3] R. H. BARTELS AND G. W. STEWART, *Solution of the Matrix Equation  $AX+XB=C$ : Algorithm 432*, Comm. ACM, 15 (1972), pp. 820–826.
- [4] P. BENNER, Z. BUJANOVIĆ, P. KÜRSCHNER, AND J. SAAK, *RADI: A low-rank ADI-type algorithm for large scale algebraic Riccati equations*, Numer. Math., 138 (2018), pp. 301–330.
- [5] P. BENNER, M. HEINKENSCHLOSS, J. SAAK, AND H. K. WEICHELDT, *An inexact low-rank Newton-ADI method for large-scale algebraic Riccati equations*, Appl. Numer. Math., 108 (2016), pp. 125–142.
- [6] P. BENNER AND P. KÜRSCHNER, *Computing Real Low-rank Solutions of Sylvester equations by the Factored ADI Method*, Comput. Math. Appl., 67 (2014), pp. 1656–1672.
- [7] P. BENNER, P. KÜRSCHNER, AND J. SAAK, *An Improved Numerical Method for Balanced Truncation for Symmetric Second Order Systems*, Math. Comput. Model. Dyn. Sys., 19 (2013), pp. 593–615.
- [8] ———, *Efficient Handling of Complex Shift Parameters in the Low-Rank Cholesky Factor ADI Method*, Numer. Algorithms, 62 (2013), pp. 225–251.
- [9] ———, *Self-Generating and Efficient Shift Parameters in ADI Methods for Large Lyapunov and Sylvester Equations*, Electr. Trans. Num. Anal., 43 (2014), pp. 142–162.
- [10] P. BENNER AND J. SAAK, *Numerical solution of large and sparse continuous time algebraic matrix Riccati and Lyapunov equations: a state of the art survey*, GAMM Mitteilungen, 36 (2013), pp. 32–52.
- [11] L. BERGAMASCHI AND Á. MARTÍNEZ, *Two-stage spectral preconditioners for iterative eigensolvers*, Numer. Lin. Alg. Appl., 24 (2017), p. e2084. e2084 nla.2084.
- [12] M. BERLJAJA AND S. GÜTTEL, *Generalized Rational Krylov Decompositions with an Application to Rational Approximation*, SIAM J. Matrix Anal. Appl., 36 (2015), pp. 894–916.
- [13] C. CANUTO, V. SIMONCINI, AND M. VERANI, *On the decay of the inverse of matrices that are sum of Kronecker products*, Linear Algebra Appl., 452 (2014), pp. 21–39.
- [14] M. CROUZEIX AND C. PALENCIA, *The Numerical Range is a  $(1 + \sqrt{2})$ -Spectral Set*, SIAM J. Matrix Anal. Appl., 38 (2017), pp. 649–655.
- [15] V. DRUSKIN, L. A. KNIZHNERMAN, AND V. SIMONCINI, *Analysis of the rational Krylov subspace and ADI methods for solving the Lyapunov equation*, SIAM J. Numer. Anal., 49 (2011), pp. 1875–1898.

- [16] V. DRUSKIN AND V. SIMONCINI, *Adaptive rational Krylov subspaces for large-scale dynamical systems*, Syst. Cont. Lett., 60 (2011), pp. 546–560.
- [17] M. A. FREITAG AND A. SPENCE, *Convergence theory for inexact inverse iteration applied to the generalised nonsymmetric eigenproblem*, Electron. Trans. Numer. Anal., 28 (2007/08), pp. 40–64.
- [18] —, *Shift-invert Arnoldi’s method with preconditioned iterative solves*, SIAM J. Matrix Anal. Appl., 31 (2009), pp. 942–969.
- [19] S. GÜTTEL, *Rational Krylov Methods for Operator Functions*, PhD thesis, Technische Universität Bergakademie Freiberg, Germany, 2010. Available online from the Qucosa server.
- [20] —, *Rational Krylov approximation of matrix functions: Numerical methods and optimal pole selection*, GAMM-Mitteilungen, 36 (2013), pp. 8–31.
- [21] I. M. JAIMOUKHA AND E. M. KASENALLY, *Krylov subspace methods for solving large Lyapunov equations*, SIAM J. Numer. Anal., 31 (1994), pp. 227–251.
- [22] P. KÜRSCHNER, *Efficient Low-Rank Solution of Large-Scale Matrix Equations*, Dissertation, Otto-von-Guericke-Universität, Magdeburg, Germany, Apr. 2016.
- [23] R. B. LEHOUCQ AND K. MEERBERGEN, *Using Generalized Cayley Transformations within an Inexact Rational Krylov Sequence Method*, SIAM J. Matrix Anal. Appl., 20 (1998), pp. 131–148.
- [24] J.-R. LI, *Model Reduction of Large Linear Systems via Low Rank System Gramians*, PhD thesis, Massachusetts Institute of Technology, September 2000.
- [25] J.-R. LI AND J. WHITE, *Low Rank Solution of Lyapunov Equations*, SIAM J. Matrix Anal. Appl., 24 (2002), pp. 260–280.
- [26] M. OPMEER, T. REIS, AND W. WOLLNER, *Finite-Rank ADI Iteration for Operator Lyapunov Equations*, SIAM J. Control Optim., 51 (2013), pp. 4084–4117.
- [27] D. PALITTA AND V. SIMONCINI, *Numerical methods for large-scale Lyapunov equations with symmetric banded data*, SIAM J. Sci. Comput., (2018). Accepted for publication.
- [28] T. PENZL, *A cyclic low rank Smith method for large sparse Lyapunov equations*, SIAM J. Sci. Comput., 21 (2000), pp. 1401–1418.
- [29] A. RUHE, *The Rational Krylov algorithm for nonsymmetric Eigenvalue problems. III: Complex shifts for real matrices*, BIT, 34 (1994), pp. 165–176.
- [30] Y. SAAD, *Numerical Solution of Large Lyapunov Equation*, in Signal Processing, Scattering, Operator Theory and Numerical Methods, M. A. Kaashoek, J. H. van Schuppen, and A. C. M. Ran, eds., Birkhäuser, 1990, pp. 503–511.

