# On the convergence of Krylov methods with low-rank truncations

Davide Palitta[*]        Patrick Kürschner[†]

September 4, 2019

## Abstract

Low-rank Krylov methods are one of the few options available in the literature to address the numerical solution of large-scale general linear matrix equations. These routines amount to well-known Krylov schemes that have been equipped with a couple of low-rank truncations to maintain a feasible storage demand in the overall solution procedure. However, such truncations may affect the convergence properties of the adopted Krylov method. In this paper we show how the truncation steps have to be performed in order to maintain the convergence of the Krylov routine. Several numerical experiments validate our theoretical findings.

## 1 Introduction

We are interested in the numerical solution of general linear matrix equations of the form

$$\sum_{i=1}^{p} A_i X B_i^T + C_1 C_2^T = 0, \tag{1}$$

where $A_i \in \mathbb{R}^{n_A \times n_A}$, $B_i \in \mathbb{R}^{n_B \times n_B}$ are large matrices that allow matrix-vector products $A_i v$, $B_i w$ to be efficiently computed for all $i = 1, \ldots, p$, and any $v \in \mathbb{R}^{n_A}$, $w \in \mathbb{R}^{n_B}$. Moreover, $C_1$, $C_2$ are supposed to be low rank, i.e., $C_1 \in \mathbb{R}^{n_A \times q}$, $C_2 \in \mathbb{R}^{n_B \times q}$, $q \ll n_A, n_B$. For sake of simplicity we consider the case of $n_A = n_B \equiv n$ in the following, so that the solution $X \in \mathbb{R}^{n \times n}$ is a square matrix, but our analysis can be applied to the rectangular case, with $n_A \neq n_B$, as well.

Many common linear matrix equations can be written as in (1). For instance, if $p = 2$ and $B_1 = A_2 = I_n$, $I_n$ identity matrix of order $n$, we get the classical Sylvester equations. Moreover, if $B_2 = A_1$, $A_2 = B_1$, and $C_1 = C_2$, the Lyapunov equation is attained. These equations are ubiquitous in signal processing and control and systems theory. See, e.g., [1, 11, 60]. Moreover, the discretization of certain elliptic PDEs yields Lyapunov and Sylvester equations. See, e.g., [14, 40].

*Generalized* Lyapunov and Sylvester equations[1] amount to a Lyapunov/Sylvester operator plus a general linear operator:

$$AXB^T + BXA^T + \sum_{i=1}^{p-2} N_i X N_i^T + CC^T = 0, \quad \text{and} \quad A_1 X B_1 + A_2 X B_2^T + \sum_{i=1}^{p-2} N_i X M_i^T + C_1 C_2^T = 0.$$

See, e.g., [8, 26]. These equations play an important role in model order reduction of bilinear and stochastic systems, see, e.g., [8, 9, 16], and many problems arising from the discretization of PDEs can be formulated as generalized Sylvester equations as well. See, e.g., [40, 44, 63].

---

[*]palitta@mpi-magdeburg.mpg.de, Research Group Computational Methods in Systems and Control Theory (CSC), Max Planck Institute for Dynamics of Complex Technical Systems, Sandtorstraße 1, 39106 Magdeburg, Germany

[†]patrick.kurschner@kuleuven.be Department of Electrical Engineering (ESAT), ESAT/STADIUS, KU Leuven, Kasteelpark Arenberg 10, 3001 Leuven, Belgium

[1]We note that also for $p = 2$, the equations we get when $B_1 \neq I_n$, $A_2 \neq I_n$ are sometimes referred to as generalized Sylvester (Lyapunov) equations. In this work the term *generalized* always refers to the case $p > 2$ consisting of a Lyapunov/Sylvester operator plus a linear operator.

General multiterm linear matrix equation of the form (1) have been attracting attention in the very recent literature because they arise in many applications like the discretization of deterministic and stochastic PDEs, see, e.g., [5,43], PDE-constrained optimization problems ( [56]), data assimilation ( [20]), matrix regression problems arising in computational neuroscience ( [31]), fluid-structure interaction problems ( [63]), and many more.

Even when the coefficient matrices $A_i$'s and $B_i$'s in (1) are sparse, the solution $X$ is, in general, dense and it cannot be stored for large scale problems. However, for particular instances of (1), as the ones above, and under certain assumptions on the coefficient matrices, a fast decay in the singular values of $X$ can be proved and, thus, the solution admits accurate low-rank approximations of the form $S_1 S_2^T \approx X$, $S_1$, $S_2 \in \mathbb{R}^{n \times t}$, $t \ll n$, so that only the low-rank factors $S_1$ and $S_2$ need to be computed and stored. See, e.g., [4,8,26,41].

For the general multiterm linear equation (1), robust low-rank approximability properties of the solution have not been established so far even though $X$ turns out to be numerically low-rank in many cases. See, e.g., [20,56]. In the rest of the paper we thus assume that the solution $X$ to (1) admits accurate low-rank approximations.

The efficient computation of the low-rank factors $S_1$ and $S_2$ is the task of the so-called low-rank methods and many different algorithms have been developed in the last decade for both generalized and standard Lyapunov and Sylvester equations. A non complete list of low-rank methods for such equations includes projection methods proposed in, e.g., [19,26,43,48,50], low-rank (bilinear) ADI iterations ( [8,10,33]), sign function methods ( [6,7]), and Riemannian optimization methods ( [29,62]). We refer the reader to [51] for a thorough presentation of low-rank techniques.

To the best of our knowledge, few options are present in the literature for the efficient numerical solution of general equations (1): A greedy low-rank method by [28], and low-rank Krylov procedures (e.g., [8,20,30,56]) which are the focus of this paper.

Krylov methods for matrix equations can be seen as standard Krylov subspace schemes applied to the $n^2 \times n^2$ linear system

$$\mathcal{A}\mathrm{vec}(X) = -\mathrm{vec}(C_1 C_2^T), \quad \mathcal{A} := \left( \sum_{i=1}^{p} B_i \otimes A_i \right) \in \mathbb{R}^{n^2 \times n^2}, \tag{2}$$

where $\otimes$ denotes the Kronecker product and $\mathrm{vec} : \mathbb{R}^{n \times n} \to \mathbb{R}^{n^2}$ is such that $\mathrm{vec}(X)$ is the vector obtained by stacking the columns of the matrix $X$ one on top of each other.

These methods construct the Krylov subspace

$$\mathbf{K}_m(\mathcal{A}, \mathrm{vec}(C_1 C_2^T)) = \mathrm{span} \left\{ \mathrm{vec}(C_1 C_2^T), \mathcal{A}\mathrm{vec}(C_1 C_2^T), \ldots, \mathcal{A}^{m-1}\mathrm{vec}(C_1 C_2^T) \right\}, \tag{3}$$

and compute an approximate solution of the form $\mathrm{vec}(X_m) = V_m y_m \approx \mathrm{vec}(X)$, where $V_m = [v_1, \ldots, v_m] \in \mathbb{R}^{n^2 \times m}$ has orthonormal columns and it is such that $\mathrm{Range}(V_m) = \mathbf{K}_m(\mathcal{A}, \mathrm{vec}(C_1 C_2^T))$ with $y_m \in \mathbb{R}^m$. The vector $y_m$ can be computed in different ways which depend on the selected Krylov method. The most common schemes are based either on a (Petrov-)Galerkin condition on the residual vector or a minimization procedure of the residual norm, see, e.g., [46].

The coefficient matrix $\mathcal{A}$ in (2) is never assembled explicitly in the construction of $\mathbf{K}_m(\mathcal{A}, \mathrm{vec}(C_1 C_2^T))$ but its Kronecker structure is exploited to efficiently perform matrix-vector products. Moreover, to keep the memory demand low, the basis vectors of $\mathbf{K}_m(\mathcal{A}, \mathrm{vec}(C_1 C_2^T))$ must be stored in low-rank format. To this end, the Arnoldi procedure to compute $V_m$ has to be equipped with a couple of low-rank truncation steps. In particular, a low-rank truncation is performed after the "matrix-vector product" $\mathcal{A}v_m$ where $v_m$ denotes the last basis vector, and during the orthogonalization process. See, e.g., [56, Section 3], [30, Section 2], [20, Section 3] and section 2.

In principle, the truncation steps can affect the convergence of the Krylov method and the well-established properties of Krylov schemes (see, e.g., [46]) may no longer hold. However, it has been numerically observed that Krylov methods with low-rank truncations often converge, even when the truncation strategy is particularly aggressive, [20,56].

In this paper we establish some theoretical foundations to explain the converge of Krylov methods with low-rank truncations. In particular, the full orthogonalization method (FOM) [46, Section 6] and the generalized minimal residual method (GMRES) proposed in [47] are analyzed.

We assume that two different truncation steps are performed within our routine and, to show that the convergence is maintained, we interpret these truncations in two distinct ways. First, the truncation performed after the matrix-vector product $\mathcal{A}v_m$ is seen as an inexact matrix-vector product and results coming from [53] are employed. Second, the low-rank truncations that take place during the orthogonalization procedure are viewed as a structured perturbation of the new basis vector that preserves orthogonality; the perturbed vector is still orthogonal with respect to the previous ones.

The following is a synopsis of the paper. In section 2 we review the low-rank formulation of FOM and GMRES and their convergence is proved in section 3. In particular, in section 3.1-3.2 the two different interpretations of the low-rank truncation steps are presented. Implementation aspects of these low-rank truncations is discussed in section 4. It is well known that Krylov methods must be equipped with effective preconditioning techniques in order to achieve a fast convergence in terms of number of iterations. Due to some peculiar aspects of our setting, the preconditioners must be carefully designed as we discuss in section 5. Short recurrence methods like CG, MINRES and BICGSTAB can be very appealing in our context due to their small memory requirements and low computational efforts per iteration. Even though their analysis can be cumbersome since the computed basis is not always orthogonal (e.g., the orthogonality may be lost due to the short recurrences), their application to the solution of (1) is discussed in section 6. Several numerical examples reported in section 7 support our theoretical analysis. The paper finishes with some conclusions given in section 8.

Throughout the paper we adopt the following notation. The matrix inner product is defined as $\langle X, Y \rangle_F = \text{trace}(Y^T X)$ so that the induced norm is $\|X\|_F = \sqrt{\langle X, X \rangle_F}$. In the paper we continuously use the identity $\text{vec}(Y)^T \text{vec}(X) = \langle X, Y \rangle_F$ so that $\|\text{vec}(X)\|_2^2 = \|X\|_F^2$. Moreover, the cyclic property of the trace operator allows for a cheap evaluation of matrix inner products with low-rank matrices. Indeed, if $M_i, N_i \in \mathbb{R}^{n \times r_i}$, $r_i \ll n$, $i = 1, 2$, $\langle M_1 N_1^T, M_2 N_2^T \rangle_F = \text{trace}(N_2 M_2^T M_1 N_1^T) = \text{trace}((M_2^T M_1)(N_1^T N_2))$ and only matrices of small dimensions $r_i$ are involved in such a computation. Therefore, even if it is not explicitly stated, we will always assume that matrix inner products with low-rank matrices are cheaply computed without assembling any dense $n \times n$ matrix. For sake of simplicity we will omit the subscript in $\| \cdot \|_F$ and write only $\| \cdot \|$.

The $k$-th singular value of a matrix $M \in \mathbb{R}^{m_1 \times m_2}$ is denoted by $\sigma_k(M)$, where the singular values are assumed to be ordered in a decreasing fashion. The condition number of $M$ is denoted by $\kappa(M) = \sigma_1(M)/\sigma_p(M)$, $p = \text{rank}(M) = \text{argmin}_i\{\sigma_i(M) \neq 0\}$.

As already mentioned, $I_n$ denotes the identity matrix of order $n$ and the subscript is omitted whenever the dimension of $I$ is clear from the context. The $i$-th canonical basis vector of $\mathbb{R}^n$ is denoted by $e_i$ while $\mathbf{0}_m$ is a vector of length $m$ whose entries are all zero.

The brackets $[\cdot]$ are used to concatenate matrices of conforming dimensions. In particular, a Matlab-like notation is adopted and $[M, N]$ denotes the matrix obtained by stacking $M$ and $N$ one next to the other whereas $[M; N]$ the one obtained by stacking $M$ and $N$ one of top of each other, i.e., $[M; N] = [M^T, N^T]^T$. The notation $\text{diag}(M, N)$ is used to denote the block diagonal matrix with diagonal blocks $M$ and $N$.

## 2  Low-rank FOM and GMRES

In this section we revise the low-rank formulation of FOM (LR-FOM) and GMRES (LR-GMRES) for the solution of the multiterm matrix equation (1).

Low-rank Krylov methods compute an approximate solution $X_m \approx X$ of the form

$$\text{vec}(X_m) = x_0 + V_m y_m. \tag{4}$$

In the following we will always assume the initial guess $x_0$ to be the zero vector $\mathbf{0}_n$ and in Remark 3.2 such a choice is motivated. Therefore, the $m$ orthonormal columns of $V_m = [v_1, \ldots, v_m] \in \mathbb{R}^{n^2 \times m}$ in (4) span the Krylov subspace (3) and $y_m \in \mathbb{R}^m$.

One of the peculiarities of low-rank Krylov methods is that the basis vectors must be stored in low-rank format. We thus write $v_j = \text{vec}(\mathcal{V}_{1,j} \mathcal{V}_{2,j}^T)$ where $\mathcal{V}_{1,j}, \mathcal{V}_{2,j} \in \mathbb{R}^{n \times s_j}$, $s_j \ll n$, for all $j = 1, \ldots, m$.

The basis $V_m$ can be computed by a reformulation of the underlying Arnoldi process (see, e.g., [46, Section 6.4]) that exploits the Kronecker structure of $\mathcal{A}$ and the low-rank format of the basis vectors. In particular, at the $m$-th iteration, the $n^2$-vector $\widehat{v} = \mathcal{A}v_m$ must be computed. For sparse matrices $A_i$, $B_i$, a naive implementation of this operation costs $\mathcal{O}(\texttt{nnz}(\mathcal{A}))$ floating point operations (flops) where $\texttt{nnz}(\mathcal{A})$ denotes the number

of nonzero entries of $\mathcal{A}$. However, it can be replaced by the linear combination $\widehat{V} = \sum_{i=1}^{p} (A_i \mathcal{V}_{1,j}) (B_i \mathcal{V}_{2,j})^T$, $\text{vec}(\widehat{V}) = \widehat{v}$, where $2ps_j$ matrix-vector products with matrices of order $n$ are performed. The cost of such operation is $\mathcal{O}((\max_i \text{nnz}(A_i) + \max_i \text{nnz}(B_i))ps_j)$ flops and it is thus much cheaper than computing $\widehat{v}$ naively via the matrix-vector product by $\mathcal{A}$ since $\text{nnz}(\mathcal{A}) = \mathcal{O}(\max_i \text{nnz}(A_i) \cdot \max_i \text{nnz}(B_i))$, $s_j$ is supposed to be small and $p$ is in general moderate. A similar argumentation carries over when (some of) the matrices $A_i$, $B_i$ are not sparse but still allow efficient matrix vector products.

