

LOW-RANK LINEAR FLUID-STRUCTURE INTERACTION DISCRETIZATIONS*

ROMAN WEINHANDL[†], PETER BENNER[‡], AND THOMAS RICHTER[§]

Abstract. Fluid-structure interaction models involve parameters that describe the solid and the fluid behavior. In simulations, there is a need to vary these parameters to examine the behavior of a fluid-structure interaction model for different solids and different fluids. A shipping company wants to know how the material, a ship's hull is made of, interacts with fluids at different Reynolds and Strouhal numbers before the building process takes place. Also, the behavior of such models for solids with different properties is considered before the prototype phase. A parameter-dependent linear fluid-structure interaction discretization provides approximations for a bundle of different parameters at one step. Such a discretization with respect to different material parameters leads to a big block diagonal system matrix that is equivalent to a matrix equation as discussed in [6]. The unknown is then a matrix which can be approximated using a low-rank approach that represents the iterate by a tensor. This paper compares a low-rank GMRES variant as first mentioned in [1] with a variant of the Chebyshev iteration. Numerical experiments show that such truncated methods applied to parameter-dependent discretizations provide approximations with relative residual norms smaller than 10^{-8} within a twentieth of the time used by individual standard approaches.

Key words. Parameter-dependent fluid-structure interaction, GMREST, ChebyshevT, low-rank, tensor

1. Introduction. A parameter-dependent linear fluid-structure interaction discretization using bilinear finite elements with a total number of $M \in \mathbb{N}$ degrees of freedom and $m \in \mathbb{N}$ parameter combinations leads to equations of the form

$$(1) \quad (A_0 + \mu_s^i A_1 + \lambda_s^i A_2 + \rho_f^i A_3)x_i = b_D \text{ for } i \in \{1, \dots, m\}$$

where the discretization matrices $A_0, A_1, A_2, A_3 \in \mathbb{R}^{M \times M}$ and the right hand side $b_D \in \mathbb{R}^M$ depends on the Dirichlet boundary conditions and the i th finite element solution $x_i \in \mathbb{R}^M$. The samples of interest are given by the shear moduli $\mu_s^i \in \mathbb{R}$, the first Lamé parameters $\lambda_s^i \in \mathbb{R}$ and the fluid densities ρ_f^i for $i \in \{1, \dots, m\}$.

Equation (1) can directly be written as the linear system

$$(2) \quad \mathcal{A} \begin{pmatrix} x_1 \\ \vdots \\ x_m \end{pmatrix} = \begin{pmatrix} b_D \\ \vdots \\ b_D \end{pmatrix}$$

where $\mathcal{A} \in \mathbb{R}^{Mm \times Mm}$ is a block diagonal matrix. Following [6], equation (2) can then be translated into the matrix equation

$$(3) \quad A_0 X + A_1 X D_1 + A_2 X D_2 + A_3 X D_3 = B$$

with $B := [b_D \mid \dots \mid b_D]$ and the diagonal matrices $D_1, D_2, D_3 \in \mathbb{R}^{M \times M}$ where the i th diagonal entry of these diagonal matrices is given by μ_s^i, λ_s^i and ρ_f^i , respectively.

*Submitted May 28, 2019.

Funding: Funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) - 314838170, GRK 2297 MathCoRe.

[†]Max-Planck Institute for Dynamics of Complex Technical Systems, Magdeburg, Germany (weinhandl@mpi-magdeburg.mpg.de).

[‡]Max-Planck Institute for Dynamics of Complex Technical Systems, Magdeburg, Germany (benner@mpi-magdeburg.mpg.de).

[§]Otto von Guericke University Magdeburg, Germany (thomas.richter@ovgu.de).

In (3) the unknown

$$X = [x_1 \mid \cdots \mid x_m] \in \mathbb{R}^{M \times m}$$

is a matrix. Now an iterative method to solve linear systems can be modified such that it uses an iterate that is a matrix. It is applied to the big system (2) but computation is kept in the matrix notation (3) by representing the iterate as a matrix instead of a vector. The methods used in this paper fix a rank $R \in \mathbb{N}$, $R \ll M, m$ and represent this iterate as a tensor. The goal is to find a low-rank approximation \hat{X} of rank R

$$\hat{X} = \sum_{j=1}^R u_j \otimes v_j^T, \quad u_j \in \mathbb{R}^M \text{ and } v_j \in \mathbb{R}^M \quad \forall j \in \{1, \dots, R\}$$

that approximates the full matrix X from (3) and therefore provides (parameter-dependent) finite element approximations for all equations in (1).

Fluid-structure interaction problems yield non symmetric system matrices. Hence the system matrix \mathcal{A} in (2) is not symmetric. The methods examined in this paper are based on the GMRES method as introduced in a truncated variant in [1] and the Chebyshev method from [8]. These methods will then be compared to a truncated method based on the Bi-CGstab method from [15] similar to [6, Algorithm 3].

In numerical experiments, the truncated approaches converged significantly faster than standard approaches applied to m individual equations of the form (1). As a consequence of the low-rank representation of the iterate, the storage needed to store the approximation provided by the truncated approaches is notably smaller than the full storage needed by standard approaches.

2. The Stationary Linear Fluid-structure Interaction Problem. Let $d \in \{2, 3\}$, $\Omega \subset \mathbb{R}^d$, $F, S \subset \Omega$ such that $\bar{F} \cup \bar{S} = \bar{\Omega}$ and $F \cap S = \emptyset$ where F represents the fluid and S the solid part. Let $\Gamma_{\text{int}} = \partial F \cap \partial S$ and $\Gamma_f^{\text{out}} \subset \partial F \setminus \partial S$ denote the boundary part where Neumann outflow conditions hold. $\Gamma_f^D = \partial F \setminus (\Gamma_f^{\text{out}} \cup \Gamma_{\text{int}})$ denotes the boundary part where Dirichlet conditions hold. Consider the Stokes fluid equations [10, Chapter 2.4.4] as a model for the fluid part and the Navier-Lamé equations [10, Problem 2.23] as a model for the solid part. Both equations are assumed to have a vanishing right hand side. If these two equations are coupled with the kinematic and the dynamic coupling conditions [10, Chapter 3.1], the weak formulation of the stationary coupled linear fluid-structure interaction problem reads

$$(4) \quad \begin{aligned} \langle \nabla \cdot v, \xi \rangle_F &= 0, \\ \mu_s \langle \nabla u + \nabla u^T, \nabla \varphi \rangle_S + \lambda_s \langle \text{tr}(\nabla u) I, \nabla \varphi \rangle_S \\ &+ \nu_f \rho_f \langle \nabla v + \nabla v^T, \nabla \varphi \rangle_F - \langle p, \nabla \cdot \varphi \rangle_F = 0, \\ \langle \nabla u, \nabla \psi \rangle_F &= 0, \end{aligned}$$

where the trial functions $v \in v_{\text{in}} + H_0^1(\Omega, \Gamma_f^D \cup \Gamma_{\text{int}})$ (velocity), where $v_{\text{in}} \in H^1(\Omega)$ is an extension of the Dirichlet data on Γ_f^D , $u \in H_0^1(\Omega)$ (deformation) and $p \in L^2(F)$ (pressure) and the test functions $\xi \in L^2(F)$ (divergence equation) $\varphi \in H_0^1(\Omega, \partial\Omega \setminus \Gamma_f^{\text{out}})$ (momentum equation) and $\psi \in H_0^1(F)$ (deformation equation) are involved. $\langle \cdot \rangle_S$ and $\langle \cdot \rangle_F$ denote the \mathcal{L}^2 scalar product on S and F , respectively. $\nu_f \in \mathbb{R}$ denotes the kinematic fluid viscosity and $\rho_f \in \mathbb{R}$ the fluid density. The shear modulus $\mu_s \in \mathbb{R}$ and the first Lamé parameter $\lambda_s \in \mathbb{R}$ determine the Poisson ratio of the solid.

DEFINITION 2.1 (The Poisson Ratio [10, Definition 2.18]).

The Poisson ratio of a solid is given by the number

$$\nu_s^p = \frac{\lambda_s}{2(\lambda_s + \mu_s)}.$$

It describes the compressibility of a solid.

3. Parameter-dependent Discretization. Assume the behavior of a linear fluid-structure interaction model for $m_1 \in \mathbb{N}$ shear moduli, $m_2 \in \mathbb{N}$ first Lamé parameters and $m_3 \in \mathbb{N}$ fluid densities is of interest. The kinematic fluid viscosity $\nu_f \in \mathbb{R}$ is assumed to be fixed. Let the samples of interest be given by the following sets.

$$\begin{aligned} \{\mu_s^{i_1}\}_{i_1 \in \{1, \dots, m_1\}} &\subset \mathbb{R}^+, \text{ a set of shear moduli,} \\ \{\lambda_s^{i_2}\}_{i_2 \in \{1, \dots, m_2\}} &\subset \mathbb{R}^+, \text{ a set of first Lamé parameters and} \\ \{\rho_f^{i_3}\}_{i_3 \in \{1, \dots, m_3\}} &\subset \mathbb{R}^+, \text{ a set of fluid densities.} \end{aligned}$$

In a bilinear finite element discretization of (4) with a mesh grid size of $N \in \mathbb{N}$ every mesh grid point corresponds to a pressure, a velocity and a deformation variable. In two dimensions the velocity and deformation are two dimensional vectors, in three dimensions they correspond to a three dimensional vector each. The total number of degrees of freedom is therefore $M = 5N$ in two dimensions and $M = 7N$ in three dimensions.

Let Ω_h be a matching mesh of the domain Ω as defined in [10, Definition 5.9] with N mesh grid points, $A_0 \in \mathbb{R}^{M \times M}$ a Q_1 discretization matrix of all functionals involved in (4) with a fixed shear modulus $\mu_s \in \mathbb{R}$, a fixed first Lamé parameter $\lambda_s \in \mathbb{R}$ and a fixed fluid density $\rho_f \in \mathbb{R}$. Moreover, let $A_1, A_2, A_3 \in \mathbb{R}^{M \times M}$ be Q_1 discretization matrices of the following functionals.

$$\begin{aligned} A_1 &\text{ discretizes } \langle \nabla u + \nabla u^T, \nabla \varphi \rangle_S, \\ A_2 &\text{ discretizes } \langle \text{tr}(\nabla u)I, \nabla \varphi \rangle_S \text{ and} \\ A_3 &\text{ discretizes } \langle \nabla v + \nabla v^T, \nabla \varphi \rangle_F. \end{aligned}$$

The parameter-dependent equation

$$(5) \quad \underbrace{(A_0 + (\mu_s^{i_1} - \mu_s)A_1 + (\lambda_s^{i_2} - \lambda_s)A_2 + \nu_f(\rho_f^{i_3} - \rho_f)A_3)}_{=: A(\mu_s^{i_1}, \lambda_s^{i_2}, \rho_f^{i_3})} x_{i_1 i_2 i_3} = b_D \text{ for} \\ (i_1, i_2, i_3) \in \{1, \dots, m_1\} \times \{1, \dots, m_2\} \times \{1, \dots, m_3\}$$

is the finite element discretization of (4) related to a shear modulus $\mu_s^{i_1}$, a first Lamé parameter $\lambda_s^{i_2}$ and a fluid density $\rho_f^{i_3}$. The finite element solution is $x_{i_1 i_2 i_3} \in \mathbb{R}^M$ and the right hand side $b_D \in \mathbb{R}^M$ depends on Dirichlet boundary conditions.

Combining all sample combinations in (5) leads to a total of $m = m_1 m_2 m_3$ equations. Written as a linear system, these equations translate to

$$(6) \quad \underbrace{\text{diag}(A(\mu_s^{i_1}, \lambda_s^{i_2}, \rho_f^{i_3}))}_{\substack{i_1 \in \{1, \dots, m_1\} \\ i_2 \in \{1, \dots, m_2\} \\ i_3 \in \{1, \dots, m_3\}}} \begin{pmatrix} x_1 \\ \vdots \\ x_{m_1 m_2 m_3} \end{pmatrix} = \begin{pmatrix} b_D \\ \vdots \\ b_D \end{pmatrix}.$$

from left is equivalent to multiplication of P_T^{-1} to $F(X)$ from left using the matrix notation from (8).

4. The Low-rank Methods. Now, a method can be applied to the big system (6). The iterate is then a vector $x \in \mathbb{R}^{Mm}$. But if the iterate is represented as a matrix instead of a vector, computation can be kept in the matrix notation from (8). For instance, the matrix-vector multiplication in such a global approach corresponds to the evaluation of the function $F(\cdot)$ from (8). The Euclidean norm of the vector

$$\begin{pmatrix} x_1 \\ \vdots \\ x_{m_1 m_2 m_3} \end{pmatrix}$$

from (6) then corresponds to the Frobenius norm of the matrix X in (8), $\|X\|_F$. There are many methods that are suitable for this approach. But since, for fluid-structure interaction problems, the matrix \mathcal{A} is not symmetric, the focus in this paper lies on methods that base on the GMRES and the Chebyshev method. As proved in Theorem 35.2 in [14], the GMRES method converges in this case, and so does the Chebyshev method, if all eigenvalues of the system matrix lie in an ellipse that does not touch the imaginary axis as proved in [8]. Also the Bi-CGSTAB method from [15] is considered for a numerical comparison.

