# THE LANCZOS METHOD FOR PARAMETERIZED SYMMETRIC LINEAR SYSTEMS WITH MULTIPLE RIGHT-HAND SIDES*

KARL MEERBERGEN† AND ZHAOJUN BAI‡

**Abstract.** The solution of linear systems with a parameter is an important problem in engineering applications, including structural dynamics, acoustics, and electronic circuit simulations, and in related model order reduction methods such as Padé via Lanczos. In this paper, we present a Lanczos-based method for solving parameterized symmetric linear systems with multiple right-hand sides. We show that for this class of applications, a simple deflation method can be used.

**Key words.** parameterized linear system, Lanczos method, recycling Ritz vectors

**AMS subject classifications.** 65F15, 65F50

**DOI.** 10.1137/08073144X

**1. Introduction.** In engineering applications including structural dynamics and acoustics, the computation of the frequency response function of a vibrating system over a given frequency range $\Omega = [\omega_{\min}, \omega_{\max}]$ can be a time-consuming operation. In applications on closed domains without damping, the frequency response function often is the solution of the parameterized linear system

$$(1.1) \qquad Ax = f \quad \text{with} \quad A = K - \omega^2 M,$$

where $K$ and $M$ are large and sparse real symmetric matrices, and $M$ is symmetric positive definite. We have chosen $\omega$ to be the frequency, but it could also be the angular frequency, the wave number, or the characteristic (dimensionless) wave number. The frequency range $\Omega$ is discretized into the set $\{\omega_1, \omega_2, \ldots, \omega_m\}$ where $m$ can be of the order 100 or 1000. This solution process is called *frequency sweeping*. Since $A$ is large, the solution of (1.1) is expensive when it is solved for the frequency points $\omega = \omega_1, \omega_2, \ldots, \omega_m$ independently. In this paper, we study frequency sweeping with multiple right-hand sides; i.e., $f$ can take different values $f_1, f_2, \ldots, f_s$ in (1.1).

A number of solution methods to this parameterized linear system have been proposed in the literature. The most famous approach in engineering is undoubtedly the modal superposition method. References can be found in textbooks on engineering, e.g., [10]. The method projects the right-hand sides and solution vectors on a basis of eigenvectors of the underlying eigenvalue problem

$$(1.2) \qquad Ku = \lambda M u,$$

where $\lambda \in [\lambda_{\min}, \lambda_{\max}]$, with $\lambda_{\min} \ll \omega_{\min}^2$ and $\lambda_{\max} \gg \omega_{\max}^2$. This method is usually experienced as efficient when the eigenvectors and eigenvalues are available, since (1.1)

is transformed into a diagonal linear system. The practical problem is that it is not always clear how to choose $\lambda_{\min}$ and $\lambda_{\max}$. For example, when $\omega_{\min} = 0$, we could use $\lambda_{\min} = 0$ and $\lambda_{\max} = \eta \omega_{\max}^2$ with $\eta \in [2, 10]$. The eigenvalues and eigenvectors of (1.2) are in practice computed by the (block) Lanczos method; see, e.g., [17] or the automated multilevel substructuring (AMLS) method [6, 5, 7] which is advocated for very large scale problems.

Another technique that received quite some attention is the *parameterized Lanczos* method. It is an iterative method for the preconditioned system

$$(1.3) \qquad (K - \sigma M)^{-1} A x = (K - \sigma M)^{-1} f \ .$$

The solution vector $x$ appears to be a truncated vector Padé series with interpolation point $\sigma$ [13, 2, 3]. See also [30, 22] in the context of parameterized iterative linear system solvers and [23] for Rayleigh damping. This method is related to the Lanczos eigenvalue solver and the conjugate gradient method. The limitation of the method is that the right-hand side vector $f$ should not have a spatial dependency on $\omega$; i.e., $f = \widetilde{f} \phi(\omega)$ where $\widetilde{f}$ does not depend on $\omega$ and $\phi$ is a scalar function. In the engineering literature, this is called the Ritz vector technique [33]. The connection between the Lanczos method and a vector Padé series is important in applications, since the frequency response function is a rational function with the eigenvalues as poles. The approximation preferably respects this rational nature. This is one reason why the Lanczos method is preferred to the minimal residual (MINRES) method [22].