Moreover, since

$$\widehat{V} = \sum_{i=1}^{p} (A_i \mathcal{V}_{1,j}) (B_i \mathcal{V}_{2,j})^T = [A_1 \mathcal{V}_{1,j}, \ldots, A_p \mathcal{V}_{1,j}][B_1 \mathcal{V}_{2,j}, \ldots, B_p \mathcal{V}_{2,j}]^T = \widehat{V}_1 \widehat{V}_2^T, \quad \widehat{V}_1, \widehat{V}_2 \in \mathbb{R}^{n \times ps_j},$$

the low-rank format is preserved in the computation of $\widehat{V}$. In order to avoid an excessive increment in the column dimensions $ps_j$ of $\widehat{V}_1, \widehat{V}_2$, it is necessary to exercise a column compression of the factors $\widehat{V}_1$ and $\widehat{V}_2$, i.e., the matrices $(\overline{V}_1, \overline{V}_2) = \texttt{trunc}(\widehat{V}_1, I, \widehat{V}_2, \varepsilon_{\mathcal{A}})$ are computed. With $\texttt{trunc}(L, M, N, \varepsilon_{\texttt{trunc}})$ we denote any routine that computes low-rank approximations of the product $LMN^T$ with a desired accuracy of order $\varepsilon_{\texttt{trunc}}$, so that, the matrices $\overline{V}_1$, $\overline{V}_2$ are such that $\|\overline{V}_1 \overline{V}_2^T - \widehat{V}_1 \widehat{V}_2^T\|/\|\widehat{V}_1 \widehat{V}_2^T\| = \varepsilon_{\mathcal{A}}$ with $\overline{V}_1$, $\overline{V}_2 \in \mathbb{R}^{n \times \overline{s}}$, $\overline{s} \leqslant ps_j$. Algorithm 1 illustrates a standard approach for such compressions that is based on thin QR-factorizations and a SVD thereafter; see, e.g., [30, Section 2.2.1], and used in the remainder of the paper. Some alternative truncation schemes are discussed in Section 4.

---

**Algorithm 1** $\texttt{trunc}(L, M, N, \varepsilon_{\texttt{trunc}})$

---

**input** : $L, N \in \mathbb{R}^{n \times r}$, $r \ll n$, $M \in \mathbb{R}^{r \times r}$, $\varepsilon_{\texttt{trunc}} > 0$
**output:** $F, G \in \mathbb{R}^{n \times \overline{k}}$, $\overline{k} \leqslant r$, $\|FG^T - LMN^T\|/\|LMN^T\| = \varepsilon_{\texttt{trunc}}$

**1** Compute skinny QR factorizations $Q_L R_L = L$, $Q_N R_N = N$
**2** Compute the SVD decomposition $U\Sigma W^T = R_L M R_N^T \in \mathbb{R}^{r \times r}$, $U = [u_1, \ldots, u_r]$, $W = [w_1, \ldots, w_r]$ $\Sigma = \text{diag}(\sigma_1, \ldots, \sigma_r)$, $\sigma_1 \geqslant \cdots \geqslant \sigma_r \geqslant 0$
**3** Find the smallest index $\overline{k}$ such that $\sqrt{\sum_{i=\overline{k}+1}^{r} \sigma_i} \leqslant \varepsilon_{\texttt{trunc}} \|\Sigma\|$
**4** Define $F := Q_L([u_1, \ldots, u_{\overline{k}}]\sqrt{\text{diag}(\sigma_1, \ldots, \sigma_{\overline{k}})})$ and $G := Q_N([w_1, \ldots, w_{\overline{k}}]\sqrt{\text{diag}(\sigma_1, \ldots, \sigma_{\overline{k}})})$

---

The vector $\text{vec}(\overline{V}_1 \overline{V}_2^T) \approx \widehat{v}$ returned by the truncation algorithm is then orthogonalized with respect to the previous basis vectors $\text{vec}(\mathcal{V}_{1,j} \mathcal{V}_{1,j}^T)$, $j = 1, \ldots, m$. Such an orthogonalization step can be implemented by performing, e.g., the modified Gram-Schmidt procedure and the low-rank format of the quantities involved can be exploited and maintained in the result. The vector formulation of the orthogonalization step is given by

$$\widetilde{v} = \text{vec}(\overline{V}_1 \overline{V}_2^T) - \sum_{j=1}^{m} \left(\text{vec}(\mathcal{V}_{1,j} \mathcal{V}_{2,j}^T)^T \text{vec}(\overline{V}_1 \overline{V}_2^T)\right) \text{vec}(\mathcal{V}_{1,j} \mathcal{V}_{2,j}^T), \tag{5}$$

and, since $\text{vec}(\mathcal{V}_{1,j} \mathcal{V}_{2,j}^T)^T \text{vec}(\overline{V}_1 \overline{V}_2^T) = \langle \mathcal{V}_{1,j} \mathcal{V}_{2,j}^T, \overline{V}_1 \overline{V}_2^T \rangle_F$, we can reformulate (5) as

$$\widetilde{V} = \overline{V}_1 \overline{V}_2^T - \sum_{j=1}^{m} h_{j,m} \mathcal{V}_{1,j} \mathcal{V}_{2,j}^T = [\overline{V}_1, \mathcal{V}_{1,1}, \ldots, \mathcal{V}_{1,m}]\Theta_m[\overline{V}_2, \mathcal{V}_{2,1}, \ldots, \mathcal{V}_{2,m}]^T, \quad h_{j,m} = \langle \mathcal{V}_{1,j} \mathcal{V}_{2,j}^T, \overline{V}_1 \overline{V}_2^T \rangle_F,$$

where $\Theta_m = \text{diag}(I_{\overline{s}}, -h_{1,m} I_{s_1}, \ldots, -h_{m,m} I_{s_m})$, $\text{vec}(\widetilde{V}) = \widetilde{v}$, and the $m$ coefficients $h_{j,m}$ are collected in the $m$-th column of an upper Hessenberg matrix $H_m \in \mathbb{R}^{m \times m}$. Obviously, the result $\widetilde{V}$ has factors with increased column dimensions such that a truncation of the matrix $[\overline{V}_1, \mathcal{V}_{1,1}, \ldots, \mathcal{V}_{1,m}]\Theta_m[\overline{V}_2, \mathcal{V}_{2,1}, \ldots, \mathcal{V}_{2,m}]^T$ becomes necessary. In particular, if $\varepsilon_{\texttt{orth}}$ is a given threshold, we compute

$$(\widetilde{V}_1, \widetilde{V}_2) = \texttt{trunc}([\overline{V}_1, \mathcal{V}_{1,1}, \ldots, \mathcal{V}_{1,m}], \Theta_m, [\overline{V}_2, \mathcal{V}_{2,1}, \ldots, \mathcal{V}_{2,m}], \varepsilon_{\texttt{orth}}). \tag{6}$$

The result in (6) is then normalized to obtained the $(m+1)$-th basis vector, namely $\mathcal{V}_{1,m+1} = \widetilde{V}_1/\sqrt{\|\widetilde{V}_1 \widetilde{V}_2^T\|}$ and $\mathcal{V}_{2,m+1} = \widetilde{V}_2/\sqrt{\|\widetilde{V}_1 \widetilde{V}_2^T\|}$. The upper Hessenberg matrix $\underline{H}_m \in \mathbb{R}^{(m+1) \times m}$ is defined such that its square principal submatrix is given by $H_m$ and $e_{m+1}^T \underline{H}_m e_m = h_{m+1,m} := \|\widetilde{V}_1 \widetilde{V}_2^T\|$.

The difference between FOM and GMRES lies in the computation of the vector $y_m$ in (4). In FOM a Galerkin condition on the residual vector

$$\text{vec}(C_1 C_2^T) + \mathcal{A}\text{vec}(X_m) \perp \mathbf{K}_m(\mathcal{A}, \text{vec}(C_1 C_2^T)). \tag{7}$$

is imposed. If no truncation steps are performed during the Arnoldi procedure, the Arnoldi relation

$$\mathcal{A}V_m = H_m V_m + h_{m+1,m}\text{vec}(\mathcal{V}_{1,m+1}\mathcal{V}_{2,m+1}^T)e_m^T, \tag{8}$$

is fulfilled and it is easy to show that imposing the Galerkin condition (7) is equivalent solving the $m \times m$ linear system

$$H_m y_m^{fom} = \beta e_1, \quad \beta = \|C_1 C_2^T\| \tag{9}$$

for $y_m = y_m^{fom}$. Moreover, in the exact setting where (8) holds, the norm of the residual vector $\text{vec}(C_1 C_2^T) + \mathcal{A}\text{vec}(X_m)$ can be cheaply computed as

$$\|\text{vec}(C_1 C_2^T) + \mathcal{A}\text{vec}(X_m)\| = h_{m+1,m}|e_m^T y_m^{fom}|.$$

See, e.g., [46, Proposition 6.7]. We show later that this is possible also when the low-rank truncations are performed and an *inexact* version of (8) is taken into account.

In GMRES, the vector $y_m = y_m^{gm}$ is computed by solving a least squares problem

$$y_m^{gm} = \underset{y_m}{\text{argmin}}\,\|\text{vec}(C_1 C_2^T) + \mathcal{A}V_m y_m\|,$$

which corresponds to the Petrov-Galerkin orthogonality condition

$$\text{vec}(C_1 C_2^T) + \mathcal{A}\text{vec}(X_m) \perp \mathcal{A} \cdot \mathbf{K}_m(\mathcal{A}, \text{vec}(C_1 C_2^T)). \tag{10}$$

If (8) holds, $y_m^{gm}$ can be computed as

$$y_m^{gm} = \underset{y_m}{\text{argmin}}\,\|\beta e_1 + \underline{H}_m y_m\|, \tag{11}$$

and, following the discussion in [46, Section 6.5.3], this reduced least squares problem can be cheaply solved by applying $m$ Givens rotations $\Omega_i$. If $\underline{U}_m = \prod_{i=1}^m \Omega_i \underline{H}_m \in \mathbb{R}^{(m+1)\times m}$ is upper triangular and $\underline{g}_m = \beta \prod_{i=1}^m \Omega_i e_1 \in \mathbb{R}^{m+1}$, then the vector $y_m^{gm}$ is given by the solution of the $m \times m$ linear system $U_m y_m^{gm} = g_m$ where $U_m$ denotes the square principal submatrix of $\underline{U}_m$ and $g_m$ collects the first $m$ components of $\underline{g}_m$. Moreover,

$$\|\text{vec}(C_1 C_2^T) + \mathcal{A}\text{vec}(X_m)\| = |e_{m+1}^T \underline{g}_m|.$$

See, e.g., [46, Proposition 6.9]. As for FOM, we will show that this is possible also in the case of GMRES equipped with low-rank truncations.

If at the $m$-th iteration the residual norm $\|\text{vec}(C_1 C_2^T) + \mathcal{A}V_m y_m\|$ is sufficiently small[2], we recover the solution $X_m$. Clearly, the full $X_m$ is not constructed efficiently as this is a large, dense matrix. However, since we have assumed that the solution $X$ to (1) admits accurate low-rank approximations, we can compute low-rank factors $S_1$, $S_2 \in \mathbb{R}^{n\times t}$, $t \ll n$, such that $S_1 S_2^T \approx X$. Also this operation can be performed by exploiting the low-rank format of the basis vectors. In particular, if $\Upsilon = \text{diag}((e_1^T y_m)I_{s_1}, \ldots, (e_m^T y_m)I_{s_m})$, then

$$(S_1, S_2) = \texttt{trunc}([\mathcal{V}_{1,1}, \ldots \mathcal{V}_{1,m}], \Upsilon, [\mathcal{V}_{1,2}, \ldots \mathcal{V}_{2,m}], \varepsilon). \tag{12}$$

The low-rank FOM and GMRES procedures are summarized in Algorithm 2. For sake of simplicity, we decide to collect the two routines in the same pseudo-algorithm as they differ only in the convergence check if a Givens rotations approach similar to the one presented for GMRES is adopted also for FOM. This allows for a cheap evaluation of the residual norm without solving the linear system (9) at each iteration.

At each iteration step $m$ of Algorithm 2 we perform three low-rank truncations[3] and these operations substantially influence the overall solution procedure. If the truncation tolerances $\varepsilon_\mathcal{A}$ and $\varepsilon_{\texttt{orth}}$ are chosen

---

[2] $y_m = y_m^{fom}$ or $y_m = y_m^{gm}$.

[3] One after the application of $\mathcal{A}$ in line 3, and two during the orthogonalization procedure in line 7, at the end of each of the two loops of the modified Gram-Schmidt method.