As mentioned, the low-rank methods this paper deals with use an iterate that is, instead of a matrix, a tensor of order two. The iterate is then given by

$$\hat{X} = \sum_{j=1}^R (u_j \otimes v_j^T) \text{ with } u_j \in \mathbb{R}^M, v_j \in \mathbb{R}^m \forall j \in \{1, \dots, R\}$$

where the tensor rank $R \in \mathbb{N}$ is kept small such that $R \ll M, m$. The goal of the method is to find a low-rank approximation \hat{X} that approximates the matrix X in (8). The methods GMREST (also mentioned in [1]) and ChebyshevT are such methods and will be explained in the following subsections. They are not just faster than the standard methods applied to m individual equations of the form (5), they also need a smaller amount of storage to store the approximation. If M and m are very big, this plays an important role since the storage amount to store \hat{X} is in $O((M+m)R)$ while the storage amount to store the full matrix X is in $O(Mm)$.

4.1. Tensor Format and Truncation. There are several formats available to represent the tensor \hat{X} . For $d = 2$, the hierarchical Tucker format ([4, Definition 11.11]) is equivalent to the Tucker format. It is based on so called minimal subspaces that are explained in [4, Chapter 6].

DEFINITION 4.1 (Tucker Format [4, Definition 8.1] for $d = 2$).

Let $V := \mathbb{R}^M \otimes \mathbb{R}^m$, $(r_1, r_2) \in \mathbb{N}^2$. For $d = 2$, the Tucker tensors of Tucker rank (r_1, r_2) are given by the set

$$T_{(r_1, r_2)}(V) := \left\{ v \in V : \left. \begin{array}{l} \text{there are subspaces } V_1 \subset \mathbb{R}^M \text{ and } V_2 \subset \mathbb{R}^m \text{ with} \\ \dim(V_1) = r_1, \dim(V_2) = r_2 \text{ and } v \in V_1 \otimes V_2. \end{array} \right\}.$$

From now on, the set $T_{(R,R)}(V)$ will be denoted by T_R . By a tensor of rank R a Tucker tensor in $T_{(R,R)}$ is addressed in the following.

As explained in [4, Chapter 13.1.4], summation of two arbitrary Tucker tensors of rank R , in general, results in a Tucker tensor of rank $2R$. But to keep a low-rank method fast the rank of the iterate has to be kept small. This induces the need for a truncation operator.

DEFINITION 4.2. *The truncation operator*

$$\mathcal{T} : \mathbb{R}^M \otimes \mathbb{R}^m \rightarrow T_R$$

maps a Tucker or a full tensor into the set of Tucker tensors of rank R .

Remark 4.3. As proved in [4, Chapter 3.2.3] it holds

$$\mathbb{R}^M \otimes \mathbb{R}^m \cong \mathbb{R}^{M \times m}.$$

Since for our purposes we consider a matrix that is represented by a tensor we assume

$$\mathcal{T} : \mathbb{R}^{M \times m} \rightarrow T_R.$$

and if

$$\hat{x} \in T_R$$

by \hat{x} the full representation of the tensor in \mathbb{R}^{Mm} is addressed for the sake of notation.

Before we proceed one more definition is needed.

DEFINITION 4.4 (Vectorization restricted to $\mathbb{R}^{M \times m}$). *The vectorization operator*

$$\text{vec} : \mathbb{R}^{M \times m} \rightarrow \mathbb{R}^{Mm}, \text{vec} \left(\begin{array}{c|c|c} v_1 & \cdots & v_m \end{array} \right) \mapsto \begin{pmatrix} v_1 \\ \vdots \\ v_m \end{pmatrix}$$

stacks matrix entries column wise into a vector. Its inverse maps to a $M \times m$ matrix so

$$\text{vec}^{-1} : \mathbb{R}^{Mm} \rightarrow \mathbb{R}^{M \times m}, \text{vec}^{-1} \left(\begin{pmatrix} v_1 \\ \vdots \\ v_m \end{pmatrix} \right) = (v_1 | \cdots | v_m).$$

Remark 4.5. The argument of the function $F(\cdot)$ from (8) is tacitly assumed to be a matrix so $F(\hat{x})$ addresses $F(\text{vec}^{-1}(\hat{x}))$ for $\hat{x} \in T_R$.

Since truncation is an operation that is applied after nearly every addition of tensors and multiple times in every iteration the format that provides the truncation with the least complexity is often the preferred one. According to [7, Algorithm 6], the htucker toolbox [7] for Matlab provides truncation with complexity $(2 \max(M, m)R^2 + 2R^4)$ if the input format is in hierarchical Tucker format. The truncation complexity of the TT toolbox [9] for Matlab that uses the Tensor Train format is in $O(2 \max(M, m)R^3)$ as stated in [9, Algorithm 2].

4.2. The GMREST and the GMRESTR Method. Consider \mathcal{A} from (6), a suitable preconditioner $\mathcal{P} = I_m \otimes P \in \mathbb{R}^{Mm \times Mm}$, a start vector $x_0 \in \mathbb{R}^{Mm}$ and

$$b := \begin{pmatrix} b_D \\ \vdots \\ b_D \end{pmatrix}, r_0 := \mathcal{P}^{-1}(b - \mathcal{A}x_0).$$

l GMRES iterations with the preconditioner \mathcal{P} applied to the system (6) minimize $\|r_0 - \mathcal{P}^{-1}\mathcal{A}z\|_2$ for $z \in \mathbb{R}^{Mm}$ over the Krylov subspace (compare [11, Chapter 6.2])

$$\mathcal{K}_l := \text{span}\{r_0, \mathcal{P}^{-1}\mathcal{A}r_0, \dots, (\mathcal{P}^{-1}\mathcal{A})^{l-1}r_0\}.$$

Algorithm 1 GMREST(l) (Preconditioned Truncated GMRES Method)

Input: Iteration number l , truncation rank R for \mathcal{T} , $F(\cdot)$ from (8), left preconditioner $P \in \mathbb{R}^{M \times M}$, right hand side $\hat{B} \in T_R$ and start matrix $\hat{X} \in T_R$

Output: Approximate solution $\hat{X} \in T_R$

Find $\hat{R} \in T_R$ such that $P\hat{R} = \mathcal{T}(\hat{B} - F(\hat{X}))$.

$$z := \left(\|\hat{R}\|_F \quad 0 \quad \cdots \quad 0 \right)^T$$

$$\hat{V}_1 := \frac{\hat{R}}{\|\hat{R}\|_F}$$

for $i = 1, \dots, l$ **do**

Find $\hat{W} \in T_R$ such that $P\hat{W} = \mathcal{T}(F(\hat{V}_i))$.

for $k = 1, \dots, i$ **do**

$$H_{k,i} := \text{trace}(\hat{V}_k^H \hat{W})$$

$$\hat{W} := \mathcal{T}(\hat{W} - H_{k,i} \hat{V}_k)$$

end for

$$H_{i+1,i} := \|\hat{W}\|_F$$

$$\hat{V}_{i+1} := \hat{W} \frac{1}{H_{i+1,i}}$$

end for

Now find a unitary matrix Q such that QH is a right upper triangular matrix via Givens rotation. Find y such that $QH y = Qz$.

$$\hat{X} = \mathcal{T}\left(\hat{X} + \sum_{j=1}^l y_j \hat{V}_j\right)$$

Algorithm 2 GMRESTR(l, d) (Preconditioned Truncated GMRES Restart Method)

Input: In addition to the inputs of [Algorithm 1](#) divider $d \in \mathbb{N}$

Output: Approximate solution $\hat{X} \in T_R$

$$d_1 := \text{floor}\left(\frac{l}{d}\right)$$

for $i = 1, \dots, d_1$ **do**

$\hat{X} = \text{GMREST}(d)$ with start matrix \hat{X}

end for

As mentioned before, from the theoretical point of view, this classical GMRES method is equivalent to the global GMRES method that uses an iterate that is a matrix instead of a vector. But if the iterate is represented by a tensor of a fixed rank R , the truncation operator \mathcal{T} generates an additional error every time it is applied to truncate the iterate or tensors involved back to rank R . With an initial guess $\hat{x}_0 := \mathcal{T}(x_0)$,

$$\hat{b} := \mathcal{T}(b) \text{ and } \hat{r}_0 := \mathcal{T}(P^{-1}[\hat{b} - F(\hat{x}_0)]),$$

l iterations of the truncated GMRES method GMREST that is coded in [Algorithm 1](#) minimize $\|\text{vec}(\mathcal{T}(\hat{r}_0 - P^{-1}F(\hat{z})))\|_2$ for $\hat{z} \in T_R$ over the truncated Krylov subspace

$$\mathcal{K}_l^{\mathcal{T}} := \text{span}\{\text{vec}(\hat{r}_0), \text{vec}(\mathcal{T}(P^{-1}F(\hat{r}_0))), \dots, \text{vec}((\mathcal{T}(P^{-1}F))^{l-1}(\hat{r}_0))\}.$$

Even the standard GMRES method can stagnate due to machine precision. This means that at the l th iteration the dimension of the numerical approximation of \mathcal{K}_l is smaller than l . As we will see later, the truncation operator brings, in addition to the machine precision error, a truncation error into play. As a result, the GMREST

Algorithm 3 ChebyshevT(l, d, c) (Preconditioned Truncated Chebyshev Method)

Input: Iteration number l , ellipse by center d and foci $d \pm c$, truncation rank R for \mathcal{T} , $F(\cdot)$ from (8), left preconditioner $P \in \mathbb{R}^{M \times M}$, right hand side $\hat{B} \in T_R$ and start matrix $\hat{X} \in T_R$

Output: Approximate solution $\hat{X} \in T_R$

Find \hat{R}_0 such that $P\hat{R}_0 = \mathcal{T}(\hat{B} - F(\hat{X}))$

$\hat{\Phi}_0 := \frac{1}{d}\hat{R}_0$

$\hat{X} = \mathcal{T}(\hat{X} + \hat{\Phi}_0)$

$t_0 := 1$

$t_1 := \frac{d}{c}$

for $i = 1, \dots, l$ **do**

$t_{i+1} := 2\frac{d}{c}t_i - t_{i-1}$

$\alpha_i := \frac{2t_i}{ct_{i+1}}$

$\beta_i := \frac{t_{i-1}}{t_{i+1}}$

Find \hat{R}_i such that $P\hat{R}_i = \mathcal{T}(\hat{B} - F(\hat{X}))$.

$\hat{\Phi}_i := \mathcal{T}(\alpha_i\hat{R}_i + \beta_i\hat{\Phi}_{i-1})$

$\hat{X} = \mathcal{T}(\hat{X} + \hat{\Phi}_i)$

end for

method can stagnate much earlier than the non truncated full approach. As in the full approach, restarting the method with the actual iterate as initial guess can be a remedy. This restarted variant of the GMREST method called GMRESTR method is coded in [Algorithm 2](#).

4.3. The ChebyshevT Method. In the complex plane, positive real values can be encircled by an ellipse that does not touch the imaginary axis. Hence, the Chebyshev method converges for non symmetric system matrices that have positive real eigenvalues.

The diagonal blocks of the preconditioned system matrix $\mathcal{P}_T^{-1}\mathcal{A}$ are

$$Bl(i_1, i_2, i_3) := P_T^{-1}(A_0 + (\mu_s^{i_1} - \mu_s)A_1 + (\lambda_s^{i_2} - \lambda_s)A_2 + \nu_f(\rho_f^{i_3} - \rho_f)A_3)$$

for $(i_1, i_2, i_3) \in \{1, \dots, m_1\} \times \{1, \dots, m_2\} \times \{1, \dots, m_3\}$.

Moreover, the parameter-dependent matrices (5) are assumed to be invertible. The eigenvalues of $\mathcal{P}_T^{-1}\mathcal{A}$ denoted by $\Lambda(\mathcal{P}_T^{-1}\mathcal{A})$ therefore coincide with the set

$$(9) \quad \bigcup_{\substack{i_1 \in \{1, \dots, m_1\} \\ i_2 \in \{1, \dots, m_2\} \\ i_3 \in \{1, \dots, m_3\}}} \Lambda(Bl(i_1, i_2, i_3)).$$

In numerical tests it turned out that $\Lambda(\mathcal{P}_T^{-1}\mathcal{A}) \subset (0, \infty)$ for the linear fluid-structure interaction problems considered and the minimum and the maximum of $\Lambda(\mathcal{P}_T^{-1}\mathcal{A})$ do not depend on the number of degrees of freedom. For a discretization the quantities

$$(10) \quad \Lambda_{\max} := \max \Lambda(\mathcal{P}_T^{-1}\mathcal{A}) \text{ and } \Lambda_{\min} := \min \Lambda(\mathcal{P}_T^{-1}\mathcal{A})$$

can therefore be computed using the representation (9) of $\Lambda(\mathcal{P}_T^{-1}\mathcal{A})$ for a small number of degrees of freedom. Since $\Lambda_{\min} > 0$ all the eigenvalues of $\mathcal{P}_T^{-1}\mathcal{A}$ can indeed be encircled by an ellipse in the complex plane that does not touch the imaginary

axis. The Chebyshev method from [8] can therefore be generalized in the same manner as the GMRES method in the subsection before and used to find a low-rank approximation \hat{X} of X in (8). The ellipse with center

$$d := \frac{\Lambda_{\min} + \Lambda_{\max}}{2} \text{ and foci } d \pm c \text{ for } c := \Lambda_{\max} - d$$

encircles all eigenvalues of $\mathcal{P}^{-1}\mathcal{A}$ and does not touch the imaginary axis, i.e. $d - c > 0$ if the eigenvalues of $\mathcal{P}_T^{-1}\mathcal{A}$ are positive. The resulting truncated Chebyshev variant ChebyshevT is coded in [Algorithm 3](#).