In contrast to the Lanczos method, the modal superposition method does not pose a condition on the right-hand side. The difficulty with the modal superposition method is that a relatively large number of eigenvectors may be required with a relatively high precision. In [5, 19], the AMLS frequency sweeping method is studied which makes a combination of modal superposition and a stationary iterative method. The method uses only the modes corresponding to the frequency range of interest; i.e., eigenvalues outside this interval are not used. Superposition on this reduced set of eigenvectors does produce the peaks in the frequency response function, but the zeros are wrong. Therefore, the AMLS frequency sweeping method uses an iterative method to improve the solution obtained by superposition on a reduced modal basis. The iterative method uses the eigenbasis as a preconditioner. Such preconditioner is also called deflation preconditioning and is related to augmented Krylov subspace methods [28]. It can be viewed as an iterative method on the subspace orthogonal to the given eigenvectors. Since the eigenvalues play an important role in the convergence of iterative methods, eliminating the eigenvalues that hinder convergence can be effective indeed.

The use of deflation preconditioning and recycling Ritz vectors is not new in the context of iterative linear system solvers [8, 28, 26, 32, 31]. In [28], a sequence of linear systems $\{A x_k = f_k \text{ for } k = 1, 2, \ldots\}$ is solved. The Ritz vectors of $A$, computed by the conjugate gradient method for solving the first system, are used in a deflation preconditioner for the second. In the Init-CG method, the recycled Ritz vectors are only deflated from the right-hand side; i.e., no deflation preconditioner is used. See [15] for an overview of methods. Recycling subspaces for linear systems with a parameter was recently introduced in [18, 11]. In both papers, the Ritz vectors for the solution of the first right-hand side are recycled for preconditioning the systems with the remaining right-hand sides. In [18], the generalized conjugate residual orthogonal (GCRO) method was extended to linear systems with a parameter. The authors prefer the situation where eigenvectors are recycled (or Ritz vectors with small residual norm), since this simplifies the method. In [11], recycling was proposed for the

generalized minimal residual (GMRES) method. The residual vectors for the additional right-hand sides are not parallel to the GMRES residual vector, which creates difficulties for recycling. A remedy was proposed for this problem.

There are important differences with this paper and the cited work. We assume it is allowed to factor a large sparse matrix. This is a reasonable assumption in vibration problems, since direct methods are commonly used in the mechanical and civil engineering communities. For improving AMLS frequency sweeping, diagonal linear systems are solved. Another difference with the cited work is that the number of iterations is typically low for frequency sweeping, i.e., typically less than a hundred.

Suppose that for a given right-hand side the preconditioned residual norms of (1.3) are small; i.e., the solutions are accurate for $\omega \in \Omega$. The key observation in the current paper is that the recycled Ritz vectors then have small residual norms, which allows for efficient deflation and recycling.

Our contributions can be summarized as follows. First, we show that the parameterized Lanczos method converges quickly when Ritz vectors are recycled. This allows us to propose a simple algorithm. Second, we show that the residual norms of the Ritz pairs associated with the interval $\Omega^2 = [\omega_{\min}^2, \omega_{\max}^2]$ computed by the Lanczos method are small. Third, we show a connection with the Padé via Lanczos method [13, 2, 3] for the model order reduction problem. If the deflated Ritz pairs have zero residual norms, we have an exact vector Padé approximation.

The paper is organized as follows. In section 2, we introduce a numerical method that applies the parameterized Lanczos method to the deflated right-hand side (as in Init-CG). In section 3, we perform a spectral analysis and we show spectral properties of the deflated linear system. In section 4, we discuss how Ritz pairs can be computed using the Lanczos method and how accurate they are. Section 5 presents a practical procedure for solving (1.1) with multiple right-hand sides recycling Ritz vectors from the first right-hand side. Section 6 shows numerical examples for applications from structural engineering acoustics and one academic example with multiple eigenvalues, which shows slower convergence. We close the paper with concluding remarks.

We summarize the used notation. The interval of $\omega$'s for which $x$ needs to be computed is denoted by $\Omega = [\omega_{\min}, \omega_{\max}]$. In our applications, $\omega_{\min} \geq 0$. We have also defined $\Omega^2 = [\omega_{\min}^2, \omega_{\max}^2]$. The transpose is denoted by $x^T$. The $M$ norm $\|x\|_M$ is defined as the induced norm from the $M$ inner product: $\sqrt{x^T M x}$.

**2. Deflation in parameterized linear systems.** In this section, we explain the ideas of deflation for solving (1.1). We start from the viewpoint of rational approximation, since $x$ is a rational function in $\omega^2$. Then, we discuss deflation of a part of the spectrum. This leads to a linear system which will be solved iteratively.