**Algorithm 2** LR-FOM and LR-GMRES

---

**input** : $A_i, B_i \in \mathbb{R}^{n \times n}$, for $i = 1, \ldots, p$, $C_1, C_2 \in \mathbb{R}^{q \times n}$, $m_{\max}$, $\varepsilon_{\mathcal{A}}$, $\varepsilon_{\mathtt{orth}}$, $\varepsilon > 0$

**output:** $S_1, S_2 \in \mathbb{R}^{n \times t}$, $t \ll n$, $S_1 S_2^T \approx X$ approximate solution to (1)

**1** Compute $\beta = \|C_1 C_2^T\|$ and set $\Omega_1 = 1$, $\underline{g}_1 = \beta e_1$, $\mathcal{V}_{1,1} = C_1 / \sqrt{\beta}$ and $\mathcal{V}_{2,1} = C_2 / \sqrt{\beta}$

    **for** $m = 1, 2, \ldots$, *till* $m_{\max}$ **do**

**2**     Set $\widehat{V}_1 = [A_1 \mathcal{V}_{1,m}, \ldots, A_p \mathcal{V}_{1,m}]$ and $\widehat{V}_2 = [B_1 \mathcal{V}_{2,m}, \ldots, B_p \mathcal{V}_{2,m}]$

**3**     Compute $(\overline{V}_1, \overline{V}_2) = \mathtt{trunc}(\widehat{V}_1, I, \widehat{V}_2, \varepsilon_{\mathcal{A}})$

**4**     Set $h_{j,m} = 0$ for $j = 1, \ldots, m$

        **for** $\ell = 1, 2$ **do**

            **for** $j = 1, \ldots, m$ **do**

**5**             Compute $h_{j,m} = h_{j,m} + \langle \mathcal{V}_{1,j} \mathcal{V}_{2,j}^T, \overline{V}_1 \overline{V}_2^T \rangle_F$ and collect it in $\underline{H}_m e_m \in \mathbb{R}^{m+1}$

            **end**

**6**         Set $\Theta_m = \mathrm{diag}(I_{\overline{s}}, -h_{1,m} I_{s_1}, \ldots, -h_{m,m} I_{s_m})$, $\overline{s} = \mathrm{rank}(\overline{V}_1)$

**7**         Compute $(\overline{V}_1, \overline{V}_2) = \mathtt{trunc}([\overline{V}_1, \mathcal{V}_{1,1}, \ldots, \mathcal{V}_{1,m}], \Theta_m, [\overline{V}_2, \mathcal{V}_{2,1}, \ldots, \mathcal{V}_{2,m}], \varepsilon_{\mathtt{orth}})$

        **end**

**8**     Set $e_{m+1}^T \underline{H}_m e_m = h_{m+1,m} = \|\overline{V}_1 \overline{V}_2^T\|$

**9**     Set $\mathcal{V}_{1,m+1} = \overline{V}_1 / \sqrt{h_{m+1,m}}$ and $\mathcal{V}_{2,m+1} = \overline{V}_2 / \sqrt{h_{m+1,m}}$

        **if** $m = 1$ **then**

**10**         Set $\underline{U}_1 = \mathrm{diag}(\Omega_1, 1) \underline{H}_1 e_1$

        **else**

**11**         Set $\underline{U}_m = [[\underline{U}_{m-1}; \mathbf{0}_m^T], \prod_{i=1}^m \mathrm{diag}(\Omega_i, I_{m+1-i}) \underline{H}_m e_m]$

        **end**

        **if FOM** and $|h_{m+1,m}(e_m^T \underline{g}_m)/(e_m^T \underline{U}_m e_m)| < \varepsilon \cdot \beta$ **then**

**12**         **Break** and go to **16**

        **end**

**13**     Compute $\Omega_{m+1} \in \mathbb{R}^{(m+1) \times (m+1)}$ such that $\underline{U}_m = \Omega_{m+1} \underline{U}_m$ is upper triangular

**14**     Set $\underline{g}_{m+1} = \mathrm{diag}(\Omega_{m+1}, 1)[\underline{g}_m; 0]$

        **if GMRES** and $|e_{m+1}^T \underline{g}_{m+1}| < \varepsilon \cdot \beta$ **then**

**15**         **Break** and go to **16**

        **end**

    **end**

**16** Set $U_m = [I_m, \mathbf{0}_m] \underline{U}_m [I_m; \mathbf{0}_m^T] \in \mathbb{R}^{m \times m}$ and $g_m = [I_m, \mathbf{0}_m] \underline{g}_m \in \mathbb{R}^m$

**17** Compute $y_m = U_m^{-1} g_m$

**18** Set $\Upsilon = \mathrm{diag}((e_1^T y_m) I_{s_1}, \ldots, (e_m^T y_m) I_{s_m})$

**19** Compute $(S_1, S_2) = \mathtt{trunc}([\mathcal{V}_{1,1}, \ldots \mathcal{V}_{1,m}], \Upsilon, [\mathcal{V}_{1,2}, \ldots \mathcal{V}_{2,m}], \varepsilon)$

---

too large, the whole Krylov method my break down. Therefore, in the following sections we discuss how to adaptively choose the truncation tolerances $\varepsilon_{\mathcal{A}}$ and $\varepsilon_{\texttt{orth}}$ to maintain convergence. Moreover, the low-rank truncation does have its own computational workload which can be remarkable, especially if the ranks of the basis vectors involved is quite large. In section 4 we discuss some computational appealing alternatives to Algorithm 1.

# 3 A convergence result

In this section we show that the convergence of LR-FOM and LR-GMRES is guaranteed if the thresholds $\varepsilon_{\mathcal{A}}$ and $\varepsilon_{\texttt{orth}}$ for the low-rank truncations in line 3 and 7 of Algorithm 2 are properly chosen and if the routine used in the truncation steps satisfies certain properties.

The truncation that takes place in line 19, after the iterative process terminated, to recover the low-rank factors of the approximate solution is not discussed. Indeed, this does not affect the convergence of the Krylov method and it is justified by assuming that the exact solution $X$ admits low-rank approximations.

## 3.1 Inexact matrix-vector products

We start by analyzing the truncation step in line 3 of Algorithm 2 assuming, for the moment, that the one in line 7 is not performed. In this way the generated basis $V_m$ is ensured to be orthogonal. In section 3.2 we will show that the truncation in line 7 of Algorithm 2 preserves the orthogonality of the constructed basis so that the results we show here still hold.

The low-rank truncation performed in line 3 of Algorithm 2 can be understood as an inexact matrix-vector product with $\mathcal{A}$. Indeed, at the $m$-th iteration, we can write

$$\widehat{V}_1 \widehat{V}_2^T = \overline{V}_1 \overline{V}_2^T + E_m,$$

where $E_m$ is the matrix discarded when $\texttt{trunc}(\widehat{V}_1, I, \widehat{V}_2, \varepsilon_{\mathcal{A}})$ is applied so that $\|E_m\|/\|\widehat{V}_1 \widehat{V}_2^T\| \leqslant \varepsilon_{\mathcal{A}}$. Therefore, we have

$$\text{vec}(\overline{V}_1 \overline{V}_2^T) = \mathcal{A}\text{vec}(\mathcal{V}_{1,m} \mathcal{V}_{2,m}^T) - \text{vec}(E_m), \quad \|\text{vec}(E_m)\| \leqslant \varepsilon_{\mathcal{A}} \cdot \|\mathcal{A}\text{vec}(\mathcal{V}_{1,m} \mathcal{V}_{2,m}^T)\|,$$

and the vector $\text{vec}(\overline{V}_1 \overline{V}_2^T)$ can thus be seen as the result of an inexact matrix-vector product by $\mathcal{A}$.

Following the discussion in [53], the Arnoldi relation (8) must be replaced with the inexact counterpart

$$\mathcal{A}V_m - [\text{vec}(E_1), \dots, \text{vec}(E_m)] = V_m H_m + h_{m+1,m}\text{vec}(\mathcal{V}_{1,m+1} \mathcal{V}_{2,m+1}^T)e_m^T, \tag{13}$$

and $\text{Range}(V_m)$ is no longer a Krylov subspace generated by $\mathcal{A}$.

The vectors $y_m^{fom}$ and $y_m^{gm}$ can be still calculated as in (9) and (11), respectively, but these are no longer equivalent to imposing the Galerkin and Petrov-Galerkin conditions (7)-(10) since the Arnoldi relation (8) no longer holds; different constraints must be taken into account.

**Proposition 3.1** (See [53]). *Let* (13) *hold and define* $W_m = \mathcal{A}V_m - [vec(E_1), \dots, vec(E_m)]$. *If* $y_m^{gm}$ *is computed as in* (11), *then* $q_m^{gm} := W_m y_m^{gm}$ *is such that*

$$q_m = \underset{q \in Range(W_m)}{\text{argmin}} \|vec(C_1 C_2^T) + q\|.$$

*Similarly, if* $y_m^{fom}$ *is computed as in* (9), *then* $q_m^{fom} := W_m y_m^{fom}$ *is such that*

$$vec(C_1 C_2^T) + q_m \perp Range(V_m).$$

Consequently, $H_m$ is not a true Galerkin projection of $\mathcal{A}$ onto $\text{range}(V_m)$. One may want to compute the vectors $y_m^{fom}$ and $y_m^{gm}$ by employing the true projection $T_m := V_m^T \mathcal{A} V_m = H_m + V_m^T[\text{vec}(E_1), \dots, \text{vec}(E_m)]$ in place of $H_m$ in (9)-(11) so that the reduced problems represent a better approximation (cf. [22]) of the original equation and the orthogonality conditions imposed are in terms of the true residual. However, the computation of $T_m$ requires to store the matrix $[\text{vec}(E_1), \dots, \text{vec}(E_m)]$ and this is impracticable as

the benefits in terms of memory demand coming from the low-rank truncations are completely lost due to the allocation of both $V_m$ and $[\text{vec}(E_1), \ldots, \text{vec}(E_m)]$. A different option is to store the matrix $\mathcal{A}V_m$ and compute an explicit projection of $\mathcal{A}$ onto the current subspace, but also this strategy leads to an unfeasible increment in the memory requirements of the overall solution process as the storage demand grows of a factor $p$. Therefore, in all the numerical experiments reported in section 7, the matrix $H_m$ arising from the orthonormalization procedure is employed in the computation of $y_m^{fom}$ and $y_m^{gm}$.

If (13) holds and $\text{vec}(X_m) = V_m y_m$ is the approximate solution to (2) computed by projection onto $\text{Range}(V_m)$, then, at the $m$-th iteration, the true residual vector can be expressed as

$$r_m = \text{vec}(C_1 C_2^T) + \mathcal{A}\text{vec}(X_m) = \text{vec}(C_1 C_2^T) + \mathcal{A}V_m y_m = \widetilde{r}_m - [\text{vec}(E_1), \ldots, \text{vec}(E_m)]y_m, \qquad (14)$$

where $\widetilde{r}_m$ is the computed residual vector.

In [53, Section 4] it has been shown that the residual gap $\delta_m := \|r_m - \widetilde{r}_m\|$ between the true residual and the computed one can be bounded by

$$\delta_m \leqslant \sum_{j=1}^{m} \|E_j\| \cdot |e_j^T y_m|.$$

Since $|e_j^T y_m|$ decreases as the the iterations proceed (see, e.g., [53, Lemma 5.1-5.2]), $\|E_m\|$ is allowed to increase while still maintaining a small residual gap and preserving the convergence of the overall solution process. This phenomenon is often referred to as *relaxation*.

**Theorem 3.1** (See [53]). *Let $\varepsilon > 0$ and let $r_m^{gm} := vec(C_1 C_2^T) + \mathcal{A}V_m y_m^{gm}$ be the true GMRES residual after $m$ iterations of the inexact Arnoldi procedure. If for every $k \leqslant m$,*

$$\|E_k\| \leqslant \frac{\sigma_m(\underline{H}_m)}{m} \frac{1}{\|\widetilde{r}_{k-1}^{gm}\|}\varepsilon, \qquad (15)$$

*then $\|r_m^{gm} - \widetilde{r}_m^{gm}\| \leqslant \varepsilon$. Moreover, if*

$$\|E_k\| \leqslant \frac{1}{m\kappa(\underline{H}_m)} \frac{1}{\|\widetilde{r}_{k-1}^{gm}\|}\varepsilon, \qquad (16)$$

*then $\|(V_{m+1}\underline{H}_m)^T r_m^{gm}\| \leqslant \varepsilon$.*

*Similarly, if $r_m^{fom} := vec(C_1 C_2^T) + \mathcal{A}V_m y_m^{fom}$ is the true FOM residual after $m$ iterations of the inexact Arnoldi procedure, and if for every $k \leqslant m$,*

$$\|E_k\| \leqslant \frac{\sigma_m(H_m)}{m} \frac{1}{\|\widetilde{r}_{k-1}^{gm}\|}\varepsilon, \qquad (17)$$

*then $\|r_m^{fom} - \widetilde{r}_m^{fom}\| \leqslant \varepsilon$ and $\|V_m^T r_m^{fom}\| \leqslant \varepsilon$.*

Notice that the bound in (17) depends on the norm of the computed GMRES residual. This can be easily computed when Algorithm 2 is performed as $\|\widetilde{r}_m^{gm}\| = |e_{m+1}^T \underline{g}_{m+1}|$ in line 14 of Algorithm 2. However, if the FOM residual $\widetilde{r}_{k-1}^{fom}$ exists for every $k \leqslant m$, $\|\widetilde{r}_{k-1}^{gm}\|$ can be replaced by $\|\widetilde{r}_{k-1}^{fom}\|$ in (17).

The quantities involved in the estimates (15)-(16)-(17) are not available at iteration $k < m$ making the latter of theoretical interest only. To have practically usable truncation thresholds, the quantities in (15)-(16)-(17) must be approximated with computable values. Following the suggestions in [53], we can replace $m$ by the maximum number $m_{\max}$ of allowed iterations, $\sigma_{m_{\max}}(\underline{H}_{m_{\max}})$ is replaced by $\sigma_{n^2}(\mathcal{A})$, and we approximate $\sigma_1(\underline{H}_{m_{\max}})$ by $\sigma_1(\mathcal{A})$ when computing $\kappa(\underline{H}_{m_{\max}})$ in (16). The extreme singular values of $\mathcal{A}$ can be computed once and for all at the beginning of the iterative procedure, e.g., by the Lanczos method that must be carefully designed to avoid the construction of $\mathcal{A}$ and exploit its Kronecker structure. Approximations of $\sigma_1(\mathcal{A})$ and $\sigma_{n^2}(\mathcal{A})$ coming, e.g., from some particular features of the problem of interest, can also be employed. To conclude, we propose to use the following practical truncation thresholds $\varepsilon_{\mathcal{A}}^{(k)}$ in line 3 of Algorithm 2 in place of $\varepsilon_{\mathcal{A}}$:

$$\|E_k\| \leqslant \varepsilon_{\mathcal{A}}^{(k)} = \begin{cases} \frac{c_1}{m_{\max}} \frac{1}{\|\widetilde{r}_{k-1}^{gm}\|}\varepsilon, & c_1 \approx \sigma_{n^2}(\mathcal{A}), \\ \frac{1}{m_{\max}c_2} \frac{1}{\|\widetilde{r}_{k-1}^{gm}\|}\varepsilon, & c_2 \approx \kappa(\mathcal{A}), \end{cases} \qquad (18)$$

for LR-GMRES, and

$$\|E_k\| \leqslant \varepsilon_{\mathcal{A}}^{(k)} = \frac{c_1}{m_{\max}} \frac{1}{\|\widetilde{r}_{k-1}^{gm}\|} \varepsilon, \tag{19}$$

for LR-FOM.

Allowing $\|E_k\|$ to grow is remarkably important in our setting, especially for the memory requirements of the overall procedure. Indeed, if the truncation step in line 3 of Algorithm 2 is not performed, the rank of the basis vectors increases very quickly as, at the $m$-th iteration, we have

$$\operatorname{rank}(\mathcal{V}_{1,m} \mathcal{V}_{2,m}^T) \leqslant q p^m.$$

Therefore, at the first iterations the rank of the basis vectors is by construction low and having a very stringent tolerance in the computation of the low-rank approximations is not an issue. When the iterations proceed, the rank of the basis vectors increases but, at the same time, the increment in the thresholds for computing low-rank approximations of such vectors leads to more aggressive truncations with consequent remarkable gains in the memory allocation.