5. Time Discretization.

5.1. The Linear Fluid-structure Interaction Problem. Let $[0, T]$ be a time interval for $T \in \mathbb{R}^+$ and $t \in [0, T]$ be the time variable. The deformation u and the velocity v now depend, in addition, on the time variable t so we write $u(x, t)$ and $v(x, t)$. With the solid density $\rho_s \in \mathbb{R}$, the non stationary Navier Lamé equations [10, Chapter 2.3.1.2] fulfill

$$\rho_s \partial_{tt} u - \operatorname{div}(\sigma) = \rho_s \partial_t v - \operatorname{div}(\sigma) = 0, \quad \partial_t u = v.$$

The time term $\rho_f \partial_t v$ coming from the Stokes fluid equations as mentioned in [10, (2.42)] is added to the left side of the momentum equation. The weak formulation of the non stationary coupled linear fluid-structure interaction problem is given by

$$\begin{aligned} & \langle \nabla \cdot v, \xi \rangle_F = 0, \\ (11) \quad & \overbrace{\rho_f \langle \partial_t v, \varphi \rangle_F}^{(*)} + \rho_s \langle \partial_t v, \varphi \rangle_S + \mu_s \langle \nabla u + \nabla u^T, \nabla \varphi \rangle_S \\ & + \lambda_s \langle \operatorname{tr}(\nabla u) I, \nabla \varphi \rangle_S + \nu_f \rho_f \langle \nabla v + \nabla v^T, \nabla \varphi \rangle_F - \langle p, \nabla \cdot \varphi \rangle_F = 0, \\ & \langle \nabla u, \nabla \psi \rangle_F = 0 \end{aligned}$$

with regularity conditions $v(x, \cdot) \in C^1([0, T])$, $\partial_t v(\cdot, t) \in H^{-1}(\Omega)$ for all $(x, t) \in \Omega \times [0, T]$ in addition to the ones in (4).

Remark 5.1. The time term in the Navier-Stokes equations can also, in some applications, be neglected. Then one is interested in a time discretization of the solid part only. In this case $(\star) = 0$ and the discretization matrix A_t^f of the following subsection can be seen as the zero matrix in $\mathbb{R}^{M \times M}$.

5.2. Time Discretization with the θ -Scheme. Let $A_t^f, A_t^s \in \mathbb{R}^{M \times M}$ be Q_1 discretization matrices.

$$A_t^f \text{ discretizes } \langle v, \varphi \rangle_F \text{ and } A_t^s \text{ discretizes } \rho_s \langle v, \varphi \rangle_S.$$

Now consider a discretization that splits the time interval $[0, T]$ into $s+1 \in \mathbb{N}$ equidistant time steps. Let the distance between two time steps be Δ_t . The starting time is $t_0 = 0$ and the following times are given by $t_i := i\Delta_t$ for $i \in \{1, \dots, s\}$. Let X^i be the approximate solution at time t_i , X^0 is given as the initial value. The given Dirichlet boundary conditions x_0^i at time t_i for all $i \in \{1, \dots, s\}$ yield the time dependent right hand side

$$B^i := x_0^i \otimes (1, \dots, 1) \text{ for } i \in \{1, \dots, s\}.$$

Consider the one-step θ -scheme explained in [10, Chapter 4.1]. Using the notation from (8) at time t_i the following equation is to be solved for X^i .

$$(12) \quad \underbrace{\frac{1}{\Delta_t} A_t^f X^i (\rho_f I_m + D_3) + \frac{1}{\Delta_t} A_t^s X^i + \theta F(X^i)}_{=: F^i(X^i)} = \underbrace{\frac{1}{\Delta_t} A_t^f X^{i-1} (\rho_f I_m + D_3) + \frac{1}{\Delta_t} A_t^s X^{i-1} - (1 - \theta) F(X^{i-1}) + \theta B^i + (1 - \theta) B^{i-1}}_{=: B^i(X^{i-1})}$$

where $\theta \in [0, 1]$. $F^i(\cdot)$ contains only two sum terms more than $F(\cdot)$ from (8). At time t_i , both Algorithm 1 and Algorithm 3 can be applied to the quasi stationary problem (12) with $F^i(\cdot)$ instead of $F(\cdot)$ and the right hand side $B^i(X^{i-1})$.

5.3. Preconditioner. At all time steps the full matrix is given by

$$\begin{aligned} \mathcal{A}^t &:= \frac{1}{\Delta_t} (\rho_f I_m + D_3) \otimes A_t^f + \frac{1}{\Delta_t} I_m \otimes A_t^s \\ &\quad + \theta (I \otimes A_0 + D_1 \otimes A_1 + D_2 \otimes A_2 + \nu_f D_3 \otimes A_3). \end{aligned}$$

The mean-based preconditioner, similar to \mathcal{P}_T from subsection 3.2, is

$$\begin{aligned} \mathcal{P}_T^t &:= I \otimes P_T^t \text{ where} \\ P_T^t &:= \frac{1}{\Delta_t} (\rho_f + \bar{\rho}_f) A_t^f + \frac{1}{\Delta_t} A_t^s + \theta (A_0 + \bar{\mu}_s A_1 + \bar{\lambda}_s A_2 + \nu_f \bar{\rho}_f A_3). \end{aligned}$$

Even though the right hand side $B^i(X^{i-1})$ changes with every time step, the system matrix does not.

6. Theoretical Error Bounds. The convergence proofs of the GMRES method from [14, Theorem 35.2] and [12, Chapter 3.4] base on the fact that the residual of the l th GMRES iterate can be represented as a product of a polynomial in \mathcal{A} and the initial residual since the l th GMRES iterate is a linear combination of the start vector x_0 and the generating elements of \mathcal{K}_l . Also the error bound of the Chebyshev method in [3] relies on the fact that the residual of the l th Chebyshev iterate is such a product. But even if one considers Algorithm 1 and Algorithm 3 in a non preconditioned version multiplication with the system matrix \mathcal{A} is always disturbed due to the error induced by the truncation operator. The GMREST method minimizes over $\mathcal{K}_l^{\mathcal{T}}$, the truncated Krylov subspace, instead of \mathcal{K}_l . In the following subsection, the basis elements of $\mathcal{K}_l^{\mathcal{T}}$ are represented explicitly taking the truncation accuracy into consideration. Let x_l be the l th GMRES iterate, \hat{x}_l be the l th GMREST iterate. The error

$$\|x_l - \hat{x}_l\|_2$$

is then estimated based on these results. In relation to Krylov subspace methods, distortions initiated by matrix vector multiplication result in so called inexact Krylov methods and have been discussed in [13]. Iterative processes that involve truncation have been discussed in a general way in [5]. For the l th Chebyshev iterate x_l and the l th ChebyshevT iterate \hat{x}_l the error

$$\|x_l - \hat{x}_l\|_2$$

is bounded in the same way in the after next subsection. These bounds show how the truncation error is propagated iteratively in [Algorithm 1](#) and [Algorithm 3](#) if the machine precision error is neglected.

Remark 6.1. If $v \in \mathbb{R}^{Mm}$, $\mathcal{T}(v)$ addresses $\mathcal{T}(\text{vec}^{-1}(v))$. Thus for the sake of notation the truncation operator \mathcal{T} from [Definition 4.2](#) is regarded as a map

$$\mathcal{T} : \mathcal{R}^{Mm} \rightarrow T_R$$

and for $v \in \mathbb{R}^{Mm}$ $\mathcal{T}(v)$ addresses the full representation of the tensor in \mathbb{R}^{Mm} .

DEFINITION 6.2 (Truncation accuracy). *The truncation operator \mathcal{T} from [Definition 4.2](#) is said to have accuracy $\epsilon > 0$ if for any $x \in \mathbb{R}^{Mm}$*

$$\hat{x} := \mathcal{T}(x) = x + \mathcal{E}_{\hat{x}} \text{ with } \mathcal{E}_{\hat{x}} \in \mathbb{R}^{Mm} \text{ and } \|\mathcal{E}_{\hat{x}}\|_2 \leq \epsilon$$

holds. $\mathcal{E}_{\hat{x}}$ is the error induced by \mathcal{T} when x is truncated.

6.1. Matrix-vector Product Evaluation Accuracy. If a tensor is multiplied with a scalar or a matrix there is no truncation needed since the tensor rank does not grow. But the evaluation of $F(\cdot)$ from [\(8\)](#) involves 4 sum terms. After an evaluation of $F(\cdot)$ with a tensor as argument, the result has to be truncated. To keep complexity low for $\hat{X} \in T_R$, the sum $\mathcal{T}(F(\hat{X}))$, in practice, is truncated consecutively

$$\begin{aligned} \mathcal{T}(F(\hat{X})) &\equiv \mathcal{T}\left(\mathcal{T}(A_0\hat{X}_0 + A_1\hat{X}D_1) + A_2\hat{X}D_2 + \nu_f A_3\hat{X}D_3\right) \\ &= \mathcal{T}\left(\mathcal{T}(A_0\hat{X}_0 + A_1\hat{X}D_1 + A_2\hat{X}D_2 + \mathcal{E}_{\hat{F}_{s_1}}) + \nu_f A_3\hat{X}D_3\right) \\ &= \mathcal{T}(A_0\hat{X}_0 + A_1\hat{X}D_1 + A_2\hat{X}D_2 + \nu_f A_3\hat{X}D_3 + \mathcal{E}_{\hat{F}_{s_1}} + \mathcal{E}_{\hat{F}_{s_2}}) \\ &= F(\hat{X}) + \mathcal{E}_{\hat{F}_{s_1}} + \mathcal{E}_{\hat{F}_{s_2}} + \mathcal{E}_{\hat{F}_{s_3}}. \end{aligned}$$

In $\mathcal{T}(F(\cdot))$ are, if the number of summands in $F(\cdot)$ is $K \in \mathbb{N}$, a total of $K - 1$ truncations hidden. For a truncation accuracy of $\epsilon > 0$ we have

$$\|\mathcal{T}(F(\hat{X})) - F(\hat{X})\|_2 \leq (K - 1)\epsilon.$$

Since K is a small number usually not bigger than 4 we will neglect this detail and simply assume

$$\|\mathcal{T}(F(\hat{X})) - F(\hat{X})\|_2 \leq \epsilon$$

in the following. To make sure that the stated error bounds are still valid the truncation accuracy would be asked to, to be exact, less than $\frac{\epsilon}{K-1}$.

6.2. GMREST Error Bounds. Let x_l be the l th standard GMRES iterate, \hat{x}_l be the l th GMREST iterate. The first question that arises is how accurate is the truncated Krylov subspace $\mathcal{K}_l^{\mathcal{T}}$ from [subsection 4.2](#)? First we derive explicit representations of the unnormalized basis elements of $\mathcal{K}_l^{\mathcal{T}}$.

LEMMA 6.3 (Basis Representation of $\mathcal{K}_l^{\mathcal{T}}$). *Assume $\dim(\mathcal{K}_l^{\mathcal{T}}) = l$ and*

$$\hat{r}_0 = \mathcal{T}(\mathcal{P}^{-1}(b - \mathcal{A}x_0)) = r_0 + \mathcal{E}_{\hat{r}_0}.$$

Let the truncation operator $\mathcal{T}(\cdot)$ have accuracy $\epsilon > 0$. The unnormalized basis elements of $\mathcal{K}_l^\mathcal{T}$ are given by

\hat{r}_0 and

$$K^{\mathcal{T}k} := (\mathcal{P}^{-1}\mathcal{A})^k r_0 + (\mathcal{P}^{-1}\mathcal{A})^k \mathcal{E}_{\hat{r}_0} + \sum_{j=1}^k (\mathcal{P}^{-1}\mathcal{A})^{j-1} \mathcal{E}_{K^{\mathcal{T}k-j+1}}$$

for all $k \in \{1, \dots, l-1\}$.