**2.1. Rational function splitting.** Let

$$(2.1) \qquad\qquad\qquad KU = MU\Lambda$$

be an eigendecomposition of (1.2), where $\Lambda$ is a diagonal matrix with diagonal elements $\lambda_j$, $j = 1, \ldots, n$, and $U^T MU = I$. By (2.1), we have that $M = U^{-T}U^{-1}$ and $K = U^{-T}\Lambda U^{-1}$. Therefore, the solution $x$ of (1.1) can be written as

$$(2.2) \qquad\qquad x = U(\Lambda - \omega^2 I)^{-1}U^T f = \sum_{j=1}^{n} u_j \frac{u_j^T f}{\lambda_j - \omega^2},$$

where $u_j$ is the $j$th column of $U$. The vector $x$ is a rational function with the eigenvalues of (1.2) as poles. As a function of $\omega^2$, $x$ has a vertical asymptote for each $\lambda_j \in \Omega^2$.

The idea in [5, 19] is to first compute the eigenvalues in $\Omega^2$ and then compute the solution vector $x$ as the sum

$$x = x^{(1)} + x^{(2)},$$

where

$$x^{(1)} = \sum_{j=1}^{p} u_j \frac{u_j^T f}{\lambda_j - \omega^2} \quad \text{and} \quad x^{(2)} = \sum_{j=p+1}^{n} u_j \frac{u_j^T f}{\lambda_j - \omega^2},$$

where $\lambda_1, \lambda_2, \ldots, \lambda_p$ are the eigenvalues of (1.2) in $\Omega^2$. The first term $x^{(1)}$ is then computed straightforwardly as a sum, whereas the second term $x^{(2)}$ is computed by an iterative process.

**2.2. Deflated linear system.** We first introduce the following notation. Let the columns of $U_p = [u_1, \ldots, u_p] \in \mathbf{R}^{n \times p}$ be a selection of eigenvectors of (1.2). We denote by $\mathbb{L}_p = \{\lambda_1, \ldots, \lambda_p\}$ the set of associated eigenvalues. Define $\Lambda_p = \text{diag}([\lambda_1, \ldots, \lambda_p])$. Recall that $KU_p = MU_p\Lambda_p$ and $U_p^T MU_p = I$. Define the shift-and-invert matrix $K_\sigma^{-1}M$ with $K_\sigma = K - \sigma M$. Then

$$K_\sigma^{-1} M U_p = U_p \Theta_p \quad \text{with} \quad \Theta_p = (\Lambda_p - \sigma I)^{-1} .$$

$\Theta_p$ is a diagonal matrix with $\theta_j = (\lambda_j - \sigma)^{-1}$ on the main diagonal. We will also use the Cayley transform, $K_\sigma^{-1} A = (K - \sigma M)^{-1}(K - \omega^2 M)$, and we have

$$K_\sigma^{-1} A U_p = U_p(\Lambda_p - \sigma I)^{-1}(\Lambda_p - \omega^2 I) = U_p(I - (\omega^2 - \sigma)\Theta_p) .$$

Let us now return to the solution of (1.1). We first precondition with $K_\sigma^{-1}$:

(2.3) $$K_\sigma^{-1} A x = b \quad \text{with} \quad b = K_\sigma^{-1} f .$$

As we shall see in section 2.3, this preconditioning is required for using the parameterized Lanczos method.

Let $P = U_p U_p^T M$ be the $M$-orthogonal projector onto the subspace $\mathcal{U}_p = \text{range}(U_p)$. Correspondingly, $P_\perp = I - U_p U_p^T M$ is an $M$ orthogonal projector onto $\mathcal{U}_p^\perp$, the $M$-orthogonal complement of $\mathcal{U}_p$.

First assume that $Pb = b$; i.e., $b$ lies in the range of $U_p$. The solution of (2.3) is then

$$x^{(1)} = A^{-1} K_\sigma U_p U_p^T M b = U_p(I - (\omega^2 - \sigma)\Theta_p)^{-1} U_p^T M b$$
$$= \sum_{j=1}^{p} u_j \frac{u_j^T M b}{1 - (\omega^2 - \sigma)\theta_j} .$$