The interpretation of the truncation in line 3 of Algorithm 2 in terms of an inexact Krylov procedure has been already proposed in [18] for the more general case of GMRES applied to (2) where $\mathcal{A}$ is a tensor and the approximate solution is represented in the tensor-train (TT) format. However, also in the tensor setting, the results in Theorem 3.1 hold if and only if the matrix $V_m$ has orthonormal columns. In general, the low-rank truncation in line 7 can destroy the orthogonality of basis. In the next section we show that $V_m$ has orthogonal columns if the truncation step is performed in an appropriate way.

We first conclude this section with a couple of remarks.

**Remark 3.2.** *We have always assumed the initial guess $x_0 \in \mathbb{R}^n$ in (4) to be zero. This choice is motivated by the discussion in [53, Section 3], [34] where the authors show how this is a good habit in the framework of inexact Krylov methods.*

**Remark 3.3.** *Since*

$$\|r_m\| \leqslant \|\widetilde{r}_m\| + \sum_{j=1}^{m} \|E_j\| \cdot |e_j^T y_m| \leqslant \|\widetilde{r}_m\| + \sum_{j=1}^{m} \varepsilon_{\mathcal{A}}^{(j)} \cdot |e_j^T y_m|,$$

*where $\varepsilon_{\mathcal{A}}^{(j)}$ denotes one of the values in (18)-(19) depending on the selected procedure, the quantity $\|\widetilde{r}_m\| + \sum_{j=1}^{m} \varepsilon_{\mathcal{A}}^{(j)} \cdot |e_j^T y_m|$ must be computed to have a reliable stopping criterion in Algorithm 2. This means that the linear system $U_m y_m = g_m$ has to be solved at each iteration $m$. This does not significantly increase the computational workload because $U_m \in \mathbb{R}^{m \times m}$ is of small dimension and already given in triangular form.*

## 3.2 Structured perturbations of the basis

In this section we show how the low-rank truncations performed during the Gram-Schmidt procedure in line 7 of Algorithm 2 preserve the orthogonality of the basis, i.e., $V_m$ is still an orthonormal matrix, and the results presented in section 3.1 are still valid.

**Proposition 3.2.** *The matrix $V_{m+1} = [vec(\mathcal{V}_{1,1}\mathcal{V}_{2,1}^T), \ldots, vec(\mathcal{V}_{1,m+1}\mathcal{V}_{2,m+1}^T)] \in \mathbb{R}^{n^2 \times (m+1)}$ computed by performing $m$ iterations of Algorithm 2 has orthonormal columns if the low-rank truncations are computed by Algorithm 1.*

*Proof.* At the $m$-th iteration, the $(m+1)$-th basis vector is computed by performing (6) and then normalizing the result. In particular, if $\Theta_m = \operatorname{diag}(I_{\overline{s}}, -h_{1,m}I_{s_1}, \ldots, -h_{m,m}I_{s_m})$, then

$$(\widetilde{V}_1, \widetilde{V}_2) = \mathtt{trunc}([\overline{V}_1, \mathcal{V}_{1,1}, \ldots, \mathcal{V}_{1,m}], \Theta_m, [\overline{V}_2, \mathcal{V}_{2,1}, \ldots, \mathcal{V}_{2,m}], \varepsilon_{\mathtt{orth}}),$$

that is

$$\widetilde{V}_1 \widetilde{V}_2^T + F_{1,m} F_{2,m}^T = \overline{V}_1 \overline{V}_2^T - \sum_{j=1}^{m} h_{j,m} \mathcal{V}_{1,j} \mathcal{V}_{2,j}^T,$$

9

where $F_{1,m}F_{2,m}^T$ is the matrix discarded during the application of Algorithm 1.

If $Q_1R_1 = [\overline{V}_1, \mathcal{V}_{1,1}, \dots, \mathcal{V}_{1,m}]$, $Q_2R_2 = [\overline{V}_2, \mathcal{V}_{2,1}, \dots, \mathcal{V}_{2,m}]$ denote the skinny QR factorizations performed during $\texttt{trunc}$ and $U\Sigma W^T = R_1\Theta_m R_2^T$, $U = [u_1, \dots, u_{\mathfrak{s}_m}]$, $W = [w_1, \dots, w_{\mathfrak{s}_m}]$, $\Sigma = \mathrm{diag}(\sigma_1, \dots, \sigma_{\mathfrak{s}_m})$, $\mathfrak{s}_m := \overline{s} + \sum_{j=1}^m s_j$, is the SVD decomposition, then

$$\widetilde{V}_1 = Q_1\left([u_1, \dots, u_{k_m}]\sqrt{\mathrm{diag}(\sigma_1, \dots, \sigma_{k_m})}\right), \quad \widetilde{V}_2 = Q_2\left([w_1, \dots, w_{k_m}]\sqrt{\mathrm{diag}(\sigma_1, \dots, \sigma_{k_m})}\right),$$

and

$$F_{1,m} = Q_1\left([u_{k_m+1}, \dots, u_{\mathfrak{s}_m}]\sqrt{\mathrm{diag}(\sigma_{k_m+1}, \dots, \sigma_{\mathfrak{s}_m})}\right), \ F_{2,m} = Q_2\left([w_{k_m+1}, \dots, w_{\mathfrak{s}_m}]\sqrt{\mathrm{diag}(\sigma_{k_m+1}, \dots, \sigma_{\mathfrak{s}_m})}\right),$$

where $k_m$ is the smallest index such that $\sqrt{\sum_{i=k_m+1}^{\mathfrak{s}_m}\sigma_i} \leqslant \varepsilon_{\texttt{orth}} \cdot \|\Sigma\|$.

By construction we have

$$\langle \mathcal{V}_{1,j}\mathcal{V}_{2,j}^T, \widetilde{V}_1\widetilde{V}_2^T + F_{1,m}F_{2,m}^T\rangle_F = \langle \mathcal{V}_{1,j}\mathcal{V}_{2,j}^T, Q_1(U\Sigma W^T)Q_2^T\rangle_F = 0, \quad \text{for all } j = 1, \dots, m,$$

and this means that $\mathcal{V}_{1,j}\mathcal{V}_{2,j}^T$ is orthogonal to the tensor space $\mathrm{Range}(Q_2W\sqrt{\Sigma}) \otimes \mathrm{Range}(Q_1U\sqrt{\Sigma})$, and therefore to any matrix of the form $Q_1U\sqrt{\Sigma}\Xi\sqrt{\Sigma}W^TQ_2^T$, $\Xi \in \mathbb{R}^{\mathfrak{s}_m \times \mathfrak{s}_m}$, for all $j = 1, \dots, m$. Since $\widetilde{V}_1\widetilde{V}_2^T = Q_1U\sqrt{\Sigma}[e_1, \dots, e_{k_m}][e_1, \dots, e_{k_m}]^T\sqrt{\Sigma}W^TQ_2^T$, it holds

$$\langle \mathcal{V}_{1,j}\mathcal{V}_{2,j}^T, \widetilde{V}_1\widetilde{V}_2^T\rangle_F = 0, \quad \text{for all } j = 1, \dots, m.$$

To conclude, $\mathcal{V}_{1,m+1}\mathcal{V}_{2,m+1}^T = \widetilde{V}_1\widetilde{V}_2^T/\|\widetilde{V}_1\widetilde{V}_2^T\|$ so that $\mathrm{vec}(\mathcal{V}_{1,m+1}\mathcal{V}_{2,m+1}^T)$ has unit norm. $\qquad\square$

As shown in the proof of Proposition 3.2, to maintain the orthogonality of the basis, it is crucial that $\widetilde{V}_1\widetilde{V}_2^T$ belongs to the space defined by $[\overline{V}_1, \mathcal{V}_{1,1}, \dots, \mathcal{V}_{1,m}]\Theta_m[\overline{V}_2, \mathcal{V}_{2,1}, \dots, \mathcal{V}_{2,m}]^T$ and this is possible thanks to the QR-SVD-based truncation we perform. In general, it may happen that the computed basis $V_m$ is no longer orthogonal if different truncation strategies are adopted.

Clearly, the truncations performed during the orthogonalization procedure consist in another source of inexactness that must be taken into account. The inexact Arnoldi relation (13) becomes

$$\mathcal{A}V_m - [\mathrm{vec}(E_1), \dots, \mathrm{vec}(E_m)] = V_mH_m + h_{m+1,m}\mathrm{vec}(\mathcal{V}_{1,m+1}\mathcal{V}_{2,m+1}^T)e_m^T + [\mathrm{vec}(F_{1,1}F_{2,1}^T), \dots, \mathrm{vec}(F_{1,m}F_{2,m}^T)],$$

and one can derive results similar to the ones in Theorem 3.1 for the inexact Arnoldi relation

$$\mathcal{A}V_m - [\mathrm{vec}(E_1 + F_{1,1}F_{2,1}^T), \dots, \mathrm{vec}(E_m + F_{1,m}F_{2,m}^T)] = V_mH_m + h_{m+1,m}\mathrm{vec}(\mathcal{V}_{1,m+1}\mathcal{V}_{2,m+1}^T)e_m^T,$$

obtaining estimates for $\|E_k + F_{1,k}F_{2,k}^T\|$. Since

$$\|E_k + F_{1,k}F_{2,k}^T\| \leqslant \|E_k\| + \|F_{1,k}F_{2,k}^T\|,$$

it may be interesting to study how to distribute the allowed inexactness between the truncation steps.

Since the rank of the iterates grows less dramatically during the orthogonalization step compared to what happens after the multiplication with $\mathcal{A}$, we allow $2\|E_k\|$ to grow in accordance with Theorem 3.1, while $\|F_{1,k}F_{2,k}^T\|$ is maintained sufficiently small. Indeed, the matrix $[\overline{V}_1, \mathcal{V}_{1,1}, \dots, \mathcal{V}_{1,m}]\Theta_m[\overline{V}_2, \mathcal{V}_{2,1} \dots, \mathcal{V}_{2,m}]^T$ in line 7 of Algorithm 2 is, in general, very rank-deficient and a significant reduction in the number of columns to be stored takes place even when the $\texttt{trunc}$ function is applied with a small threshold.

In particular, at the $m$-th iteration, we can set

$$\varepsilon_{\texttt{orth}} = \min\{\|E_k\|, \varepsilon/(m_{\max})\}, \tag{20}$$

where $\varepsilon$ is the desired accuracy of the final solution in terms of relative residual norm. This means that $\|E_k + F_{1,k}F_{2,k}^T\|$ fulfills the estimates in (15)-(16)-(17) and the convergence is thus preserved.

The vectors $y_m^{fom}$ and $y_m^{gm}$ can be still computed as in (9)-(11) and Proposition 3.1 holds also when the low-rank truncation in line 7 of Algorithm 2 are performed.

**Proposition 3.3.** *Let* (3.2) *hold and define* $W_m = \mathcal{A}V_m - [vec(E_1), \ldots, vec(E_m)]$. *If* $y_m^{gm}$ *is computed as in* (11), *where* $\underline{H}_m$ *stems from the low-rank Arnoldi procedure illustrated in Algorithm 2 with low-rank truncations are performed by Algorithm 1, then* $q_m^{gm} := W_m y_m^{gm}$ *is such that*

$$q_m = \underset{q \in Range(W_m)}{\mathrm{argmin}} \|vec(C_1 C_2^T) + q\|.$$

*Similarly, if* $y_m^{fom}$ *is computed as in* (9) *where* $H_m$ *is the principal square submatrix of the aforementioned* $\underline{H}_m$, *then* $q_m^{fom} := W_m y_m^{fom}$ *is such that*

$$vec(C_1 C_2^T) + q_m \perp Range(V_m).$$

*Proof.* We only need to prove that $V_m^T[vec(F_{1,1} F_{2,1}^T), \ldots, vec(F_{1,m} F_{2,m}^T)] = 0$ as the rest of the proof comes from [53, Proposition 3.2-3.3].

Using the same arguments of the proof of Proposition 3.2, we can show that $F_{1,j} F_{2,j}^T$ is orthogonal to $\mathcal{V}_{1,i} \mathcal{V}_{2,i}^T$ for all $j, i = 1, \ldots, m$, $i + 1 \neq j$. Therefore, the only nonzero components of

$$V_m^T[vec(F_{1,1} F_{2,1}^T), \ldots, vec(F_{1,m} F_{2,m}^T)],$$

are in the first subdiagonal. These entries are of the form $\langle \mathcal{V}_{1,\ell+1} \mathcal{V}_{2,\ell+1}^T, F_{1,\ell} F_{2,\ell}^T \rangle_F$ and we show they are zero for every $\ell = 1, \ldots, m-1$. We have $\mathcal{V}_{1,\ell+1} \mathcal{V}_{2,\ell+1}^T = \widetilde{V}_1 \widetilde{V}_2^T / \|\widetilde{V}_1 \widetilde{V}_2^T\|$ and following the proof of Proposition 3.2, $\widetilde{V}_1$, $\widetilde{V}_2$, $F_{1,\ell}$ and $F_{2,\ell}$ can be written as

$$\widetilde{V}_1 = Q_1 \left( [u_1, \ldots, u_{k_\ell}] \sqrt{\mathrm{diag}(\sigma_1, \ldots, \sigma_{k_\ell})} \right), \quad \widetilde{V}_2 = Q_2 \left( [w_1, \ldots, w_{k_\ell}] \sqrt{\mathrm{diag}(\sigma_1, \ldots, \sigma_{k_\ell})} \right),$$

and

$$F_{1,\ell} = Q_1 \left( [u_{k_\ell+1}, \ldots, u_{\mathfrak{s}_\ell}] \sqrt{\mathrm{diag}(\sigma_{k_\ell+1}, \ldots, \sigma_{\mathfrak{s}_\ell})} \right), \ F_{2,\ell} = Q_2 \left( [w_{k_\ell+1}, \ldots, w_{\mathfrak{s}_\ell}] \sqrt{\mathrm{diag}(\sigma_{k_\ell+1}, \ldots, \sigma_{\mathfrak{s}_\ell})} \right),$$

where $Q_1 R_1 = [\overline{V}_1, \mathcal{V}_{1,1}, \ldots \mathcal{V}_{1,\ell}]$, $Q_2 R_2 = [\overline{V}_2, \mathcal{V}_{2,1}, \ldots \mathcal{V}_{2,\ell}]$ are skinny QR factorizations and the SVD decomposition is given by $U \Sigma W^T = R_1 \Theta_\ell R_2^T$, $U = [u_1, \ldots, u_{\mathfrak{s}_\ell}]$, $W = [w_1, \ldots, w_{\mathfrak{s}_\ell}]$, $\Sigma = \mathrm{diag}(\sigma_1, \ldots, \sigma_{\mathfrak{s}_\ell})$, $\mathfrak{s}_\ell = \mathrm{rank}(\overline{V}_1) + \sum_{i=1}^{\ell} s_i$, and $k_\ell$ is the smallest index such that $\sqrt{\sum_{i=k_\ell+1}^{\mathfrak{s}_\ell} \sigma_i} \leqslant \varepsilon_{\mathtt{orth}} \cdot \|\Sigma\|$.