Proof. (by induction)

For $k = 1$

$$\begin{aligned} K^{\mathcal{T}1} &= \mathcal{T}(\mathcal{P}^{-1}F(\hat{r}_0)) = \mathcal{T}(\mathcal{P}^{-1}\mathcal{A}\hat{r}_0) = \mathcal{T}(\mathcal{P}^{-1}\mathcal{A}(r_0 + \mathcal{E}_{\hat{r}_0})) \\ &= \mathcal{P}^{-1}\mathcal{A}r_0 + \mathcal{P}^{-1}\mathcal{A}\mathcal{E}_{\hat{r}_0} + \mathcal{E}_{K^{\mathcal{T}1}} \end{aligned}$$

and $k-1 \Rightarrow k$ since

$$\begin{aligned} K^{\mathcal{T}k} &= \mathcal{T}((\mathcal{P}^{-1}F)^k(\hat{r}_0)) = \mathcal{T}(\mathcal{P}^{-1}F(K^{\mathcal{T}k-1})) = \mathcal{T}(\mathcal{P}^{-1}\mathcal{A}K^{\mathcal{T}k-1}) \\ &= \mathcal{T}\left(\mathcal{P}^{-1}\mathcal{A}\left((\mathcal{P}^{-1}\mathcal{A})^{k-1}r_0 + (\mathcal{P}^{-1}\mathcal{A})^{k-1}\mathcal{E}_{\hat{r}_0} + \sum_{j=1}^{k-1} (\mathcal{P}^{-1}\mathcal{A})^{j-1} \mathcal{E}_{K^{\mathcal{T}k-j}}\right)\right) \\ &= (\mathcal{P}^{-1}\mathcal{A})^k r_0 + (\mathcal{P}^{-1}\mathcal{A})^k \mathcal{E}_{\hat{r}_0} + \sum_{j=1}^{k-1} (\mathcal{P}^{-1}\mathcal{A})^j \mathcal{E}_{K^{\mathcal{T}k-j}} + \mathcal{E}_{K^{\mathcal{T}k}} \\ &= (\mathcal{P}^{-1}\mathcal{A})^k r_0 + (\mathcal{P}^{-1}\mathcal{A})^k \mathcal{E}_{\hat{r}_0} + \sum_{j=1}^k (\mathcal{P}^{-1}\mathcal{A})^{j-1} \mathcal{E}_{K^{\mathcal{T}k-j+1}}. \quad \square \end{aligned}$$

Remark 6.4 (Truncation Error of \hat{r}_0). Consider the line

$$\text{Find } \hat{R} \text{ such that } P\hat{R} = \mathcal{T}(\hat{B} - F(\hat{X}))$$

of [Algorithm 1](#). In vector notation

$$(13) \quad \mathcal{P}\hat{r}_0 = \mathcal{T}(b - F(\hat{x}_0))$$

is solved for \hat{r}_0 . Usually the initial vector x_0 is chosen such that it can be represented by a tensor of low rank. So we assume $\hat{x}_0 = \mathcal{T}(x_0) = x_0$. If the linear system (13) is solved before truncation we have

$$\|\mathcal{E}_{\hat{r}_0}\|_2 \leq \epsilon.$$

But this is rarely implemented this way. In practice, the right side of (13) is truncated before the linear system is solved for \hat{r}_0 . In this case

$$\|\mathcal{E}_{\hat{r}_0}\|_2 \leq \epsilon \|\mathcal{P}^{-1}\|_2$$

holds. The following statements refer to the former case. The second case is always treated separately in a successive remark.

Remark 6.5 (Truncation of \hat{W} and Orthogonality). Consider the line

$$\hat{W} := \mathcal{T}(\hat{W} - H_{k,i}\hat{V}_k)$$

in [Algorithm 1](#). In the lemma above, this truncation is neglected. When the k th basis element \hat{V}_k is set up there are k extra additions involved due to this line. Let $\mathcal{E}_{\hat{W}}$ be the truncation error that occurs when this line is executed. To be precise, the basis elements are then given by

$$\hat{r}_0 \text{ and}$$

$$K^{\mathcal{T}^k} := (\mathcal{P}^{-1}\mathcal{A})^k r_0 + (\mathcal{P}^{-1}\mathcal{A})^k \mathcal{E}_{\hat{r}_0} + \sum_{j=1}^k (\mathcal{P}^{-1}\mathcal{A})^{j-1} \mathcal{E}_{K^{\mathcal{T}_{k-j+1}}} + k \mathcal{E}_{\hat{W}}$$

for $k \in \{1, \dots, l-1\}$.

The basis elements the exact Arnoldi iteration yields are also not

$$(14) \quad \{r_0, \mathcal{P}^{-1}\mathcal{A}, \dots, (\mathcal{P}^{-1}\mathcal{A})^{l-1}r_0\}.$$

But we neglect machine precision. If truncation is also neglected, the GMRES and the GMREST methods compute exactly the same basis elements for our problems. We incorporate the error made at the orthogonalization of the basis elements, in the truncated case, into $\mathcal{E}_{\hat{W}}$ and address by [\(14\)](#) the normalized basis elements that result from the Arnoldi iteration. In other words, we tacitly assume that the basis elements [\(14\)](#) of \mathcal{K}_l are orthonormal, write them in the representation [\(14\)](#) and incorporate the error we made at orthogonalization into $\mathcal{E}_{\hat{W}}$. This is just one result of the assumption that machine precision is zero.

LEMMA 6.6 (Error Bound for Truncated Basis Elements). *Let $\sigma_{\mathcal{P}} := \|\mathcal{P}^{-1}\mathcal{A}\|_2$. Under the assumptions of [Lemma 6.3](#) for*

$$e_k := \begin{cases} \|\hat{r}_0 - r_0\|_2 & \text{if } k = 0 \\ \|K^{\mathcal{T}^k} - (\mathcal{P}^{-1}\mathcal{A})^k r_0\|_2 & \text{if } k \in \{1, \dots, l-1\} \end{cases}$$

it holds that

$$e_k \leq \epsilon \sum_{j=1}^{k+1} \sigma_{\mathcal{P}}^{j-1} \text{ for } k \in \{0, \dots, l-1\}.$$

Proof. Case $k = 0$ is clear. For $k \geq 1$ [Lemma 6.3](#) can be used.

$$\begin{aligned} e_k &= \|(\mathcal{P}^{-1}\mathcal{A})^k \mathcal{E}_{\hat{r}_0} + \sum_{j=1}^k (\mathcal{P}^{-1}\mathcal{A})^{j-1} \mathcal{E}_{K^{\mathcal{T}_{k-j+1}}}\|_2 \\ &= \|(\mathcal{P}^{-1}\mathcal{A})^k \mathcal{E}_{\hat{r}_0} + \mathcal{E}_{K^{\mathcal{T}^k}} + \sum_{j=1}^{k-1} (\mathcal{P}^{-1}\mathcal{A})^j \mathcal{E}_{K^{\mathcal{T}_{k-j}}}\|_2 \\ &\leq \sigma_{\mathcal{P}}^k \epsilon + \epsilon + \sum_{j=1}^{k-1} \sigma_{\mathcal{P}}^j \epsilon = \epsilon \sum_{j=1}^{k+1} \sigma_{\mathcal{P}}^{j-1} \end{aligned} \quad \square$$

Remark 6.7. In the second case mentioned in [Remark 6.4](#)

$$e_k \leq \epsilon \left(\sum_{j=1}^k \sigma_{\mathcal{P}}^{j-1} + \|\mathcal{P}^{-1}\|_2 \sigma_{\mathcal{P}}^k \right)$$

holds for $k \in \{0, \dots, l-1\}$ since then

$$\|(\mathcal{P}^{-1}\mathcal{A})^k \mathcal{E}_{\hat{r}_0}\|_2 \leq \epsilon \|\mathcal{P}^{-1}\|_2 \sigma_{\mathcal{P}}^k$$

instead. The convention

$$\sum_{j \in \emptyset} \sigma_{\mathcal{P}}^{j-1} = 0$$

is used.

Remark 6.8. If one takes the line mentioned in [Remark 6.5](#) into consideration the bounds translate to

$$e_k \leq \epsilon \left(\sum_{j=1}^{k+1} \sigma_{\mathcal{P}}^{j-1} + k \right) \text{ or}$$

$$e_k \leq \epsilon \left(\sum_{j=1}^k \sigma_{\mathcal{P}}^{j-1} + \|\mathcal{P}^{-1}\|_2 \sigma_{\mathcal{P}}^k + k \right),$$

respectively, for $k \in \{0, \dots, l-1\}$.

The standard GMRES minimizes over the Krylov subspace \mathcal{K}_l . In terms of [Remark 6.5](#) the standard GMRES method finds coefficients $c_i \in \mathbb{R}$ for $i \in \{1, \dots, l\}$ such that

$$x_l = x_0 + c_1 r_0 + c_2 \mathcal{P}^{-1} \mathcal{A} r_0 + \dots + c_l (\mathcal{P}^{-1} \mathcal{A})^{l-1} r_0.$$

In the same way we can write

$$\hat{x}_l = \hat{x}_0 + d_1 \hat{r}_0 + d_2 K^{\mathcal{T}_1} + \dots + d_l K^{\mathcal{T}_{l-1}}$$

where the coefficients d_i for $i \in \{1, \dots, l\}$ refer to the coefficients found by the Arnoldi iteration in the GMREST method. This allows to state the following theorem.

THEOREM 6.9 (Approximation Error of GMREST). *Let x_l be the l th iterate of the standard GMRES method, \hat{x}_l be the l th iterate of the GMREST method. It holds*

$$\|\hat{x}_l - x_l\|_2 \leq \epsilon \sum_{j=1}^l \sum_{i=1}^{j+1} |d_j| \sigma_{\mathcal{P}}^{i-1} + \sum_{j=1}^l |c_j - d_j| + \epsilon(l-1).$$

Proof.

$$\begin{aligned} \|\hat{x}_l - x_l\|_2 &= \|\hat{x}_0 - x_0 + d_1 \hat{r}_0 - c_1 r_0 + d_2 K^{\mathcal{T}_1} - c_2 \mathcal{P}^{-1} \mathcal{A} r_0 + \dots \\ &\quad + d_l K^{\mathcal{T}_{l-1}} - c_l (\mathcal{P}^{-1} \mathcal{A})^{l-1} r_0\|_2 \\ &\leq |d_1| e_0 + |d_2| e_1 + \dots + |d_l| e_{l-1} + |c_1 - d_1| \underbrace{\|r_0\|_2}_{(\star)} \\ &\quad + |c_2 - d_2| \underbrace{\|\mathcal{P}^{-1} \mathcal{A} r_0\|_2}_{(\star)} + \dots + |c_l - d_l| \underbrace{\|(\mathcal{P}^{-1} \mathcal{A})^{l-1} r_0\|_2}_{(\star)} = (*) \end{aligned}$$

We assume that the standard GMRES method does an accurate orthogonalization of the Krylov subspace \mathcal{K}_l (see [Remark 6.5](#)). To be precise, the basis elements at (\star) are not the matrix products

$$(\mathcal{P}^{-1} \mathcal{A})^{j-1} r_0 \text{ for } j \in \{1, \dots, l\}$$

but this is neglected. By the elements (\star) we address the orthonormal basis elements of \mathcal{K}_l . They all have an Euclidean norm of 1. Therefore

$$\begin{aligned} (*) &= \sum_{j=1}^l (|d_j|e_{j-1} + |c_j - d_j|) \\ &\leq \epsilon \sum_{j=1}^l \sum_{i=1}^{j+1} |d_j| \sigma_{\mathcal{P}}^{i-1} + \sum_{j=1}^l |c_j - d_j| \end{aligned}$$

holds. The additional sum term $\epsilon(l-1)$ comes from the last successive sum in the method where the approximation is built. \square

Remark 6.10. For the second case mentioned in [Remark 6.4](#) the bound translates to

$$\|\hat{x}_l - x_l\|_2 \leq \epsilon \sum_{j=1}^l |d_j| \left(\sum_{i=1}^j \sigma_{\mathcal{P}}^{i-1} + \|\mathcal{P}^{-1}\|_2 \sigma_{\mathcal{P}}^j \right) + \epsilon(l-1).$$

Remark 6.11. In addition, with the error from the orthogonalization one obtains

$$\begin{aligned} \|\hat{x}_l - x_l\|_2 &\leq \epsilon \sum_{j=1}^l |d_j| \left(\sum_{i=1}^{j+1} \sigma_{\mathcal{P}}^{i-1} + j \right) + \sum_{j=1}^l |c_j - d_j| + \epsilon(l-1) \text{ and} \\ \|\hat{x}_l - x_l\|_2 &\leq \epsilon \sum_{j=1}^l |d_j| \left(\sum_{i=1}^j \sigma_{\mathcal{P}}^{i-1} + \|\mathcal{P}^{-1}\|_2 \sigma_{\mathcal{P}}^j + j \right) + \sum_{j=1}^l |c_j - d_j| \\ &\quad + \epsilon(l-1) \end{aligned}$$

in the other case.

6.3. ChebyshevT Error Bounds. In a similar way an error bound for the ChebyshevT method coded in [Algorithm 3](#) will be derived in this subsection. Let x_l denote the l th iterate of the standard Chebyshev method and \hat{x}_l denote the l th iterate of the ChebyshevT method.

Remark 6.12. The i th residual is given by the solution r_i to

$$\mathcal{P}r_i = b - \mathcal{A}x_i.$$

The truncation of r_i yields

$$\hat{r}_i = \mathcal{T}(r_i) = r_i + \mathcal{E}_{\hat{r}_i} \text{ with } \|\mathcal{E}_{\hat{r}_i}\|_2 \leq \epsilon$$

if the truncation operator \mathcal{T} is assumed to have accuracy ϵ . In analogy to [Remark 6.4](#) the two cases $\|\mathcal{E}_{\hat{r}_i}\|_2 \leq \epsilon$ (case 1) and $\|\mathcal{E}_{\hat{r}_i}\|_2 \leq \epsilon \|\mathcal{P}^{-1}\|_2$ (case 2) will be distinguished in the following.

The start vector and the right hand side are assumed to be of low rank, namely

$$\hat{x}_0 = \mathcal{T}(x_0) = x_0 \text{ and } \hat{b} = \mathcal{T}(b) = b.$$

In the same way as in the preceding subsection the norm $\|\hat{x}_l - x_l\|_2$ is to be estimated. $\mathcal{E}_{\hat{x}_l^a}$ denotes the total error

$$\mathcal{E}_{\hat{x}_l^a} := \hat{x}_l - x_l$$

not to be confused with $\mathcal{E}_{\hat{x}_l}$, the truncation error with norm ϵ that occurs when truncating \hat{x}_l . The iterative Chebyshev method is a three term recursion. Thus, the Chebyshev iterates itself can be represented by a recursive formula.