- [31] J. SAAK, *Efficient Numerical Solution of Large Scale Algebraic Matrix Equations in PDE Control and Model Order Reduction*, PhD thesis, TU Chemnitz, July 2009. Available from <http://nbn-resolving.de/urn:nbn:de:bsz:ch1-200901642>.
- [32] J. SABINO, *Solution of Large-Scale Lyapunov Equations via the Block Modified Smith Method*, PhD thesis, Rice University, Houston, Texas, June 2007. Available from: [http://www.caam.rice.edu/tech\\_reports/2006/TR06-08.pdf](http://www.caam.rice.edu/tech_reports/2006/TR06-08.pdf).
- [33] V. SIMONCINI, *Variable accuracy of matrix-vector products in projection methods for eigencomputation*, SIAM J. Numer. Anal., 43 (2005), pp. 1155–1174.
- [34] —, *A new iterative method for solving large-scale Lyapunov matrix equations*, SIAM J. Sci. Comput., 29 (2007), pp. 1268–1288.
- [35] —, *Analysis of the rational Krylov subspace projection method for large-scale algebraic Riccati equations*, SIAM J. Matrix Anal. Appl., 37 (2016), pp. 1655–1674.
- [36] —, *Computational methods for linear matrix equations*, SIAM Rev., 58 (2016), pp. 377–441.
- [37] V. SIMONCINI AND D. B. SZYLD, *Theory of inexact Krylov subspace methods and applications to scientific computing*, SIAM J. Sci. Comput., 25 (2003), pp. 454–477.
- [38] V. SIMONCINI, D. B. SZYLD, AND M. MONSALVE, *On two numerical methods for the solution of large-scale algebraic Riccati equations*, IMA J. Numer. Anal., 34 (2014), pp. 904–920.
- [39] K. SOODHALTER, *A block MINRES algorithm based on the banded Lanczos method*, Numer. Algorithms, 69 (2015), pp. 473–494.
- [40] —, *Block Krylov Subspace Recycling for Shifted Systems with Unrelated Right-Hand Sides*, SIAM J. Sci. Comput., 38 (2016), pp. A302–A324.
- [41] K. SUN, *Model order reduction and domain decomposition for large-scale dynamical systems*, PhD thesis, Rice University, Houston, 2008. Available from <http://search.proquest.com/docview/304507831>.
- [42] J. VAN DEN ESHOF AND G. SLEIJPEN, *Inexact Krylov Subspace Methods for Linear Systems*, SIAM J. Matrix Anal. Appl., 26 (2004), pp. 125–153.
- [43] E. L. WACHSPRESS, *The ADI Model Problem*, Springer New York, 2013.
- [44] T. WOLF,  *$\mathcal{H}_2$  Pseudo-Optimal Model Order Reduction*, PhD thesis, Technische Universität München, 2015.
- [45] T. WOLF AND H. K.-F. PANZER, *The ADI iteration for Lyapunov equations implicitly performs  $H_2$  pseudo-optimal model order reduction*, Internat. J. Control, 89 (2016), pp. 481–493.

- [46] T. WOLF, H. K. F. PANZER, AND B. LOHMANN, *On the residual of large-scale Lyapunov equations for Krylov-based approximate solutions*, in American Control Conference, ACC 2013, Washington, DC, USA, June 17-19, 2013, 2013, pp. 2606–2611.