Since $Q_1$, $Q_2$, $W$ and $U$ are orthogonal matrices, we have

$$\begin{aligned}
\langle \mathcal{V}_{1,\ell+1} \mathcal{V}_{2,\ell+1}^T, F_{1,\ell} F_{2,\ell}^T \rangle_F =& \mathrm{trace} \big( \mathrm{diag}(\sigma_1, \ldots, \sigma_{k_\ell}) [u_1, \ldots, u_{k_\ell}]^T [u_{k_\ell+1}, \ldots, u_{\mathfrak{s}_\ell}] \\
& \mathrm{diag}(\sigma_{k_\ell+1}, \ldots, \sigma_{\mathfrak{s}_\ell}) [w_{k_\ell+1}, \ldots, w_{\mathfrak{s}_\ell}]^T [w_1, \ldots, w_{k_\ell}] \big) \\
=& 0,
\end{aligned}$$

and we get the result. $\qquad\square$

The true relative residual norm can be written as

$$r_m = \widetilde{r}_m - [vec(E_1 + F_{1,1} F_{2,1}^T), \ldots, vec(E_m + F_{1,m} F_{2,m}^T)] y_m,$$

and following the discussion in Remark 3.3 we have

$$\|r_m\| \leqslant \|\widetilde{r}_m\| + \sum_{j=1}^{m} \|E_j\| \cdot |e_j^T y_m| + \sum_{j=1}^{m} \|F_{1,j} F_{2,j}^T\| \cdot |e_j^T y_m| \leqslant \|\widetilde{r}_m\| + \sum_{j=1}^{m} \left( \varepsilon_{\mathcal{A}}^{(j)} + \frac{m}{m_{\max}} \varepsilon \right) |e_j^T y_m|, \quad (21)$$

so that the right-hand side in the above expression must be computed to check convergence.

# 4 Alternative truncation strategies

As we discussed above, to keep the low-rank Krylov methods computationally feasible, the quantities involved in the solution process have to be compressed so that their rank, i.e., the sizes of the low-rank factors, is kept small. Let $NML^T$ with factors $N, L \in \mathbb{R}^{n \times m}$, $M \in \mathbb{R}^{m \times m}$, be the quantity to be compressed, and assume that $\text{rank}(NML^T) = m$. So far we have used a direct approach using QR and SVD decompositions in Algorithm 1 which essentially computes a partial SVD of $NML^T$ corresponding to all $m$ nonzero singular values. This whole procedure relies heavily on dense linear algebra computations and can, hence, become quite expensive. This is especially due to the QR decompositions which will be expensive if the rectangular factors $N, L$ have many columns. Moreover, if $NML^T$ has a very small numerical numerical rank, say $k \ll m$, then Algorithm 1 will generate a substantial computational overhead because $m - k$ singular vectors will be thrown away. Nevertheless, thanks to the complete knowledge of all singular values, this procedure is able to correctly assess the truncation error in the Frobenius norm so that the required accuracy of the truncation is always met.

Following the discussion in, e.g., [12, 38, 56], a more economical alternative could be to compute only a partial SVD $NML^T \approx U_k \Sigma_k W_k^T$ associated to the $k$ singular values that are larger than the given truncation threshold. If also the $(k + 1)$-th singular value is computed, one has the truncation error in the 2-norm: $\|NML^T - U_k \Sigma_k W_k^T\|_2 \leqslant \sigma_{k+1}(NML^T)$. Obviously, the results of the previous section are still valid if this form of truncation is used. Approximations of the dominant singular values and corresponding singular vectors can be computed by iterative methods for large-scale SVD computations as, e.g., Lanczos bidiagonalization (see, e.g., [3, 32, 55]) or Jacobi-Davidson methods; see [25]. To apply these methods, only matrix vector products $N(M(L^T x))$ and $L(M^T(N^T x))$ are required. For achieving the compression goal one could, e.g., compute $k_{\max} \geqslant k$ triplets and, if required, neglect any singular vectors corresponding to singular value below a certain threshold. However, we do in general not know in advance how many singular values will be larger than a given threshold. Picking a too small value of $k_{\max}$ can lead to very inaccurate truncations that do not satisfy the required thresholds (15)–(17), (20) and, therefore, endanger the convergence of the low-rank Krylov method. Some of aforementioned iterative SVD methods converge theoretically monotonically, i.e., the singular values are found in a decreasing sequence starting with the largest one. Hence, the singular value finding iteration can be kept running until a sufficiently small singular value approximation, e.g., $\widetilde{\sigma} < \varepsilon_{\text{trunc}} \|NML^T\|_2$, is detected. In the practical situations within low-rank Krylov methods, the necessary number of singular triplets can be $\mathcal{O}(10^2)$ or larger and it may be difficult to ensure that the iterative SVD algorithms do not miss some of the largest singular values or that no singular values are detected several times. Due to the sheer number of occurrences where compression is required in Algorithm 2, preliminary tests with iterative SVD methods did not yield any substantial savings compared to the standard approach in Algorithm 1.

Compression algorithms based on randomized linear algebra might offer further alternative approaches with reduced computational times. See, e.g., [15, 23, 27].

# 5 Preconditioning

It is well-known that Krylov methods require preconditioning in order to obtain a fast convergence in terms of number of iterations and low-rank Krylov methods are no exception. However, due to the peculiarity of our framework, the preconditioner operator must possess some supplementary features with respect to standard preconditioners for linear systems. Indeed, in addition to be effective in reducing the number of iterations at a reasonable computational cost, the preconditioner operator must not dramatically increase the memory requirements of the solution process.

Given a nonsingular operator $\mathcal{P}$ or its inverse $\mathcal{P}^{-1}$, if we employ right preconditioning, the original systems (2) is transformed into

$$\mathcal{A}\mathcal{P}^{-1}\overline{x} = -\text{vec}(C_1 C_2^T), \quad \text{vec}(X) = \mathcal{P}^{-1}\overline{x}, \tag{22}$$

so that, at each iteration $m$, we have to apply $\mathcal{P}^{-1}$ to the current basis vector $\text{vec}(\mathcal{V}_{1,m}\mathcal{V}_{2,m}^T)$. Note that we restrict ourselves here to right preconditioning because this has the advantage that one can still monitor the true unpreconditioned residuals without extra work within the Krylov routine. Of course, in principle also left and two-sided preconditioning can be used.

The preconditioning operation must be able to exploit the low-rank format of $\mathcal{V}_{1,m}\mathcal{V}_{2,m}^T$. Therefore, a naive operation of the form $\mathcal{P}^{-1}\text{vec}(\mathcal{V}_{1,m}\mathcal{V}_{2,m}^T)$ is not admissible in our context as this would require the allocation of the dense $n \times n$ matrix $\mathcal{V}_{1,m}\mathcal{V}_{2,m}^T$. One way to overcome this numerical difficulty is to employ a preconditioner operator $\mathcal{P}$ which allows for a representation in terms of a Kronecker sum, namely

$$\mathcal{P} = \sum_{i=1}^{\ell} P_i \otimes T_i. \tag{23}$$

This means that the operation $z_m = \mathcal{P}^{-1}\text{vec}(\mathcal{V}_{1,m}\mathcal{V}_{2,m}^T)$ is equivalent to solving the matrix equation

$$\sum_{i=1}^{\ell} T_i Y_m P_i^T - \mathcal{V}_{1,m}\mathcal{V}_{2,m}^T = 0, \quad \text{vec}(Y_m) = z_m. \tag{24}$$

In our setting, the operator $\mathcal{P}$ often amounts to an approximation to $\mathcal{A}$ in (2) obtained by either dropping some terms in the series or replacing some of them by a multiple of the identity. See, e.g., [40,43,57]. Another option that has not been fully explored in the matrix equation literature so far is the case of polynomial preconditioners (see, e.g., [35, 61]) where $\mathcal{P}^{-1}$ resembles a fixed low-degree polynomial evaluated in $\mathcal{A}$. Alternatively, we can formally set $\mathcal{P} = \mathcal{A}$ in (23) and inexactly solve equation (24) by few iterations of another Krylov method (e.g., Algorithm 2) leading to an inner-outer Krylov method; see, e.g., [52].

Clearly, equation (24) must be easy to solve. For instance, if $\ell = 1$, then $Y_m = (T_1^{-1}\mathcal{V}_{1,m})(P_1^{-1}\mathcal{V}_{2,m})^T$ and an exact application of the preconditioner can be carried out. Similarly, when $\ell = 2$ and a fixed number of ADI iterations are performed at each Krylov iteration $m$, then it is easy to show that we are still working in an exact preconditioning framework. See, e.g. [8, 16]. In all these cases, the results presented in the previous sections still hold provided $\mathcal{A}$ is replaced by the preconditioned matrix $\mathcal{A}\mathcal{P}^{-1}$.

Equation (24) is often iteratively solved and, in general, this procedure leads to the computation of a low-rank approximation $\mathcal{Z}_{1,m}\mathcal{Z}_{2,m}^T$ to $Y_m$ that has to be interpreted as a variable preconditioning scheme with a different preconditioning operator at each outer iteration. In this cases, a flexible variant of Algorithm 2 must be employed which consists in a standard flexible Krylov procedure equipped with the low-rank truncations presented in the previous sections. See, e.g., [54, Section 10] for some details about flexible Krylov methods and [45, 46, Section 9.4.1] for a discussion about flexible GMRES.

We must mention that the employment of a flexible procedure doubles, at least, the memory requirements of the solution process. Indeed, both the *preconditioned* and *unpreconditioned* bases must be stored and $\text{rank}(\mathcal{Z}_{1,m}\mathcal{Z}_{2,m}^T) \geqslant \text{rank}(\mathcal{V}_{1,m}\mathcal{V}_{2,m}^T)$ for all $m$. This aspect must be taken into account when designing the preconditioner. See Example 29.

At a first glance, the presence of a variable preconditioning procedure can complicate the derivations illustrated in sections 3.1-3.2 for the safe selection of the low-rank truncation thresholds that guarantee the convergence of the solution method. Indeed, if at iteration $m$, $\mathcal{Z}_{1,m}\mathcal{Z}_{2,m}^T$ is the result of the preconditioning step (24), we still want to truncate the matrix $[A_1\mathcal{Z}_{1,m}, \ldots, A_p\mathcal{Z}_{1,m}][B_1\mathcal{Z}_{2,m}, \ldots, B_p\mathcal{Z}_{2,m}]^T$ in order to moderate the storage demand and one may wonder if the inexactness of step (24) plays a role in such a truncation. Thanks to the employment of a flexible strategy, we are going to show how the tolerances for the low-rank truncations, namely $\varepsilon_{\mathcal{A}}$ and $\varepsilon_{\text{orth}}$ in Algorithm 2, can be still computed as illustrated in sections 3.1-3.2.

Flexible Krylov methods are characterized not only by having a preconditioner that changes at each iteration, but also from the fact that the solution is recovered by means of the preconditioned basis. In particular,

$$\text{vec}(X_m) = Z_m y_m, \quad Z_m := [\text{vec}(\mathcal{Z}_{1,1}\mathcal{Z}_{2,1}^T), \ldots, \text{vec}(\mathcal{Z}_{1,m}\mathcal{Z}_{2,m}^T)],$$

see, e.g., [45]; this is a key ingredient in our analysis.

We start our discussion by considering flexible Krylov methods with no truncations. For this class of solvers the relation

$$\mathcal{A}Z_m = V_m H_m + h_{m+1,m}\text{vec}(\mathcal{V}_{1,m+1}\mathcal{V}_{2,m+1}^T)e_m^T, \tag{25}$$

holds, see, e.g., [46, Equation (9.22)], and $\text{span}\{\text{vec}(\mathcal{Z}_{1,1}\mathcal{Z}_{2,1}^T), \ldots, \text{vec}(\mathcal{Z}_{1,m}\mathcal{Z}_{2,m}^T)\}$ is not a Krylov subspace in general. Therefore, also for the flexible Krylov methods with no low-rank truncations we must consider

13

constrains different from the ones in (7)-(10) and results similar to the ones in Proposition 3.1 with $W_m = \mathcal{A}Z_m$ hold. See, e.g., [46, Proposition 9.2].

If we now introduce a low-rank truncation of the matrix

$$[A_1 \mathcal{Z}_{1,m}, \ldots, A_p \mathcal{Z}_{1,m}][B_1 \mathcal{Z}_{2,m}, \ldots, B_p \mathcal{Z}_{2,m}]^T,$$

at each iteration $m$, that is we compute

$$(\overline{V}_1, \overline{V}_2) = \texttt{trunc}([A_1 \mathcal{Z}_{1,m}, \ldots, A_p \mathcal{Z}_{1,m}], I, [B_1 \mathcal{Z}_{2,m}, \ldots, B_p \mathcal{Z}_{2,m}], \varepsilon_{\mathcal{A}}), \qquad (26)$$

then the relation (25) becomes

$$\mathcal{A}Z_m - [\text{vec}(E_1), \ldots, \text{vec}(E_m)] = V_m H_m + h_{m+1,m} \text{vec}(\mathcal{V}_{1,m+1} \mathcal{V}_{2,m+1}^T) e_m^T, \qquad (27)$$

where the matrices $E_k$'s are the ones discarded when (26) is performed. If $\|E_k\|$ satisfies the inequalities in Theorem 3.1, then the convergence of the low-rank flexible Krylov procedure is still guaranteed in the sense that the residual norm keeps decreasing as long as span$\{\text{vec}(\mathcal{Z}_{1,1} \mathcal{Z}_{2,1}^T), \ldots, \text{vec}(\mathcal{Z}_{1,m} \mathcal{Z}_{2,m}^T)\}$ grows. However, the matrix $H_m$ no longer represents an approximation of $\mathcal{A}$ onto the current subspace and the approximation of $\sigma_{m_{\max}}(\underline{H}_{m_{\max}})$ and $\sigma_1(\underline{H}_{m_{\max}})$ in the right-hand side of (15)-(16)-(17) by the corresponding singular values of $\mathcal{A}$ may no longer be effective. In our numerical experience, approximating $\sigma_{m_{\max}}(\underline{H}_{m_{\max}})$ and $\sigma_1(\underline{H}_{m_{\max}})$ by the smallest and largest singular values of the preconditioned matrix $\mathcal{A}\mathcal{P}^{-1}$, i.e., mimicking what is done in case of exact applications of $\mathcal{P}$, provides satisfactory results. Obtaining computable approximations to $\sigma_{m_{\max}}(\underline{H}_{m_{\max}})$ and $\sigma_1(\underline{H}_{m_{\max}})$ for the inner-outer approach is not straightforward. In this case, a practical approach may be to still approximate $\sigma_{m_{\max}}(\underline{H}_{m_{\max}})$ and $\sigma_1(\underline{H}_{m_{\max}})$ by $\sigma_{n^2}(\mathcal{A})$ and $\sigma_1(\mathcal{A})$, respectively. These approximations may be very rough as they completely neglect the role of the preconditioner so that they may lead to quite conservative truncation thresholds. However, at the moment, we do not see any another possible alternatives.