LEMMA 6.13 (Representation of the ChebyshevT Iterates). *Let the scalars*

$$\alpha_i, \beta_i \in \mathbb{R} \text{ for } i \in \{1, \dots, l\} \text{ and } \hat{\Phi}_i \in T_R \text{ for } i \in \{0, \dots, l\}$$

be given as defined in [Algorithm 3](#). Φ_i denote the non truncated full matrices corresponding to $\hat{\Phi}_i$ if [Algorithm 3](#) is applied and any truncation is neglected. Thus $\mathcal{T}(\hat{\Phi}_i) = \hat{\Phi}_i$ for $i \in \{0, \dots, l\}$. If

$$\hat{r}_0 = r_0 + \mathcal{E}_{\hat{r}_0}$$

it holds that

$$\begin{aligned} \mathcal{E}_{\hat{x}_0^a} &= 0, \\ \hat{x}_1 &= x_1 + \underbrace{\frac{1}{d}\mathcal{E}_{\hat{r}_0} + \mathcal{E}_{\hat{x}_1}}_{\mathcal{E}_{\hat{x}_1^a}}, \\ \hat{x}_2 &= x_2 + \underbrace{\mathcal{E}_{\hat{\Phi}_1} + \mathcal{E}_{\hat{x}_1^a} + \mathcal{E}_{\hat{x}_2} + \alpha_1(\mathcal{E}_{\hat{r}_1} - \mathcal{P}^{-1}\mathcal{A}\mathcal{E}_{\hat{x}_1^a}) + \frac{\beta_1}{d}\mathcal{E}_{\hat{r}_0}}_{\mathcal{E}_{\hat{x}_2^a}} \text{ and} \\ \hat{x}_l &= x_l + \mathcal{E}_{\hat{\Phi}_{l-1}} + \mathcal{E}_{\hat{x}_{l-1}^a} + \mathcal{E}_{\hat{x}_l} + \sum_{j=1}^{l-1} \alpha_j \left(\prod_{i=1}^{l-j-1} \beta_{i+j} \right) (\mathcal{E}_{\hat{r}_j} - \mathcal{P}^{-1}\mathcal{A}\mathcal{E}_{\hat{x}_j^a}) \\ &\quad + \left(\prod_{j=1}^{l-1} \beta_j \right) \frac{1}{d} \mathcal{E}_{\hat{r}_0} + \sum_{j=1}^{l-2} \left(\prod_{i=1}^{l-j-1} \beta_{i+j} \right) \mathcal{E}_{\hat{\Phi}_j} \text{ for } l \geq 3 \end{aligned}$$

where $\mathcal{E}_{\hat{x}_j^a} := \hat{x}_j - x_j$ for $j \in \{0, \dots, l\}$. We use the convention

$$\prod_{j \in \emptyset} \beta_j = 1.$$

If a truncation operator of accuracy $\epsilon > 0$ is used certainly $\|\mathcal{E}_{\hat{x}_i}\|_2 \leq \epsilon$ but not necessarily $\|\mathcal{E}_{\hat{x}_i^a}\|_2 \leq \epsilon$ holds for $i \in \{0, \dots, l\}$. The error induced by the truncation operator that truncates $\hat{\Phi}_i$ is denoted by $\mathcal{E}_{\hat{\Phi}_i}$ for $i \in \{0, \dots, l\}$.

Proof. $l = 1$

Provided that $\hat{x}_0 = x_0 = x_0 + \mathcal{E}_{\hat{x}_0^a} \Rightarrow \mathcal{E}_{\hat{x}_0^a} = 0$.

$$\hat{x}_1 = \mathcal{T}(\hat{x}_0 + \frac{1}{d}\hat{r}_0) = \mathcal{T}(x_0 + \frac{1}{d}r_0 + \frac{1}{d}\mathcal{E}_{\hat{r}_0}) = \underbrace{x_0 + \frac{1}{d}r_0}_{=x_1} + \underbrace{\frac{1}{d}\mathcal{E}_{\hat{r}_0} + \mathcal{E}_{\hat{x}_1}}_{=\mathcal{E}_{\hat{x}_1^a}}$$

$l = 2$

$$\begin{aligned} \hat{x}_2 &= \mathcal{T}(\hat{x}_1 + \hat{\Phi}_1) = \mathcal{T}(\hat{x}_1 + \mathcal{T}(\alpha_1\hat{r}_1 + \beta_1\hat{\Phi}_0)) = \mathcal{T}(\hat{x}_1 + \mathcal{T}(\alpha_1\hat{r}_1 + \frac{\beta_1}{d}\hat{r}_0)) \\ &= \mathcal{T}\left(x_1 + \mathcal{E}_{\hat{x}_1^a} + \mathcal{T}(\alpha_1(r_1 + \mathcal{E}_{\hat{r}_1} - \mathcal{P}^{-1}\mathcal{A}\mathcal{E}_{\hat{x}_1^a}) + \frac{\beta_1}{d}(r_0 + \mathcal{E}_{\hat{r}_0}))\right) = (\star) \end{aligned}$$

since

$$\hat{r}_1 = \mathcal{T}(\mathcal{P}^{-1}(b - \mathcal{A})\hat{x}_1) = \mathcal{T}(\mathcal{P}^{-1}(b - \mathcal{A}x_1 - \mathcal{A}\mathcal{E}_{\hat{x}_1^a})) = r_1 - \mathcal{P}^{-1}\mathcal{A}\mathcal{E}_{\hat{x}_1^a} + \mathcal{E}_{\hat{r}_1}.$$

The machine precision error is neglected so whether the preconditioner is applied before or after the truncation operator does not affect the term $\mathcal{P}^{-1}\mathcal{A}\mathcal{E}_{\hat{x}_1^a}$.

$$\begin{aligned} (\star) &= \mathcal{T}\left(\underbrace{x_1 + \alpha_1 r_1 + \frac{\beta_1}{d}r_0}_{=x_2} + \alpha_1(\mathcal{E}_{\hat{r}_1} - \mathcal{P}^{-1}\mathcal{A}\mathcal{E}_{\hat{x}_1^a}) + \frac{\beta_1}{d}\mathcal{E}_{\hat{r}_0} + \mathcal{E}_{\hat{x}_1^a} + \mathcal{E}_{\hat{\Phi}_1}\right) \\ &= x_2 + \alpha_1(\mathcal{E}_{\hat{r}_1} - \mathcal{P}^{-1}\mathcal{A}\mathcal{E}_{\hat{x}_1^a}) + \underbrace{\frac{\beta_1}{d}\mathcal{E}_{\hat{r}_0} + \mathcal{E}_{\hat{x}_1^a} + \mathcal{E}_{\hat{\Phi}_1} + \mathcal{E}_{\hat{x}_2}}_{=\mathcal{E}_{\hat{x}_2^a}} \end{aligned}$$

The proof for $l \geq 3$ goes by induction. For the initial step $l = 3$ we need

$$\begin{aligned} \hat{\Phi}_0 &= \frac{1}{d}\hat{r}_0 = \frac{1}{d}(r_0 + \mathcal{E}_{\hat{r}_0}) = \Phi_0 + \frac{1}{d}\mathcal{E}_{\hat{r}_0}, \\ \hat{\Phi}_1 &= \mathcal{T}(\alpha_1\hat{r}_1 + \beta_1\hat{\Phi}_0) = \mathcal{T}(\alpha_1 r_1 + \beta_1\Phi_0 + \alpha_1(\mathcal{E}_{\hat{r}_1} - \mathcal{P}^{-1}\mathcal{A}\mathcal{E}_{\hat{x}_1^a}) + \frac{\beta_1}{d}\mathcal{E}_{\hat{r}_0}) \\ &= \Phi_1 + \alpha_1(\mathcal{E}_{\hat{r}_1} - \mathcal{P}^{-1}\mathcal{A}\mathcal{E}_{\hat{x}_1^a}) + \frac{\beta_1}{d}\mathcal{E}_{\hat{r}_0} + \mathcal{E}_{\hat{\Phi}_1} \end{aligned}$$

and

$$\begin{aligned} \hat{\Phi}_2 &= \mathcal{T}(\alpha_2\hat{r}_2 + \beta_2\hat{\Phi}_1) = \alpha_2 r_2 + \beta_2\Phi_1 + \alpha_2(\mathcal{E}_{\hat{r}_2} - \mathcal{P}^{-1}\mathcal{A}\mathcal{E}_{\hat{x}_2^a}) \\ &\quad + \alpha_1\beta_2(\mathcal{E}_{\hat{r}_1} - \mathcal{P}^{-1}\mathcal{A}\mathcal{E}_{\hat{x}_1^a}) + \frac{\beta_1\beta_2}{d}\mathcal{E}_{\hat{r}_0} + \beta_2\mathcal{E}_{\hat{\Phi}_1} + \mathcal{E}_{\hat{\Phi}_2} \\ (15) \quad &= \Phi_2 + \alpha_2(\mathcal{E}_{\hat{r}_2} - \mathcal{P}^{-1}\mathcal{A}\mathcal{E}_{\hat{x}_2^a}) + \alpha_1\beta_2(\mathcal{E}_{\hat{r}_1} - \mathcal{P}^{-1}\mathcal{A}\mathcal{E}_{\hat{x}_1^a}) \\ &\quad + \frac{\beta_1\beta_2}{d}\mathcal{E}_{\hat{r}_0} + \beta_2\mathcal{E}_{\hat{\Phi}_1} + \mathcal{E}_{\hat{\Phi}_2}. \end{aligned}$$

Therefore

$$\begin{aligned} \hat{x}_3 &= \mathcal{T}(\hat{x}_2 + \hat{\Phi}_2) = \mathcal{T}(x_2 + \Phi_2 + \mathcal{E}_{\hat{x}_2^a} + \alpha_2(\mathcal{E}_{\hat{r}_2} - \mathcal{P}^{-1}\mathcal{A}\mathcal{E}_{\hat{x}_2^a}) + \alpha_1\beta_2(\mathcal{E}_{\hat{r}_1} - \mathcal{P}^{-1}\mathcal{A}\mathcal{E}_{\hat{x}_1^a}) \\ &\quad + \frac{\beta_1\beta_2}{d}\mathcal{E}_{\hat{r}_0} + \beta_2\mathcal{E}_{\hat{\Phi}_1} + \mathcal{E}_{\hat{\Phi}_2}) = x_3 + \mathcal{E}_{\hat{\Phi}_2} + \mathcal{E}_{\hat{x}_2^a} + \mathcal{E}_{\hat{x}_3} + \alpha_1\beta_2(\mathcal{E}_{\hat{r}_1} - \mathcal{P}^{-1}\mathcal{A}\mathcal{E}_{\hat{x}_1^a}) \\ &\quad + \alpha_2(\mathcal{E}_{\hat{r}_2} - \mathcal{P}^{-1}\mathcal{A}\mathcal{E}_{\hat{x}_2^a}) + \frac{\beta_1\beta_2}{d}\mathcal{E}_{\hat{r}_0} + \beta_2\mathcal{E}_{\hat{\Phi}_1} = x_3 + \mathcal{E}_{\hat{\Phi}_2} + \mathcal{E}_{\hat{x}_2^a} + \mathcal{E}_{\hat{x}_3} \\ &\quad + \sum_{j=1}^2 \alpha_j \left(\prod_{i=1}^{3-j-1} \beta_{i+j} \right) (\mathcal{E}_{\hat{r}_j} - \mathcal{P}^{-1}\mathcal{A}\mathcal{E}_{\hat{x}_j^a}) + \left(\prod_{j=1}^2 \beta_j \right) \frac{1}{d}\mathcal{E}_{\hat{r}_0} + \beta_2\mathcal{E}_{\hat{\Phi}_1}. \end{aligned}$$

To conclude $l-1 \rightarrow l$ we first prove that

$$\begin{aligned} \hat{\Phi}_{l-1} &= \Phi_{l-1} + \mathcal{E}_{\hat{\Phi}_{l-1}} + \sum_{j=1}^{l-1} \alpha_j \left(\prod_{i=1}^{l-j-1} \beta_{i+j} \right) (\mathcal{E}_{\hat{r}_j} - \mathcal{P}^{-1}\mathcal{A}\mathcal{E}_{\hat{x}_j^a}) \\ (16) \quad &\quad + \left(\prod_{j=1}^{l-1} \beta_j \right) \frac{1}{d}\mathcal{E}_{\hat{r}_0} + \sum_{j=1}^{l-2} \left(\prod_{i=1}^{l-j-1} \beta_{i+j} \right) \mathcal{E}_{\hat{\Phi}_j} \end{aligned}$$

under the assumption that this equation holds for $\hat{\Phi}_{l-2}$ on the left side. For $\hat{\Phi}_2$ this is true since from (15) we have that

$$\begin{aligned}\hat{\Phi}_2 &= \Phi_2 + \mathcal{E}_{\hat{\Phi}_2} + \sum_{j=1}^2 \alpha_j \left(\prod_{i=1}^{2-j} \beta_{i+j} \right) (\mathcal{E}_{\hat{r}_j} - \mathcal{P}^{-1} \mathcal{A} \mathcal{E}_{\hat{x}_j^a}) \\ &\quad + \left(\prod_{j=1}^2 \beta_j \right) \frac{1}{d} \mathcal{E}_{\hat{r}_0} + \beta_2 \mathcal{E}_{\hat{\Phi}_1}.\end{aligned}$$

The induction step for (16) goes as follows.