The introduction of the low-rank truncations that lead to (27) implies that the constrained imposed on the residual vector are no longer in terms of the space spanned by $Z_m$ and the results presented in Proposition 3.1 with $W_m = \mathcal{A}Z_m - [\text{vec}(E_1), \ldots, \text{vec}(E_m)]$ hold.

In flexible Krylov methods, the orthogonalization procedure involves only the unpreconditioned basis $V_m$ so that the truncation step in line 7 of Algorithm 2 is not really affected by the preconditioning procedure and the results in Proposition 3.2-3.3 are still valid. The truncation threshold $\varepsilon_{\texttt{orth}}$ can be still selected as proposed in section 3.2.

# 6  Short recurrence methods

Short recurrence Krylov methods can be very appealing in our context as only a fixed, usually small, number of basis vectors have to be stored. In case of symmetric problems, i.e., equation (1) where all the coefficient matrices $A_i$'s and $B_i$'s are symmetric, the low-rank MINRES algorithm proposed in [39] can be employed in the solution process.

If $\mathcal{A}$ in (2) is also positive definite, the low-rank CG method illustrated in [24] is a valid candidate for the solution of equation (1). Notice that, in general, it is not easy to characterize the spectral distribution of $\mathcal{A}$ in terms of the spectrum of the coefficient matrices $A_i$'s and $B_i$'s. However, it can be shown that if $A_i$ and $B_i$ are positive definite for all $i$, then also $\mathcal{A}$ is positive definite.

Short recurrence methods can be appealing also in case of a nonsymmetric $\mathcal{A}$ and low-rank variants of BICGSTAB ( [59]), QMR ( [21]) or other methods can be employed to solve equation (1).

See, e.g., [8,56] for an implementation of low-rank MINRES, CG and BICGSTAB.

In all the short recurrence Krylov methods, the constructed basis $V_m$ is not orthogonal in practice and this loss of orthogonality must be taken into account in the bounds for the allowed inexactness proposed in Theorem 3.1. In [53, Section 6], the authors propose to incorporate the smallest singular values of the computed basis, namely $\sigma_m(V_m)$, in the right-hand side of (15)-(16)-(17) to guarantee the convergence of the method. However, no practical approximation to $\sigma_m(V_m)$ is proposed in [53].

A different approach that can be pursued is the one illustrated in [13]. In this paper the authors propose to select bounds of the form

$$\|E_k\| \leqslant \min\{\alpha_k\varepsilon, 1\}, \quad \alpha_k = \frac{1}{\min\{\|\widetilde{r}_k\|, 1\}}, \tag{28}$$

where $\widetilde{r}_k$ is the current computed residual vector, and in [58] the authors studied the effects of such a choice on the convergence of a certain class of inexact Krylov methods. In particular, in [58] it is shown how the residual gap $\delta_m$ remains small if $\|E_k\|$ fulfills (28) for all $k \leqslant m$. Even though the true residual and the computed one are close, this does not imply that the residual norm is actually always small and we thus have to assume that the norm of the computed residual goes to zero as it is done in [58].

# 7  Numerical examples

In this section we present some numerical results that confirm the theoretical analysis derived in the previous sections. To this end we consider some general multiterm linear matrix equation of the form (1) stemming from the discretization of certain deterministic and stochastic PDEs.

We apply the LR-GMRES variant of Algorithm 2 in the solution process and we always select Algorithm 1 for the low-rank truncations.

We report the number of performed iterations, the rank of the computed solution, the computational time needed to calculate such a solution together with the relative residual norm achieved, and the storage demand. For the latter, we document the number of columns $\mathfrak{s} = \sum_{j=1}^{m+1} s_j$ of the matrix $[\mathcal{V}_{1,1}, \ldots, \mathcal{V}_{1,m+1}]$, where $m$ is the number of iterations needed to converge. Similarly, if a flexible strategy is adopted, we also report the number of columns $\mathfrak{z}$ of $[\mathcal{Z}_{1,1}, \ldots, \mathcal{Z}_{1,m}]$.

This means that, for equations of the form (1) where $n_A = n_B = n$, we have to allocate $2\mathfrak{s}$ ($2(\mathfrak{s} + \mathfrak{z})$) vectors of length $n$. If $n_A \neq n_B$, the memory requirements amount to $\mathfrak{s}$ ($\mathfrak{s} + \mathfrak{z}$) vectors of length $n_A$ and $\mathfrak{s}$ ($\mathfrak{s} + \mathfrak{z}$) vectors of length $n_B$.

The solution process is stopped as soon as the upper bound on the residual norm in (21), normalized by $\|C_1 C_2^T\|_F$, gets smaller than $10^{-6}$.

As already mentioned, we always assume that the exact solution $X$ admits accurate low-rank approximations. Nevertheless, if $S_1, S_2$ are the low-rank factors computed by Algorithm 2, we report also the real relative residual norm $\|\sum_{i=1}^{p} A_i S_1 S_2^T B_i^T + C_1 C_2^T\|_F / \|C_1 C_2^T\|_F$ in the following to confirm the reliability of our numerical procedure. Once again, the real residual norm can be computed at low cost by exploiting the low rank of $S_1 S_2^T$ and the cyclic property of the trace operator.

All results were obtained with Matlab R2017b ( [37]) on a Dell machine with 2.4GHz processors and 250 GB of RAM.

**Example 7.1.** We consider a slight modification of Example 4 in [40]. In particular, the continuous problem we have in mind is the convection-diffusion equation

$$\begin{aligned} -\nu\Delta u + \vec{w} \cdot \nabla u &= 1, \quad \text{in } D = (0,1)^2, \\ u &= 0, \quad \text{on } \partial D, \end{aligned} \tag{29}$$

where $\nu > 0$ is the viscosity parameter and the convection vector $\vec{w}$ is given by $\vec{w} = (\phi_1(x)\psi_1(y), \phi_2(x)\psi_2(y)) = ((1 - (2x+1)^2)y, -2(2x+1)(1-y^2))$. The centered finite differences discretization of equation (29) yields the following matrix equation

$$\nu TX + \nu XT + \Phi_1 BX\Psi_1 + \Phi_2 XB^T\Psi_2 - \mathbf{1}\mathbf{1}^T = 0, \tag{30}$$

where $T \in \mathbb{R}^{n \times n}$ is the negative discrete laplacian, $B \in \mathbb{R}^{n \times n}$ corresponds to the discretization of the first derivative, $\Phi_i$ and $\Psi_i$ are diagonal matrices collecting the nodal values of the corresponding functions $\phi_i$, $\psi_i$, $i = 1, 2$, and $\mathbf{1} \in \mathbb{R}^n$ is the vector of all ones. See [40] for more details.

Even though equation (30) amounts to a generalized Sylvester equation, the solution schemes available in the literature and tailored to this kind of problems cannot be applied to equation (30) in general. Indeed, to the best of our knowledge, all the existing methods for large-scale generalized equations rely on a splitting of the overall discrete operator of the form $\mathcal{M} + \mathcal{N}$, $\mathcal{M}(X) = \nu TX + \nu XT$, $\mathcal{N}(X) = \Phi_1 BX\Psi_1 + \Phi_2 XB^T\Psi_2$,

which is supposed to be convergent. See, e.g., [8, 26, 48]. However, the latter property may be difficult to meet in case of the convection-diffusion equation, especially for dominant convection.

We thus have to interpret (30) as a general multiterm matrix equation of the form (1) and we solve it by the preconditioned LR-GMRES. Following the discussion in [40], we use the operator

$$\mathcal{L}: \quad \begin{array}{ccc} \mathbb{R}^{n \times n} & \to & \mathbb{R}^{n \times n} \\ X & \mapsto & (\nu T + \overline{\psi}_1 \Psi_1 B) X + X(\nu T + \overline{\phi}_2 B^T \Psi_2), \end{array}$$

as preconditioner, where $\overline{\psi}_1, \overline{\phi}_2 \in \mathbb{R}$ are the mean values of $\psi_1(y)$ and $\phi_2(x)$ on $(0, 1)$, respectively.

At each LR-GMRES iteration, we approximately invert $\mathcal{L}$ by performing 10 iterations of the extended Krylov subspace method for Sylvester equation[4] derived in [14]. Since this scheme gives a different preconditioner every time it is called, we must employ the flexible variant of LR-GMRES. To avoid an excessive increment in the memory requirements due to the allocation of both the preconditioned and unpreconditioned bases, we do not apply $\mathcal{L}$ to the current basis vector, i.e., at iteration $k$, we do not compute $\mathcal{Z}_{1,k} \mathcal{Z}_{2,k}^T \approx \mathcal{L}^{-1}(\mathcal{V}_{1,k} \mathcal{V}_{2,k}^T)$. We first truncate the low-rank factors $\mathcal{V}_{1,k}, \mathcal{V}_{2,k}$, namely we compute $(\widehat{\mathcal{V}}_{1,k}, \widehat{\mathcal{V}}_{2,k}) = \text{trunc}(\mathcal{V}_{1,k}, I, \mathcal{V}_{2,k}, \varepsilon_{\texttt{precond}})$, and then define $\mathcal{Z}_{1,k}, \mathcal{Z}_{2,k}$ such that $\mathcal{Z}_{1,k} \mathcal{Z}_{2,k}^T \approx \mathcal{L}^{-1}(\widehat{\mathcal{V}}_{1,k} \widehat{\mathcal{V}}_{2,k}^T)$. This procedure leads to a lower storage demand of the overall solution process and to less time consuming preconditioning steps. On the other hand, the effectiveness of the preconditioner in reducing the total iteration count may get weakened, especially for large $\varepsilon_{\texttt{precond}}$. In the results reported in the following we have always set $\varepsilon_{\texttt{precond}} = 10^{-3}$.

In Table 1 we report the results for different values of $n$ and $\nu$.

Table 1: Example 7.1. Results for different values of $n$ and $\nu$.

| $\nu$ | $n$ | It. | rank($S_1 S_2^T$) | Time (s) | Memory | | Conv. Checks | |
| | | | | | $V_m$ | $Z_m$ | (21)/$\|C_1 C_2^T\|_F$ | Real Res. |
|---|---|---|---|---|---|---|---|---|
| 0.5 | 5000 | 8 | 58 | 2.872e1 | 1174 | 915 | 4.078e-7 | 2.974e-7 |
| | 10000 | 8 | 59 | 8.352e1 | 1543 | 1079 | 4.242e-7 | 3.144e-7 |
| | 15000 | 8 | 69 | 1.812e2 | 2075 | 1239 | 9.492e-7 | 6.401e-7 |
| 0.1 | 5000 | 15 | 66 | 1.256e2 | 3284 | 1880 | 7.803e-7 | 4.509e-7 |
| | 10000 | 15 | 71 | 4.687e2 | 4566 | 2364 | 7.798e-7 | 4.497e-7 |
| | 15000 | 15 | 81 | 1.169e3 | 6152 | 2800 | 8.623e-7 | 4.519e-7 |
| 0.05 | 5000 | 20 | 77 | 4.067e2 | 5957 | 2980 | 8.533e-7 | 2.644e-7 |
| | 10000 | 20 | 82 | 1.486e3 | 7896 | 3624 | 8.558e-7 | 2.640e-7 |
| | 15000 | 20 | 88 | 3.467e3 | 9867 | 4093 | 8.691e-7 | 2.656e-7 |

We notice that the number of iterations is very robust with respect to the problem dimension $n$, and thus the mesh-size. Unfortunately, this does not lead to a storage demand that is also independent of $n$. The rank of the basis vectors, i.e., the number of columns of the matrices $[\mathcal{V}_{1,1}, \ldots, \mathcal{V}_{1,m+1}]$ and $[\mathcal{Z}_{1,1}, \ldots, \mathcal{Z}_{1,m}]$ increases with the problem size. This trend is probably inherited from some intrinsic properties of the continuous problem. Indeed, the rank of the computed solution also grows with $n$ suggesting the idea that the rank of the exact solution increases with the problem size as well. Therefore, we are applying low-rank techniques to a problem whose low-rank approximability deteriorates for large $n$ and an increment in the memory requirements of our procedures is thus inevitable. A similar behavior is observed when decreasing the viscosity parameter $\nu$ as well.
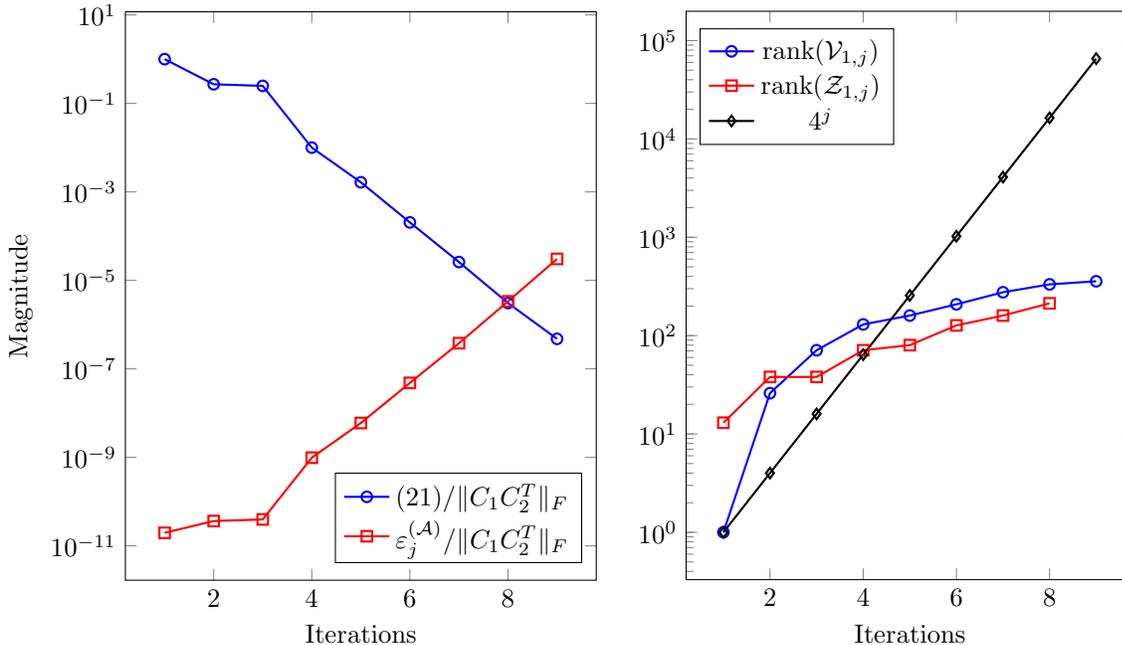
A growth in the rank of the basis vectors determines also a remarkable increment in the computational time as illustrated in Table 1. Indeed, the computational cost of basically all the steps of Algorithm 2, from the Arnoldi procedure and the low-rank truncations, to the preconditioning phase, depends on the rank of the basis vectors.