$$\begin{aligned}\hat{\Phi}_{l-1} &= \mathcal{T}(\alpha_{l-1} \hat{r}_{l-1} + \beta_{l-1} \hat{\Phi}_{l-2}) = \alpha_{l-1} r_{l-1} + \beta_{l-1} \Phi_{l-2} + \mathcal{E}_{\hat{\Phi}_{l-1}} \\ &\quad + \alpha_{l-1} (\mathcal{E}_{\hat{r}_{l-1}} - \mathcal{P}^{-1} \mathcal{A} \mathcal{E}_{\hat{x}_{l-1}^a}) + \beta_{l-1} \mathcal{E}_{\hat{\Phi}_{l-2}} \\ &\quad + \beta_{l-1} \sum_{j=1}^{l-2} \alpha_j \left(\prod_{i=1}^{l-j-2} \beta_{i+j} \right) (\mathcal{E}_{\hat{r}_j} - \mathcal{P}^{-1} \mathcal{A} \mathcal{E}_{\hat{x}_j^a}) \\ &\quad + \beta_{l-1} \left(\prod_{j=1}^{l-2} \beta_j \right) \frac{1}{d} \mathcal{E}_{\hat{r}_0} + \beta_{l-1} \sum_{j=1}^{l-3} \left(\prod_{i=1}^{l-j-2} \beta_{i+j} \right) \mathcal{E}_{\hat{\Phi}_j}\end{aligned}$$

$$\begin{aligned}
 &= \Phi_{l-1} + \mathcal{E}_{\hat{\Phi}_{l-1}} + \sum_{j=1}^{l-1} \alpha_j \left(\prod_{i=1}^{l-j-1} \beta_{i+j} \right) (\mathcal{E}_{\hat{r}_j} - \mathcal{P}^{-1} \mathcal{A} \mathcal{E}_{\hat{x}_j^a}) \\
 &+ \left(\prod_{j=1}^{l-1} \beta_j \right) \frac{1}{d} \mathcal{E}_{\hat{r}_0} + \sum_{j=1}^{l-2} \left(\prod_{i=1}^{l-j-1} \beta_{i+j} \right) \mathcal{E}_{\hat{\Phi}_j}
 \end{aligned}$$

With this it follows that

$$\begin{aligned}
 \hat{x}_l &= \mathcal{T}(\hat{x}_{l-1} + \hat{\Phi}_{l-1}) = x_l + \mathcal{E}_{\hat{\Phi}_{l-1}} + \mathcal{E}_{\hat{x}_{l-1}^a} + \mathcal{E}_{\hat{x}_l} \\
 &+ \sum_{j=1}^{l-1} \alpha_j \left(\prod_{i=1}^{l-j-1} \beta_{i+j} \right) (\mathcal{E}_{\hat{r}_j} - \mathcal{P}^{-1} \mathcal{A} \mathcal{E}_{\hat{x}_j^a}) \\
 &+ \left(\prod_{j=1}^{l-1} \beta_j \right) \frac{1}{d} \mathcal{E}_{\hat{r}_0} + \sum_{j=1}^{l-2} \left(\prod_{i=1}^{l-j-1} \beta_{i+j} \right) \mathcal{E}_{\hat{\Phi}_j}. \quad \square
 \end{aligned}$$

THEOREM 6.14 (ChebyshevT Approximation Error). *Let $\epsilon_R > 0$ be such that $\|\mathcal{E}_{\hat{r}_i}\|_2 \leq \epsilon_R$ for all $i \in \{1, \dots, l\}$ and $\sigma_{\mathcal{P}} := \|\mathcal{P}^{-1} \mathcal{A}\|_2$. Under the assumptions of [Lemma 6.13](#) the following error bounds hold for a truncation operator of accuracy $\epsilon > 0$.*

$$\begin{aligned}
 e_1 &:= \|\hat{x}_l - x_l\|_2 = \|\mathcal{E}_{\hat{x}_l^a}\|_2 \leq \epsilon + \frac{1}{|d|} \epsilon_R, \\
 e_2 &:= \|\hat{x}_2 - x_2\|_2 \leq (3 + |\alpha_1| \sigma_{\mathcal{P}}) \epsilon + \left(|\alpha_1| + \frac{1 + |\beta_1| + |\alpha_1| \sigma_{\mathcal{P}}}{|d|} \right) \epsilon_R \text{ and} \\
 e_l &:= \|\hat{x}_l - x_l\|_2 \leq (1 + |\alpha_{l-1}| \sigma_{\mathcal{P}}) e_{l-1} + \sum_{j=1}^{l-2} |\alpha_j| e_j \sigma_{\mathcal{P}} \prod_{i=1}^{l-j-1} |\beta_{i+j}| \\
 &+ \left(2 + \sum_{j=1}^{l-2} \prod_{i=1}^{l-j-1} |\beta_{i+j}| \right) \epsilon \\
 &+ \left(\sum_{j=1}^{l-1} |\alpha_j| \prod_{i=1}^{l-j-1} |\beta_{i+j}| + \frac{\prod_{j=1}^{l-1} |\beta_j|}{|d|} \right) \epsilon_R \text{ for } l \geq 3.
 \end{aligned}$$

Proof. $l = 1$

$$e_1 = \|\mathcal{E}_{\hat{x}_1^a}\|_2 \leq \epsilon + \frac{1}{|d|} \epsilon_R$$

$l = 2$

$$\begin{aligned}
 e_2 &= \|\mathcal{E}_{\hat{\Phi}_1} + \frac{1}{d} \mathcal{E}_{\hat{r}_0} + \mathcal{E}_{\hat{x}_1} + \mathcal{E}_{\hat{x}_2} + \alpha_1 (\mathcal{E}_{\hat{r}_1} - \mathcal{P}^{-1} \mathcal{A} (\frac{1}{d} \mathcal{E}_{\hat{r}_0} + \mathcal{E}_{\hat{x}_1})) + \frac{\beta_1}{d} \mathcal{E}_{\hat{r}_0}\|_2 \\
 &\leq \|\mathcal{E}_{\hat{\Phi}_1} + \mathcal{E}_{\hat{x}_1} + \mathcal{E}_{\hat{x}_2}\|_2 + |\alpha_1| \sigma_{\mathcal{P}} \|\mathcal{E}_{\hat{x}_1}\|_2 + |\alpha_1| \|\mathcal{E}_{\hat{r}_1}\|_2 + (1 + |\alpha_1| \sigma_{\mathcal{P}} + |\beta_1|) \|\frac{\mathcal{E}_{\hat{r}_0}}{d}\|_2 \\
 &\leq (3 + |\alpha_1| \sigma_{\mathcal{P}}) \epsilon + \left(|\alpha_1| + \frac{1 + |\beta_1| + |\alpha_1| \sigma_{\mathcal{P}}}{|d|} \right) \epsilon_R
 \end{aligned}$$

$l \geq 3$

$$\begin{aligned}
e_l &= \left\| \mathcal{E}_{\hat{\Phi}_{l-1}} + \mathcal{E}_{\hat{x}_{l-1}^a} + \mathcal{E}_{\hat{x}_l} + \sum_{j=1}^{l-1} \alpha_j \left(\prod_{i=1}^{l-j-1} \beta_{i+j} \right) (\mathcal{E}_{\hat{r}_j} - \mathcal{P}^{-1} \mathcal{A} \mathcal{E}_{\hat{x}_j^a}) \right. \\
&\quad \left. + \left(\prod_{j=1}^{l-1} \beta_j \right) \frac{1}{d} \mathcal{E}_{\hat{r}_0} + \sum_{j=1}^{l-2} \left(\prod_{i=1}^{l-j-1} \beta_{i+j} \right) \mathcal{E}_{\hat{\Phi}_j} \right\|_2 \\
&\leq \underbrace{\|\mathcal{E}_{\hat{x}_{l-1}^a}\|_2}_{=e_{l-1}} + |\alpha_{l-1}| \sigma_{\mathcal{P}} \|\mathcal{E}_{\hat{x}_{l-1}^a}\|_2 + \sum_{j=1}^{l-2} |\alpha_j| \sigma_{\mathcal{P}} \|\mathcal{E}_{\hat{x}_j^a}\|_2 \prod_{i=1}^{l-j-1} |\beta_{i+j}| \\
&\quad + \|\mathcal{E}_{\hat{\Phi}_{l-1}} + \mathcal{E}_{\hat{x}_l}\|_2 + \sum_{j=1}^{l-2} \|\mathcal{E}_{\hat{\Phi}_j}\|_2 \prod_{i=1}^{l-j-1} |\beta_{i+j}| \\
&\quad \sum_{j=1}^{l-1} |\alpha_j| \|\mathcal{E}_{\hat{r}_j}\|_2 \prod_{i=1}^{l-j-1} |\beta_{i+j}| + \frac{\prod_{j=1}^{l-1} |\beta_j|}{|d|} \|\mathcal{E}_{\hat{r}_0}\|_2 \\
&\leq (1 + |\alpha_{l-1}| \sigma_{\mathcal{P}}) e_{l-1} + \sum_{j=1}^{l-2} |\alpha_j| e_j \sigma_{\mathcal{P}} \prod_{i=1}^{l-j-1} |\beta_{i+j}| \\
&\quad + (2 + \sum_{j=1}^{l-2} \prod_{i=1}^{l-j-1} |\beta_{i+j}|) \epsilon \\
&\quad + \left(\sum_{j=1}^{l-1} |\alpha_j| \prod_{i=1}^{l-j-1} |\beta_{i+j}| + \frac{\prod_{j=1}^{l-1} |\beta_j|}{|d|} \right) \epsilon_R \quad \square
\end{aligned}$$

7. Numerical Evaluation of the Error Bounds. In algorithm and software implementations, the accuracy of a truncation operator depends on the truncation rank and not vice versa. If one chooses a rank R , the iterate of the GMREST or the ChebyshevT method is truncated to, the accuracy of the truncation operator is still unknown. Most truncation techniques like the HOSVD for hierarchical Tucker tensors ([4, Chapter 8.3 and Chapter 10.1.1]) or the TT-rounding for TT tensors ([9, Algorithm 1 and 2]) provide quasi optimality for tensors of order $d > 2$. For tensors of order $d = 2$ they even provide optimality in the sense that the result of the truncation of a matrix to rank R is indeed the best rank R approximation of the matrix. Nonetheless, since the singular value decay of the argument to be truncated is, in general, not known, the truncation operator will be simulated for a numerical evaluation of the error bounds of [Theorem 6.9](#) and [Theorem 6.14](#). Using the Matlab routine `rand()`, a vector

$$\tilde{z} \in \mathbb{R}^{Mm}$$

with entries that are uniformly distributed on the interval $(0, 1)$ is constructed first. The argument $x \in \mathbb{R}^{Mm}$ is then truncated using the truncation simulator

$$(17) \quad \mathcal{T}_s(x) := x + \frac{\epsilon}{\|\tilde{z}\|_2} \tilde{z}.$$

Of course, \tilde{z} is computed anew every time $\mathcal{T}_s(\cdot)$ is applied. For this subsection, all the computations are therefore made in the full format and whenever a truncation operator is applied, the truncation simulator $\mathcal{T}_s(\cdot)$ is evaluated. The main advantage of this strategy is that

$$\|\mathcal{T}_s(x) - x\|_2 = \epsilon \forall x \in \mathbb{R}^{Mm}.$$

A truncation operator based on the singular value decomposition does not provide such a reliable behavior. Let

$$\{\sigma_i\}_{i \in \{1, \dots, \min\{M, m\}\}}, \sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{\min\{M, m\}}$$

be the singular values of $\text{vec}^{-1}(x)$ and $\exists k \in \{1, \dots, \min\{M, m\} - 1\}$ such that

$$\sigma_k = 10^{-4} \text{ and } \sigma_{k+1} = 10^{-10}.$$

A truncation operator with accuracy $\epsilon = 10^{-5}$ based on the singular value decomposition would provide an approximation of x with accuracy 10^{-10} in this example.

7.1. GMREST Error Bound. We apply the preconditioner after the truncation operator so we use the second bound mentioned in [Remark 6.11](#) that reads

$$\begin{aligned} \|\hat{x}_l - x_l\|_2 &\leq \epsilon \sum_{j=1}^l |d_j| \left(\sum_{i=1}^j \sigma_{\mathcal{P}}^{i-1} + \|\mathcal{P}^{-1}\|_2 \sigma_{\mathcal{P}}^j + j \right) + \sum_{j=1}^l |c_j - d_j| \\ &\quad + \epsilon(l-1). \end{aligned}$$

This theoretical error bound is compared with

$$(18) \quad \|x_l - \hat{x}_l\|_2$$

where x_l denotes the l th GMRES iterate and \hat{x}_l the l th GMREST iterate. As just explained, everything is computed in the full format and every time a truncation is involved (which affects the GMREST iterate \hat{x}_l only) \mathcal{T}_s from [\(17\)](#) is evaluated. The 3d jetty from [subsection 8.1](#) is considered with a number of degrees of freedom of $M = 4095$ and a three parameter discretization with a total of $m = 8000$ parameter combinations as used in [subsection 8.2](#). We use the estimate $\sigma_{\mathcal{P}} \approx d + c$ with c, d from [subsection 8.4](#). In addition, the basis element error bound from [Remark 6.8](#) (the second case) is plotted for a truncation accuracy of $\epsilon = 10^{-12}$. If one starts with a matrix whose entries are all set to 1, the error bound [\(18\)](#) states that $\|x_{10} - \hat{x}_{10}\|_2$ is not bigger than $\approx 10^{-2}$, which can be seen in [Figure 1a](#). The reason for such a tolerant bound is that the first coefficients d_1, d_2, \dots are very big if the initial guess is bad. But if both methods are restarted with \hat{x}_6 as start matrix these coefficients become smaller as shown in [Figure 1b](#). Also the relative residual norm of the GMRES iterate, $\frac{\|B-F(x)\|_F}{\|B\|_F}$, and the one of the GMREST iterate, $\frac{\|B-F(\hat{x})\|_F}{\|B\|_F}$, are plotted. So even though $\|x_l - \hat{x}_l\|_2$ stagnates the truncated approach still provides nearly the same accuracy as the full approach does.

The dominating terms are

$$\epsilon \sum_{j=1}^l |d_j| \|\mathcal{P}^{-1}\|_2 \sigma_{\mathcal{P}}^j \text{ and } \epsilon \sum_{j=1}^l |d_j| \sum_{i=1}^j \sigma_{\mathcal{P}}^{i-1}.$$

As pointed out above, d_1, d_2, \dots are big for a bad initial guess. Then, in addition, ϵ can not compensate the (exponential) growth of $\sigma_{\mathcal{P}}^j$ for $j \in \{1, \dots, l\}$. Notice that $\sigma_{\mathcal{P}} \approx 1.6$ in this example. Since $\|\mathcal{P}^{-1}\|_2$ is very big, namely $\approx 3 \cdot 10^4$ in this example, the former of these two terms is bigger for the first 10 iterations.

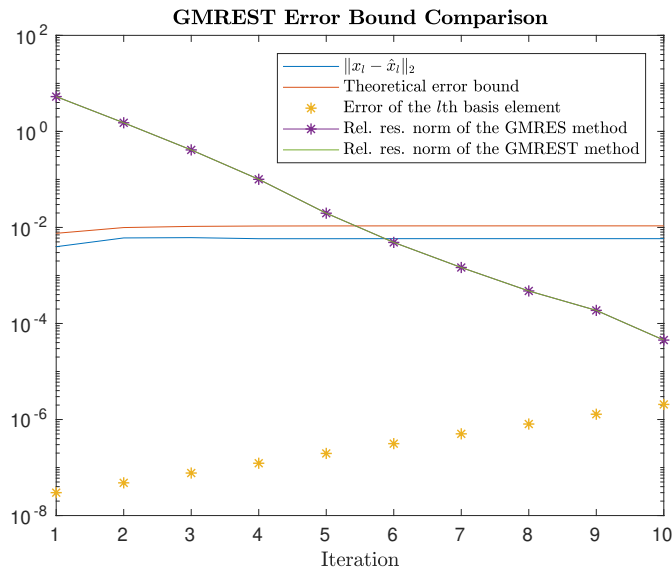
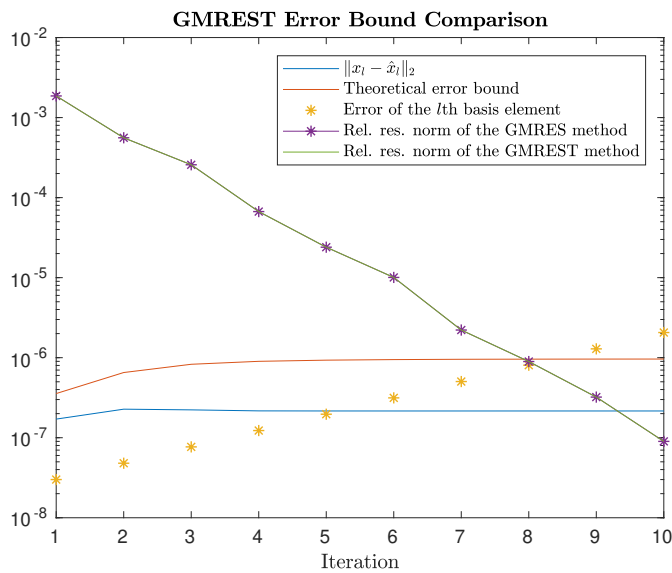
(a) $\epsilon = 10^{-12}$. All entries of the start matrix are set to 1.(b) $\epsilon = 10^{-12}$. \hat{x}_6 is used as start matrix.

FIG. 1. A numerical evaluation of the theoretical GMREST error bound.

7.2. ChebyshevT Error Bound. In this subsection, the approximation error from [Theorem 6.14](#) is numerically examined. Let x_l be the l th Chebyshev iterate and \hat{x}_l be the l th ChebyshevT iterate. For the ChebyshevT method, similar to the GMRES comparison, \mathcal{T}_s is used as truncation operator.

Remark 7.1. If preconditioned with \mathcal{P}_T , in theory,

$$\epsilon_R \leq \epsilon \|\mathcal{P}_T^{-1}\|_2$$

holds since the implementation used works like the case 2 mentioned in [Remark 6.12](#). But, due to the bad condition of the preconditioner and machine precision, in practice, ϵ_R is often bigger. The line

$$\text{Find } \hat{R}_i \text{ such that } P\hat{R}_i = \mathcal{T}(\hat{B} - F(\hat{X})).$$

is executed in every iteration in [Algorithm 3](#) which sometimes leads to real errors that are higher than the error bound from [Theorem 6.14](#). In contrast to this, the line

$$\text{Find } \hat{W} \in T_R \text{ such that } P\hat{W} = \mathcal{T}(F(\hat{V}_i)).$$

in [Algorithm 1](#) is less vulnerable. To circumvent this problem, for the error bound of [Theorem 6.14](#), the error ϵ_R is computed explicitly. This value instead of $\epsilon \|\mathcal{P}_T^{-1}\|_2$ is then used to compute the theoretical error bounds that are compared with the real errors.

We use the same configuration as used in the preceding subsection. Even though the theoretical error bound literally explodes, for $\epsilon = 10^{-12}$, the truncated method converges roughly as good as the non truncated method until iteration 10 for the 3d jetty model from [subsection 8.1](#) as shown in [Figure 2a](#). But the convergence of the ChebyshevT method deteriorates remarkably after 4 iterations for $\epsilon = 10^{-6}$ if compared to the full approach (see [Figure 2b](#)).

The two terms in the error bound that are not multiplied with ϵ are

$$(19) \quad (1 + |\alpha_{l-1}| \sigma_{\mathcal{P}}) e_{l-1} \text{ and}$$

$$(20) \quad \sum_{j=1}^{l-2} |\alpha_j| e_j \sigma_{\mathcal{P}} \prod_{i=1}^{\overbrace{l-j-1}^{(*)}} |\beta_{i+j}|.$$

The coefficients β_i in the Chebyshev method have norms that are smaller than 1. This becomes clear if one considers the recursive computation formula for the Chebyshev polynomials (see [\[8\]](#) (2.4)) evaluated at $\frac{d}{c}$ with $|c| < |d|$. The coefficients β_i are then given as a fraction where the numerator has a norm that is smaller than the denominator. The product (\star) becomes smaller the higher the iteration number is and therefore the term [\(20\)](#) becomes negligibly small, at least if it is compared with the term [\(19\)](#). For our configuration, $c = 0.6$. Hence the coefficients α_i are bigger than 1 on the other hand (see [\[8\]](#) (2.24)). For $\sigma_{\mathcal{P}} \approx c + d = 1.6$

$$1 + |\alpha_{l-1}| \sigma_{\mathcal{P}} \geq 2.6.$$

This explains why the first term in [Theorem 6.14](#), the term [\(19\)](#), dominates the error bound and makes it explode.

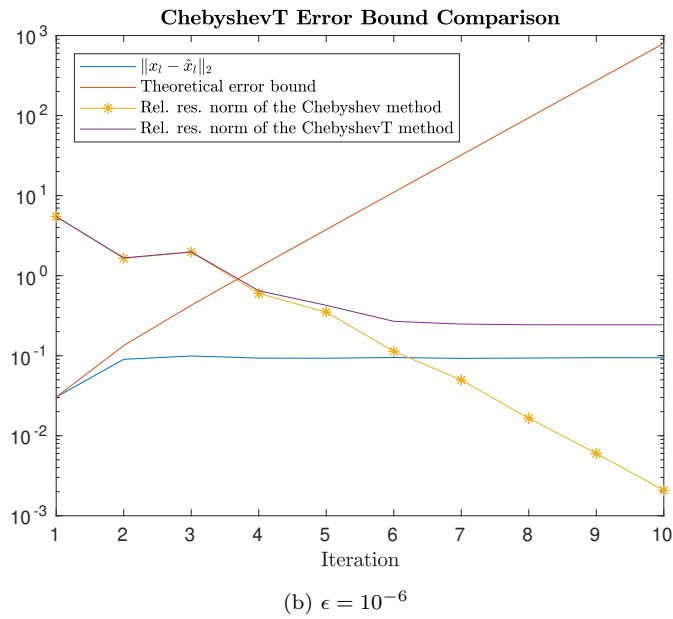
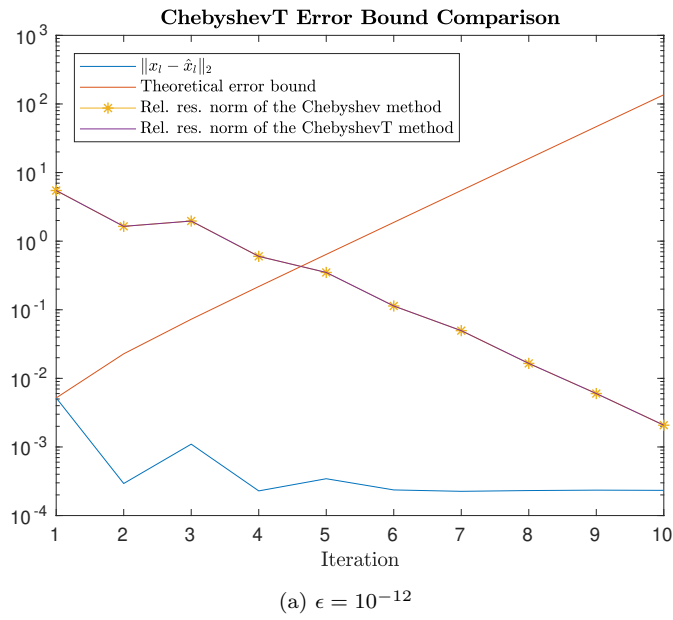


FIG. 2. A numerical evaluation of the error bound from *Theorem 6.14*. All entries of the start matrix are set to 1.

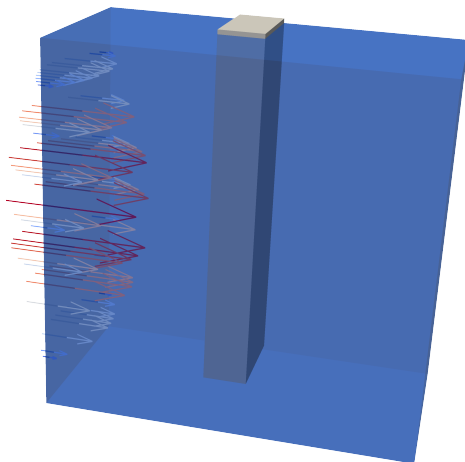


FIG. 3. The initial configuration of the jetty where the Dirichlet boundary conditions simulate an inflow from left.

8. Numerical Examples.

8.1. A Three Dimensional Jetty in a Channel. The geometric configuration of a 3d jetty in a channel is given by

$$\Omega := (0, 8) \times (0, 8) \times (0, 4), S := (3, 4) \times (0, 8) \times (0, 2) \text{ and } F := \Omega \setminus \bar{S}.$$

With the velocity

$$\begin{pmatrix} v_1 \\ v_2 \\ v_3 \end{pmatrix} \in \mathbb{R}^3, \text{ the deformation } \begin{pmatrix} u_1 \\ u_2 \\ u_3 \end{pmatrix} \in \mathbb{R}^3 \text{ and coordinates } (x, y, z) \in \bar{\Omega}$$

the left Dirichlet inflow is given by

$$v = \begin{pmatrix} \frac{1}{2}y(8-y)z(4-z) \\ 0 \\ 0 \end{pmatrix} \text{ if } x = 0.$$

The geometric configuration is illustrated in [Figure 3](#). On the right, for $x = 8$, do nothing boundary outflow conditions [[10](#), Chapter 2.4.2] hold. The surface is at $y = 8$. There, v_2 and u_2 vanish only. Everywhere else on $\partial(\Omega)$ the velocity and the deformation fulfill zero Dirichlet boundary conditions.

To stabilize the Stokes equations on the fluid, stabilized Stokes elements [[10](#), Lemma 4.47] are used.

8.2. Three Parameter Discretization. Problem [\(4\)](#) is discretized with respect to

- 20 shear moduli $\mu_s^{i_1} \in [30000, 50000]$,
- 20 first Lamé parameters $\lambda_s^{i_2} \in [100000, 200000]$ and
- 20 fluid densities $\rho_f^{i_3} \in [50, 200]$.

The kinematic fluid viscosity is fixed to $\nu_f = 0.01$. The shear modulus and first Lamé parameter ranges cover solids with Poisson ratios between $\frac{1}{3}$ (e.g. concrete)

and 0.43478 (e.g. clay). The total number of equations is $m = 20^3 = 8000$ and the number of degrees of freedom is $M = 192423$.

In the following computations, Matlab 2017b on a CentOS 7.6.1810 64bit with 2 AMD EPYC 7501 and 512GB of RAM is used. The htucker Matlab toolbox [7] is used to realize the Tucker format T_R . The preconditioners are decomposed into a permuted LU decomposition using the Matlab builtin command `lu()`. All methods start with a start matrix whose entries are all set to 1.