We also underline the fact that the true relative residual norm turns out to be always smaller than the normalized computed bound (21) validating the reliability of (21) as convergence check.

In Figure 1 (left) we report the normalized bound (21) together with the truncation threshold $\varepsilon_{\mathcal{A}}^{(j)}/\|C_1 C_2^T\|_F$ for the case $n = 5000$ and $\nu = 0.5$. We can appreciate how the tolerance for the low-rank truncations increases as the residual norm decreases. As already mentioned, this is a key element to obtain a solution

---

[4]A Matlab implementation is available at `http://www.dm.unibo.it/~ simoncin/software.html`.

Figure 1: Example 7.1, $n = 5000$, $\nu = 0.5$. Left: Normalized bound (21) and $\varepsilon_j^{(\mathcal{A})}/\|C_1 C_2^T\|_F$ for $j = 1, \ldots, 9$. Right: Rank of the matrix representing the $j$-th vector of the preconditioned and unpreconditioned basis.

procedure with a feasible storage demand. Moreover, in Figure 1 (right) we document the increment in the rank of the vectors of the preconditioned and unpreconditioned bases as the iterations proceed. We also plot the rank of the unpreconditioned basis we would obtain if no truncations (and no preconditioning steps) were performed, i.e., $4^j$. We can see how we would obtain full-rank basis vectors after very few iterations with consequent impracticable memory requirements of the overall solution process.

To conclude, in Figure 2, we report the inner product between the last basis vector we have computed and the previous ones, namely we report $\langle \mathcal{V}_{1,9} \mathcal{V}_{2,9}^T, \mathcal{V}_{1,j} \mathcal{V}_{2,j}^T \rangle_F$ for $j = 1, \ldots, 9$. This numerically confirms that the strategy illustrated in section 3.2 is able to maintain the orthogonality of the basis.

**Example 7.2.** In the second example we consider the algebraic problem stemming from the discretization of stochastic steady-state diffusion equations. In particular, given a sufficiently regular spatial domain $D$ and a sample space $\Omega$ associated with the probability space $(\Omega, \mathcal{F}, \mathbb{P})$, we seek an approximation to the function $u : D \times \Omega \to \mathbb{R}$ which is such that $\mathbb{P}$-almost surely

$$
\begin{aligned}
-\nabla \cdot (a(x,\omega)\nabla u(x,\omega)) &= f(x), &&\text{in } D, \\
u(x,\omega) &= 0, &&\text{on } \partial D.
\end{aligned}
\tag{31}
$$

We consider $D = [-1, 1]^2$ and we suppose $a$ to be a random field of the form

$$
a(x,\omega) = a_0(x) + \sum_{i=1}^{r} a_i(x)\sigma_i(\omega),
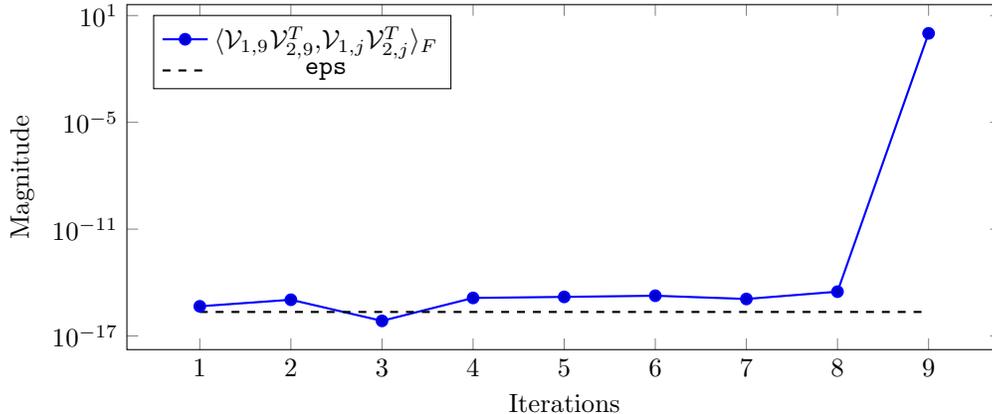$$

where $\sigma_i : \Omega \to \Gamma_i \subset \mathbb{R}$ are real-valued independent random variables (RVs).

In our case, $a(x,\omega)$ is a truncated Karhunen-Loève (KL) expansion

$$
a(x,\omega) = \mu(x) + \theta \sum_{i=1}^{r} \sqrt{\lambda_i}\phi_i(x)\sigma_i(\omega).
\tag{32}
$$

See, e.g., [36] for more details.

17

Figure 2: Example 7.1, $n = 5000$, $\nu = 0.5$. $\langle \mathcal{V}_{1,9}\mathcal{V}_{2,9}^T, \mathcal{V}_{1,j}\mathcal{V}_{2,j}^T \rangle_F$ for $j = 1,\ldots,9$. `eps` denotes machine precision.



The stochastic Galerkin method discussed in, e.g., $[2, 17, 42, 43, 57]$, lead to a discrete problem that can be written as a matrix equation of the form

$$K_0 X G_0^T + \sum_{i=1}^{r} K_i X G_i^T = f_0 g_0^T, \tag{33}$$

where $K_i \in \mathbb{R}^{n_x \times n_x}$, $G_i \in \mathbb{R}^{n_\sigma \times n_\sigma}$, and $f_0 \in \mathbb{R}^{n_x}$, $g_0 \in \mathbb{R}^{n_\sigma}$. See, e.g., $[42, 43]$.

We solve equation (33) by LR-GMRES and the following operators

$$
\begin{array}{llll}
\mathcal{P}_{\texttt{mean}}: & \mathbb{R}^{n_x \times n_\sigma} \rightarrow \mathbb{R}^{n_x \times n_\sigma} & \mathcal{P}_{\texttt{Ullmann}}: & \mathbb{R}^{n_x \times n_\sigma} \rightarrow \mathbb{R}^{n_x \times n_\sigma} \\
& X \mapsto K_0 X, & & X \mapsto K_0 X \overline{G}^T, \quad \overline{G} := \sum_{i=0}^{r} \frac{\text{trace}(K_i^T K_0)}{\text{trace}(K_0^T K_0)} G_i,
\end{array}
$$

are selected as preconditioners. $\mathcal{P}_{\texttt{mean}}$ is usually referred to as mean-based preconditioner, see, e.g., $[42, 43]$ and the references therein, while Ullmann proposed $\mathcal{P}_{\texttt{Ullmann}}$ in $[57]$.

Both $\mathcal{P}_{\texttt{mean}}$ and $\mathcal{P}_{\texttt{Ullmann}}$ are very well-suited for our framework as their application amount to the solution of a couple of linear systems so that the rank of the current basis vector does not increase. See the discussion in section 5. Moreover, supposing that these linear systems can be solved exactly by, e.g., a sparse direct solver, there is no need to employ flexible GMRES so that only one basis has to be stored. In particular, in all our tests, we precompute once and for all the LU factors of the matrices[5] which define the selected preconditioner so that only triangular systems are solved during the LR-GMRES iterations.

We generate instances of (33) with the help of the S-IFISS[6] package version 1.04; see $[49]$. The S-IFISS routine `stoch_diff_testproblem_pc` is executed to generate two instances of (33). The first equation (`Data 1`) is obtained by using a spatial discretization with $2^7$ points in each dimension, $r = 2$ RVs in (32) which are approximated by polynomial chaos expansions of length $\ell = 100$ leading to $n_x = 16129$, $n_\sigma = 5151$, and $r + 1 = 3$. The second instance (`Data 2`) was generated with $2^8$ grid points, $r = 5$, and chaos expansions of length $\ell = 10$ resulting in $n_x = 65025$, $n_\sigma = 3003$, and $r + 1 = 6$.

Table 2 summarizes the results and apparently problem `Data 2` is much more challenging than `Data 1`. This is meanly due to the number of terms in (33). Indeed, the effectiveness of the preconditioners may deteriorate as $r$ increases even though the actual capability of $\mathcal{P}_{\texttt{mean}}$ and $\mathcal{P}_{\texttt{Ullmann}}$ in reducing the iteration count is related to the coefficients of the KL expansion (32). See, e.g., $[42, \text{Theorem } 3.8]$ and $[57, \text{Corollary } 5.4]$. Moreover, $r + 1$ terms are involved in the products in line 2 of Algorithm 2 and a sizable $r$ leads, in general, to a faster growth in the rank of the basis vectors so that a larger number of columns are retained during the truncation step in line 3. As a result, the computational cost of our iterative scheme increases as well leading to a rather time consuming routine.

---

[5] The computational time of such decompositions is always included in the reported results.

[6] Available at `https://personalpages.manchester.ac.uk/staff/david.silvester/ifiss/sifiss.html`

Table 2: Example 7.2. Results of preconditioned LR-GMRES applied to different test problems. `Data 1`: $n_x = 16129$, $n_\sigma = 5151$, $r + 1 = 3$, `Data 2`: $n_x = 65025$, $n_\sigma = 3003$, $r + 1 = 6$.

| Prec. | Its | rank($S_1 S_2^T$) | Mem. | Conv. Checks $(21)/\|C_1 C_2^T\|_F$ | Real Res. | Time (s) |
|---|---|---|---|---|---|---|
| \multicolumn{7}{c}{Data 1} ||||||| 
| $\mathcal{P}_{\texttt{Ullmann}}$ | 9 | 44 | 220 | 3.703e-7 | 3.551e-7 | 1.204e1 |
| $\mathcal{P}_{\texttt{mean}}$ | 13 | 64 | 507 | 7.636e-7 | 7.369e-7 | 2.521e1 |
| \multicolumn{7}{c}{Data 2} ||||||| 
| $\mathcal{P}_{\texttt{Ullmann}}$ | 15 | 791 | 10266 | 5.611e-7 | 5.359e-7 | 8.847e4 |
| $\mathcal{P}_{\texttt{mean}}$ | 20 | 806 | 14912 | 8.118e-7 | 7.703e-7 | 1.626e5 |

If the discrete operator stemming from the discretization of (31) is well posed, then it is also symmetric positive definite and the CG method can be employed in the solution process. See, e.g., [42, Section 3]. We thus try to apply the (preconditioned) low-rank variant of CG (LR-CG) to the matrix equation (33). To this end, we adopt the LR-CG implementation proposed in [8]. With the notation of [8, Algorithm 1] we truncate all the iterates $X_{k+1}$, $R_{k+1}$, $P_{k+1}$ and $Q_{k+1}$. In particular, the threshold for the truncation of $X_{k+1}$ is set to $10^{-12}$ while the value on the right-hand side of (28) is used at the $k$-th LR-CG iteration for the low-rank truncation of all the other iterates. We want to point out that in the LR-CG implementation proposed in [8], the residual matrix $R_{k+1}$ is explicitly calculated by means of the current approximate solution $X_{k+1}$. We compute the residual norm before truncating $R_{k+1}$ so that what we are actually evaluating is the true residual norm and not an upper bound thereof.

The results are collected in Table 3 where the column "Mem." reports the maximum number of columns that had to be stored in the low-rank factors of all the iterates $X_{k+1}$, $R_{k+1}$, $P_{k+1}$, $Q_{k+1}$, and $Z_{k+1}$.
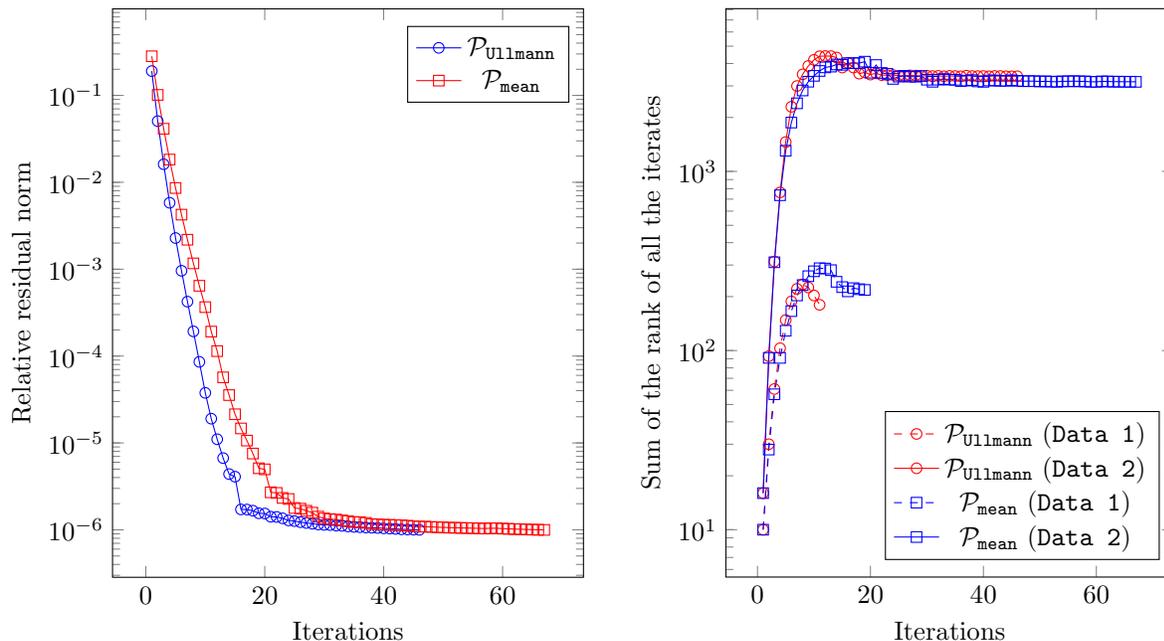
Table 3: Example 7.2. Results of preconditioned LR-CG applied to different test problems. `Data 1`: $n_x = 16129$, $n_\sigma = 5151$, $r + 1 = 3$, `Data 2`: $n_x = 65025$, $n_\sigma = 3003$, $r + 1 = 6$.

| Prec. | Its | rank($S_1 S_2^T$) | Mem. | Real Res. | Time (s) |
|---|---|---|---|---|---|
| \multicolumn{6}{c}{Data 1} |||||| 
| $\mathcal{P}_{\texttt{Ullmann}}$ | 11 | 41 | 234 | 9.517e-7 | 1.921e0 |
| $\mathcal{P}_{\texttt{mean}}$ | 19 | 52 | 288 | 9.629e-7 | 3.369e0 |
| \multicolumn{6}{c}{Data 2} |||||| 
| $\mathcal{P}_{\texttt{Ullmann}}$ | 46 | 483 | 4404 | 9.976e-7 | 9.642e2 |
| $\mathcal{P}_{\texttt{mean}}$ | 67 | 450 | 4096 | 9.981e-7 | 1.325e3 |

Except for `Data 1` with $\mathcal{P}_{\texttt{Ullmann}}$ as a preconditioner where LR-GMRES and LR-CG show similar results especially in terms of memory requirements, LR-CG allows for a much lower storage demand with a consequent reduction in the total computational efforts while achieving the prescribed accuracy. However, for `Data 2`, LR-CG requires a rather large number of iterations to converge regardless of the adopted preconditioner. This is due to a very small reduction of the residual norm, almost a stagnation, from one iteration to the following one we observe in the final stage of the algorithm. See Figure 3 (left). This issue may be fixed by employing a more robust, possibly more conservative, threshold for the low-rank truncations. Alternatively, a condition of the form $\|X_k - X_{k+1}\|_F \leqslant \varepsilon$ can be included in the convergence check as proposed in [43].