8.3. GMREST. A standard GMRES approach is compared with the GMRESTR method from [Algorithm 2](#).

By standard GMRES approach the standard GMRES method applied to $m = 8000$ different equations of the form (5) is meant. It is once restarted after 8 iterations so it uses a total of 16 iterations per equation. For all 8000 separate standard GMRES methods 5 preconditioners given by

$$(21) \quad A_0 + (D_1)_{i,i}A_1 + (D_2)_{i,i}A_2 + (D_3)_{i,i}A_3 \text{ for } i \in \{800, 2400, 4000, 5600, 7200\}$$

are set up where the diagonal matrices $\{D_j\}_{j \in \{1,2,3\}}$ are the ones from (7).

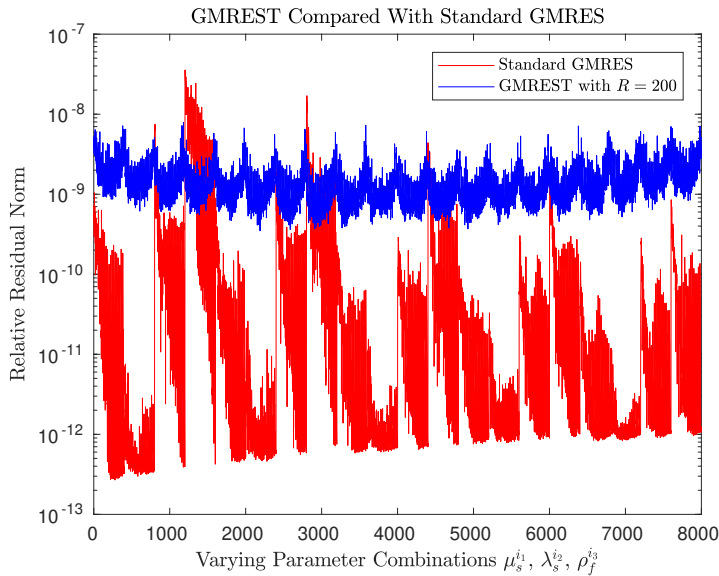


FIG. 4. The standard GMRES method applied to 8000 separate equations (residual norms in red) is compared with the GMREST method (residual norms in blue).

The GMRESTR method uses 6 iterations per restart and is restarted 3 times. The mean based preconditioner \mathcal{P}_T is used. The times to compute the preconditioners (one in the case of GMRESTR, 5 in the case of standard GMRES) can be found in the column Precon. in [Table 1](#). The method itself took the time that is listed in the column Comp. and the column Total is then the sum of these times. Both methods result in 8000 approximations. Each of these approximations (x axis) then provide a certain accuracy (y axis) that is plotted in [Figure 4](#). The standard GMRES method applied to 8000 equations in this way provides accuracies that are plotted in red within about 5432 minutes. The approximations the GMRESTR method provides have accuracies that are plotted in blue. The GMRESTR method took only about

TABLE 1
GMRESTR Compared With Standard GMRES

Method	Approx. Storage	Computation Times (in Minutes)		
		Precon.	Comp.	Total
GMRESTR ($R = 200$)	$O[(M + m + R)R]$ $\approx 306.12\text{MB}$	1.24	179.88	181.12
Standard GMRES (8000 times)	$O(Mm)$ $\approx 11744.56\text{MB}$	6.63	5426.23	5432.86

181 minutes to compute these approximations as one can see in [Table 1](#). Also the storage that is needed to store the approximation varies significantly. The rank 200 approximation, in the Tucker format, requires only about 306MB whereas the full matrix requires about 11744MB.

8.4. ChebyshevT. Before the Chebyshev method can be applied, the extreme eigenvalues of $\mathcal{P}_T^{-1}\mathcal{A}$ have to be estimated as explained in subsection 4.3. An estimation of Λ_{\max} and Λ_{\min} from (10) involves the estimation of extreme eigenvalues for m different matrices if the representation (9) is considered. But we restrict to an estimation of

$$\bar{\Lambda}_{\max} = \max_{\substack{i_1 \in \{1, m_1\} \\ i_2 \in \{1, m_2\} \\ i_3 \in \{1, m_3\}}} \Lambda(Bl(i_1, i_2, i_3)) \approx \Lambda_{\max} \text{ and } \bar{\Lambda}_{\min} = \min_{\substack{i_1 \in \{1, m_1\} \\ i_2 \in \{1, m_2\} \\ i_3 \in \{1, m_3\}}} \Lambda(Bl(i_1, i_2, i_3)) \approx \Lambda_{\min}.$$

With the mean based preconditioner \mathcal{P}_T , this leads to $d = 1$ and $c = 0.6$ in this configuration. The time needed to compute $\bar{\Lambda}_{\max}$ and $\bar{\Lambda}_{\min}$ on a coarse grid with a number of degrees of freedom of $M = 735$ is listed in the column Est. in Table 2.

In the same manner as in the preceding subsection, for a comparison, a standard Chebyshev approach is applied to 8000 equations of the form (5). The standard

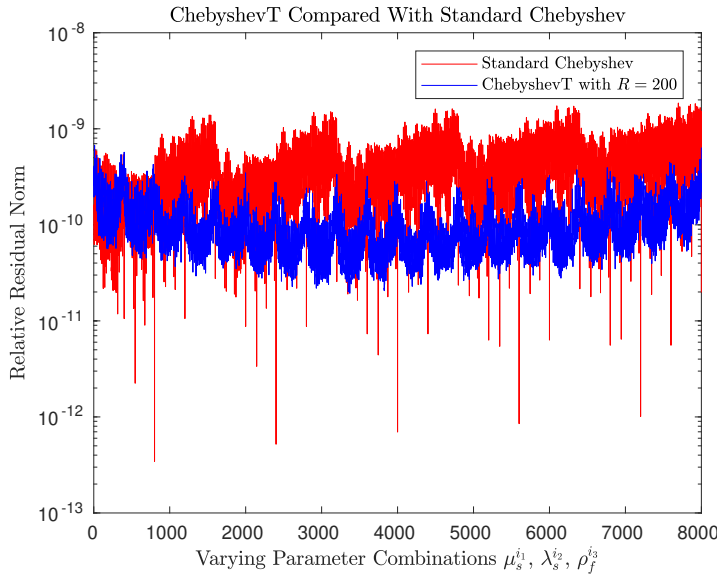


FIG. 5. The standard Chebyshev method applied to 8000 separate equation (residual norms in red) is compared with the ChebyshevT method (residual norms in blue).

TABLE 2
ChebyshevT Compared With Standard Chebyshev

Method	Storage needed by Approx.	Computation Times (in Minutes)			
		Est.	Precon.	Comp.	Total
ChebyshevT ($R = 200$)	$O[(M + m + R)R]$ $\approx 306.12\text{MB}$	0.013	1.24	177.99	179.243
Standard Chebyshev (8000 times)	$O(Mm)$ $\approx 11744.56\text{MB}$	0.013	6.63	5490.85	5497.493

Chebyshev method uses 20 iterations at each equation and, in total, the same 5 preconditioners (21) as the standard GMRES uses. The ChebyshevT method iterates, in total, 24 times and uses \mathcal{P}_T , the mean based preconditioner. The ChebyshevT

method is restarted 3 times with 6 iterations per restart. Compared to this, 24 iterations without restart take about the same time but provide approximation accuracies that are slightly worse.

9. Comparison With the Bi-CGstab method. Another method that also works for non symmetric matrices is the Bi-CGstab method [15]. It was not considered in the first place because it can break down under some circumstances as explained in [2, Chapter 2.3.8]. The preconditioned truncated variant similar to [6, Algorithm 3] but strictly based on [15] is compared with the GMRESTR and the ChebyshevT method. The truncated Bi-CGstab method is applied with 6 iterations per restart. If once restarted, in total, the method iterates 12 times. The resulting approximation accuracy is indeed better than the one obtained when iterating 12 times directly without any restart.

TABLE 3
Computation Time Comparison of the Truncated Approaches

Method (R=200)	Computation Times (in Minutes)			
	Est.	Precon.	Comp.	Total
ChebyshevT	0.013	1.24	177.99	179.243
GMREST	-	1.24	179.88	181.12
Truncated Bi-CGstab	-	1.24	302.94	304.18

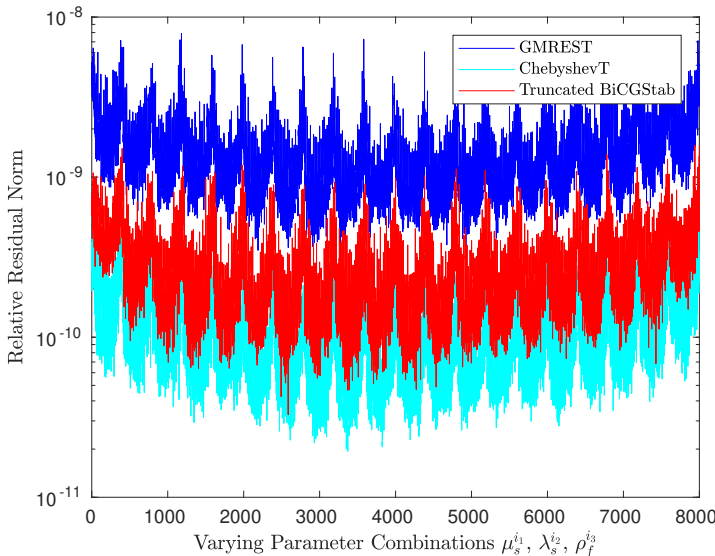


FIG. 6. The approximation accuracies for the GMREST (blue), the ChebyshevT (cyan) and the truncated Bi-CGstab (red).

To avoid early stagnation, the residual at step i is computed directly

$$\hat{R}_i = \mathcal{T}(\hat{B} - F(\hat{X})).$$

10. Conclusions. The truncated methods discussed in this paper provide approximations with relative residual norms smaller than 10^{-8} within less than a twen-

tieth of the time needed by the correspondent standard approaches that solve the m equations individually. This raises the question how these methods perform when applied to nonlinear problems.

Since the truncation error affects, in addition to the machine precision error, the accuracy of the Arnoldi orthogonalization, the GMREST method should preferably be applied in a restarted version. Mostly, the ChebyshevT method is a bit faster and a bit more accurate than the GMREST method. But the main disadvantage of the ChebyshevT method is that the ellipse that contains the eigenvalues of $\mathcal{P}^{-1}\mathcal{A}$ described by the foci $d \pm c$ has to be approximated newly every time the parameter configuration changes. In this matter, the GMREST method can be seen as a method that is a bit more flexible if compared to the ChebyshevT method.

Also the ChebyshevT and the truncated BiCGstab methods can and preferably should be applied in a restarted manner. If not restarted, the methods stagnate after a few iterations already. The reason is a numerical issue initiated by the bad condition of the mean-based preconditioner.

There is still investigation needed regarding the error bounds. If the GMREST method is applied, the coefficients c_j are not known. The ChebyshevT bound is rather of theoretical nature. The method seems to converge too fast such that the truncation error does not really play a role in the cases examined.

Acknowledgments. Funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) - 314838170, GRK 2297 MathCoRe.

REFERENCES

- [1] J. BALLANI AND L. GRASEDYCK, *A projection method to solve linear systems in tensor format*, Numer. Linear Algebra Appl., 20 (2013), pp. 27–43.
- [2] R. BARRETT, M. BERRY, T. F. CHAN, AND ET AL., *Templates for the solution of linear systems: building blocks for iterative methods*, Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 1994.
- [3] D. CALVETTI, G. H. GOLUB, AND L. REICHEL, *An adaptive Chebyshev iterative method for non-symmetric linear systems based on modified moments*, Numer. Math., 67 (1994), pp. 21–40.
- [4] W. HACKBUSCH, *Tensor spaces and numerical tensor calculus*, vol. 42 of Springer Series in Computational Mathematics, Springer, Heidelberg, 2012.
- [5] W. HACKBUSCH, B. N. KHOROMSKIJ, AND E. E. TYRTYSHNIKOV, *Approximate iterations for structured matrices*, Numer. Math., 109 (2008), pp. 365–383.
- [6] D. KRESSNER AND C. TOBLER, *Low-rank tensor Krylov subspace methods for parametrized linear systems*, SIAM J. Matrix Anal. Appl., 32 (2011), pp. 1288–1316.
- [7] ———, *Algorithm 941: htucker—a Matlab toolbox for tensors in hierarchical Tucker format*, ACM Trans. Math. Software, 40 (2014), pp. Art. 22, 22.
- [8] T. A. MANTEUFFEL, *The Tchebychev iteration for nonsymmetric linear systems*, Numer. Math., 28 (1977), pp. 307–327.
- [9] I. V. OSELEDETS, *Tensor-train decomposition*, SIAM J. Sci. Comput., 33 (2011), pp. 2295–2317.
- [10] T. RICHTER, *Fluid-structure interactions*, vol. 118 of Lecture Notes in Computational Science and Engineering, Springer, Cham, 2017. Models, analysis and finite elements.
- [11] Y. SAAD, *Iterative methods for sparse linear systems*, Society for Industrial and Applied Mathematics, Philadelphia, PA, second ed., 2003.
- [12] Y. SAAD AND M. H. SCHULTZ, *GMRES: a generalized minimal residual algorithm for solving nonsymmetric linear systems*, SIAM J. Sci. Statist. Comput., 7 (1986), pp. 856–869.
- [13] V. SIMONCINI AND D. B. SZYLD, *Theory of inexact Krylov subspace methods and applications to scientific computing*, SIAM J. Sci. Comput., 25 (2003), pp. 454–477.
- [14] L. N. TREFETHEN AND D. BAU, III, *Numerical linear algebra*, Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 1997.
- [15] 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 (1992), pp. 631–644.