We conclude by mentioning a somehow surprising behavior of LR-CG. In particular, in the first iterations the rank of all the iterates increases as expected, while it starts decreasing from a certain $\bar{k}$ on until it reaches an almost constant value. See Figure 3 (right). This trend allows for a feasible storage demand also when many iterations are performed as for `Data 2`. We think that such a phenomenon deserves further studies.

Figure 3: Example 7.2. Left: LR-CG relative residual norm for `Data 2`. Right: Sum of the rank of all the LR-CG iterates $X_{k+1}$, $R_{k+1}$, $P_{k+1}$, $Q_{k+1}$, and $Z_{k+1}$ as the iterations proceed.



# 8   Conclusions

Low-rank Krylov methods are one of the few options for solving general linear matrix equations of the form (1), especially for large problem dimensions. An important step of these procedures consist in truncating the rank of the basis vectors to maintain a feasible storage demand of the overall solution process. In principle, such truncations can severely impact on the converge of the adopted Krylov routine.

In this paper we have shown how to perform the low-rank truncations in order to maintain the convergence of the selected Krylov procedure. In particular, our analysis points out that not only the thresholds employed for the truncations are important, but also the actual procedure adopted for the low-rank truncations plays a fundamental role. Indeed, such a routine must be able to preserve the orthogonality of the computed basis.

## Acknowledgments

## References

[1] A. C. Antoulas. *Approximation of large-scale dynamical systems*, volume 6 of *Advances in Design and Control*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2005.

[2] I. Babuška, R. Tempone, and G. E. Zouraris. Galerkin finite element approximations of stochastic elliptic partial differential equations. *SIAM J. Numer. Anal.*, 42(2):800–825, 2004.

[3] J. Baglama and L. Reichel. Augmented implicitly restarted Lanczos bidiagonalization methods. *SIAM J. Sci. Comput.*, 27(1):19–42, 2005.

[4] J. Baker, M. Embree, and J. Sabino. Fast singular value decay for Lyapunov solutions with nonnormal coefficients. *SIAM J. Matrix Anal. Appl.*, 36(2):656–668, 2015.

[5] M. Baumann, R. Astudillo, Y. Qiu, E. Y. M. Ang, M. B. van Gijzen, and R.-É. Plessix. An MSSS-preconditioned matrix equation approach for the time-harmonic elastic wave equation at multiple frequencies. *Computational Geosciences*, 22(1):43–61, Feb 2018.

[6] U. Baur. Low rank solution of data-sparse Sylvester equations. *Numer. Linear Algebra Appl.*, 15(9):837–851, 2008.

[7] U. Baur and P. Benner. Factorized solution of Lyapunov equations based on hierarchical matrix arithmetic. *Computing*, 78(3):211–234, 2006.

[8] P. Benner and T. Breiten. Low rank methods for a class of generalized Lyapunov equations and related issues. *Numer. Math.*, 124(3):441–470, 2013.

[9] P. Benner and T. Damm. Lyapunov equations, energy functionals, and model order reduction of bilinear and stochastic systems. *SIAM J. Control Optim.*, 49(2):686–711, 2011.

[10] P. Benner, R.-C. Li, and N. Truhar. On the ADI method for Sylvester equations. *J. Comput. Appl. Math.*, 233(4):1035–1045, 2009.

[11] P. Benner, M. Ohlberger, A. Cohen, and K. Willcox. *Model Reduction and Approximation*. Society for Industrial and Applied Mathematics, Philadelphia, PA, 2017.

[12] P. Benner, A. Onwunta, and M. Stoll. Low-rank solution of unsteady diffusion equations with stochastic coefficients. *SIAM/ASA Journal on Uncertainty Quantification*, 3(1):622–649, 2015.

[13] A. Bouras and V. Frayssé. Inexact matrix-vector products in Krylov methods for solving linear systems: a relaxation strategy. *SIAM J. Matrix Anal. Appl.*, 26(3):660–678, 2005.

[14] T. Breiten, V. Simoncini, and M. Stoll. Low-rank solvers for fractional differential equations. *Electron. Trans. Numer. Anal.*, 45:107–132, 2016.

[15] M. Che and Y. Wei. Randomized algorithms for the approximations of Tucker and the tensor train decompositions. *Advances in Computational Mathematics*, 45(1):395–428, Feb 2019.

[16] T. Damm. Direct methods and ADI-preconditioned Krylov subspace methods for generalized Lyapunov equations. *Numer. Linear Algebra Appl.*, 15(9):853–871, 2008.

[17] M. K. Deb, I. M. Babuška, and J. T. Oden. Solution of stochastic partial differential equations using Galerkin finite element techniques. *Comput. Methods Appl. Mech. Engrg.*, 190(48):6359–6372, 2001.

[18] S. V. Dolgov. TT-GMRES: solution to a linear system in the structured tensor format. *Russian J. Numer. Anal. Math. Modelling*, 28(2):149–172, 2013.

[19] V. Druskin and V. Simoncini. Adaptive rational Krylov subspaces for large-scale dynamical systems. *Systems Control Lett.*, 60(8):546–560, 2011.

[20] M. A. Freitag and D. L. H. Green. A low-rank approach to the solution of weak constraint variational data assimilation problems. *J. Comput. Phys.*, 357:263–281, 2018.

[21] R. W. Freund and N. M. Nachtigal. QMR: a quasi-minimal residual method for non-Hermitian linear systems. *Numer. Math.*, 60(3):315–339, 1991.

[22] S. Güttel. *Rational Krylov methods for operator functions*. PhD thesis, Technische Universität Bergakademie Freiberg, Germany, 2010. Available online from the Qucosa server.

[23] N. Halko, P. Martinsson, and J. Tropp. Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions. *SIAM Review*, 53(2):217–288, 2011.

[24] M. R. Hestenes and E. Stiefel. Methods of conjugate gradients for solving linear systems. *J. Research Nat. Bur. Standards*, 49:409–436 (1953), 1952.

[25] M. E. Hochstenbach. A Jacobi–Davidson type SVD method. *SIAM J. Sci. Comput.*, 23(2):606–628, 2001.

[26] E. Jarlebring, G. Mele, D. Palitta, and E. Ringh. Krylov methods for low-rank commuting generalized sylvester equations. *Numerical Linear Algebra with Applications*, 25(6), 2018. e2176.

[27] D. Kressner and L. Periša. Recompression of Hadamard products of tensors in Tucker format. *SIAM Journal on Scientific Computing*, 39(5):A1879–A1902, 2017.

[28] D. Kressner and P. Sirković. Truncated low-rank methods for solving general linear matrix equations. *Numer. Linear Algebra Appl.*, 22(3):564–583, 2015.

[29] D. Kressner, M. Steinlechner, and B. Vandereycken. Preconditioned low-rank Riemannian optimization for linear systems with tensor product structure. *SIAM J. Sci. Comput.*, 38(4):A2018–A2044, 2016.

[30] D. Kressner and C. Tobler. Low-rank tensor Krylov subspace methods for parametrized linear systems. *SIAM J. Matrix Anal. Appl.*, 32(4):1288–1316, 2011.

[31] P. Kürschner, S. Dolgov, K. D. Harris, and P. Benner. Greedy low-rank algorithm for spatial connectome regression. e-print 1808.05510, arXiv, 2018. math.NA.

[32] R. Larsen. Lanczos bidiagonalization with partial reorthogonalization. *DAIMI Report Series*, 27(537), 1998.

[33] J.-R. Li and J. White. Low-rank solution of Lyapunov equations. *SIAM Rev.*, 46(4):693–713, 2004.

[34] J. Liesen and Z. Strakos. *Krylov subspace methods: Principles and analysis*. Oxford University Press, 2012.

[35] Q. Liu, R. B. Morgan, and W. Wilcox. Polynomial preconditioned GMRES and GMRES-DR. *SIAM J. Sci. Comput.*, 37(5):S407–S428, 2015.

[36] G. J. Lord, C. E. Powell, and T. Shardlow. *An introduction to computational stochastic PDEs*. Cambridge Texts in Applied Mathematics. Cambridge University Press, New York, 2014.

[37] MATLAB. *version 9.3.0 (R2017b)*. The MathWorks Inc., Natick, Massachusetts, 2017.

[38] A. Onwunta. *Low-rank iterative solvers for stochastic Galerkin linear systems*. Dissertation, Otto-von-Guericke-Universität, Magdeburg, Germany, 2016.

[39] C. C. Paige and M. A. Saunders. Solutions of sparse indefinite systems of linear equations. *SIAM J. Numer. Anal.*, 12(4):617–629, 1975.

[40] D. Palitta and V. Simoncini. Matrix-equation-based strategies for convection-diffusion equations. *BIT*, 56(2):751–776, 2016.

[41] T. Penzl. Eigenvalue decay bounds for solutions of Lyapunov equations: the symmetric case. *Systems Control Lett.*, 40(2):139–144, 2000.

[42] C. E. Powell and H. C. Elman. Block-diagonal preconditioning for spectral stochastic finite-element systems. *IMA J. Numer. Anal.*, 29(2):350–375, 2009.

[43] C. E. Powell, D. Silvester, and V. Simoncini. An efficient reduced basis solver for stochastic Galerkin matrix equations. *SIAM J. Sci. Comput.*, 39(1):A141–A163, 2017.

[44] E. Ringh, G. Mele, J. Karlsson, and E. Jarlebring. Sylvester-based preconditioning for the waveguide eigenvalue problem. *Linear Algebra Appl.*, 542:441–463, 2018.

[45] Y. Saad. A flexible inner-outer preconditioned GMRES algorithm. *SIAM J. Sci. Comput.*, 14(2):461–469, 1993.

[46] Y. Saad. *Iterative methods for sparse linear systems.* SIAM, Society for Industrial and Applied Mathematics, Philadelphia, PA, 2nd edition, 2003.

[47] Y. Saad and M. H. Schultz. GMRES: a generalized minimal residual algorithm for solving nonsymmetric linear systems. *SIAM J. Sci. Statist. Comput.*, 7(3):856–869, 1986.

[48] S. D. Shank, V. Simoncini, and D. B. Szyld. Efficient low-rank solution of generalized Lyapunov equations. *Numer. Math.*, 134(2):327–342, 2016.

[49] D. J. Silvester, A. Bespalov, and C. E. Powell. *S-IFISS version 1.04*, 2017.

[50] V. Simoncini. A new iterative method for solving large-scale Lyapunov matrix equations. *SIAM J. Sci. Comput.*, 29(3):1268–1288, 2007.

[51] V. Simoncini. Computational methods for linear matrix equations. *SIAM Rev.*, 58(3):377–441, 2016.

[52] V. Simoncini and D. B. Szyld. Flexible inner-outer Krylov subspace methods. *SIAM J. Numer. Anal.*, 40(6):2219–2239 (2003), 2002.

[53] V. Simoncini and D. B. Szyld. Theory of inexact Krylov subspace methods and applications to scientific computing. *SIAM J. Sci. Comput.*, 25(2):454–477, 2003.

[54] V. Simoncini and D. B. Szyld. Recent computational developments in Krylov subspace methods for linear systems. *Numer. Linear Algebra Appl.*, 14(1):1–59, 2007.

[55] M. Stoll. A Krylov-Schur approach to the truncated SVD. *Linear Algebra Appl.*, 436(8):2795–2806, 2012.

[56] M. Stoll and T. Breiten. A low-rank in time approach to PDE-constrained optimization. *SIAM J. Sci. Comput.*, 37(1):B1–B29, 2015.

[57] E. Ullmann. A Kronecker product preconditioner for stochastic Galerkin finite element discretizations. *SIAM J. Sci. Comput.*, 32(2):923–946, 2010.

[58] J. van den Eshof and G. L. G. Sleijpen. Inexact Krylov subspace methods for linear systems. *SIAM J. Matrix Anal. Appl.*, 26(1):125–153, 2004.

[59] H. A. van der Vorst. Bi-CGSTAB: a fast and smoothly converging variant of Bi-CG for the solution of nonsymmetric linear systems. *SIAM J. Sci. Statist. Comput.*, 13(2):631–644, 1992.

[60] P. M. Van Dooren. Structured linear algebra problems in digital signal processing. In *Numerical linear algebra, digital signal processing and parallel algorithms (Leuven, 1988)*, volume 70 of *NATO Adv. Sci. Inst. Ser. F Comput. Systems Sci.*, pages 361–384. Springer, Berlin, 1991.

[61] M. B. van Gijzen. A polynomial preconditioner for the GMRES algorithm. *J. Comput. Appl. Math.*, 59(1):91–107, 1995.

[62] B. Vandereycken and S. Vandewalle. A Riemannian optimization approach for computing low-rank solutions of Lyapunov equations. *SIAM J. Matrix Anal. Appl.*, 31(5):2553–2579, 2010.

[63] R. Weinhandl, P. Benner, and T. Richter. Low-rank Linear Fluid-structure Interaction Discretizations. *arXiv e-prints*, May 2019. ArXiv: 1905.11000.