We can prove that

$$x^{(1)} = \sum_{j=1}^{p} u_j \frac{u_j^T f}{\lambda_j - \sigma} .$$

For general $b$, we use $x^{(1)}$ as an initial guess for an iterative procedure for solving (2.3). So, the problem now is to find $x^{(2)}$ so that $x = x^{(1)} + x^{(2)}$ with $x^{(2)}$ the solution of

$$K_\sigma^{-1} A(x^{(1)} + x^{(2)}) = b,$$
$$K_\sigma^{-1} A x^{(2)} = b - K_\sigma^{-1} A x^{(1)},$$

and with $K_\sigma^{-1} A x^{(1)} = Pb$, we have

(2.4) $$K_\sigma^{-1} A x^{(2)} = b - Pb = P_\perp b \ .$$

In [19], a stationary iterative solver was used for each value of $\omega$ for which $x$ is computed. Experiments showed that only a few iterations (less than 10) for each $\omega_j$ are usually sufficient for convergence, where $x^{(2)}(\omega_{j-1})$ is used as the starting vector for computing $x^{(2)}(\omega_j)$. The explanation relies on the fact that $x^{(2)}$ does not have vertical asymptotes in $\Omega^2$ and that $x^{(2)}$ is a much smoother function than $x^{(1)}$. When the number of $\omega$'s is a few hundred, the total cost is still significant. Instead of a stationary solver, we can use the parameterized Lanczos method [22] since $P_\perp b$ is independent of $\omega$. This reduces the cost even more since $x$ is computed for all $\omega$'s at once with a marginal additional cost per $\omega$. We will discuss this in sections 2.3 and 2.4.

**2.3. Lanczos method.** In this section, we present the parameterized Lanczos method for the solution of (1.1), first without deflation and then with deflation. The method starts with the spectral transformation Lanczos procedure using $M$ orthogonalization [20, 21, 12].

ALGORITHM 2.1 (Lanczos procedure).
1. Let $v_0 = 0$ and set $\beta_0 = 0$.
2. Solve $K_\sigma b = f$ for $b$.
3. Let $v_1 = b/\|b\|_M$.
4. For $j = 1, 2, \ldots, k$ do:
   4.1.  Solve $K_\sigma w_j = M v_j$ for $w_j$.
   4.2.  Compute $\widehat{w}_j = w_j - v_{j-1} \beta_{j-1}$.
   4.3.  Compute $\alpha_j = v_j^T M \widehat{w}_j$.
   4.4.  Compute $\widetilde{w}_j = \widehat{w}_j - v_j \alpha_j$.
   4.5.  Let $\beta_j = \|\widetilde{w}_j\|_M$ and $v_{j+1} = \widetilde{w}_j/\beta_j$.

The computation of $w_j$ in step 4.1 requires a linear system solve with $K_\sigma$. In frequency response function computations in structural dynamics usually a direct solver is used. Alternatively, the AMLS method leads to a diagonal $K_\sigma$.

Define the tridiagonal matrix

$$T_k = \begin{bmatrix} \alpha_1 & \beta_1 & & \\ \beta_1 & \ddots & \ddots & \\ & \ddots & \ddots & \beta_{k-1} \\ & & \beta_{k-1} & \alpha_k \end{bmatrix}$$

and $V_k = [v_1, \ldots, v_k]$. The following equations readily follow from Algorithm 2.1:

(2.5) $$K_\sigma^{-1} M V_k = V_k T_k + v_{k+1} \beta_k e_k^T,$$
$$V_{k+1}^T M V_{k+1} = I \ .$$

Equation (2.5) is called the Lanczos recurrence relation.

The parameterized Lanczos method for the solution of (1.1) was first proposed in papers on model order reduction by the Padé via Lanczos method [13, 2, 3] and later studied in the context of frequency response computation by the shifted Lanczos method [14, 29, 30, 22]. This method solves the preconditioned system (2.3). This preconditioner is needed to be able to use the same Krylov space for all values of $\omega$,

as we now explain. From (2.5), it follows that, for all $\omega$,

$$K_\sigma^{-1}(K - \omega^2 M)V_k = V_k(I - (\omega^2 - \sigma)T_k) - (\omega^2 - \sigma)v_{k+1}\beta_k e_k^T$$

with $V_k$ and $T_k$ satisfying (2.5). This is the Lanczos recurrence relation for $K_\sigma^{-1}A$ for solving (2.3) by the Lanczos method (or CG). An approximate solution $\widetilde{x}$ of (1.1) is given by

(2.6) $$\widetilde{x} = V_k z,$$

where $z$ is the solution of the linear system

$$(I - (\omega^2 - \sigma)T_k)z = e_1 \|b\|_M \ .$$

This requires the solution of a $k \times k$ tridiagonal linear system. If $k$ is small, its cost is low. The preconditioned residual is

$$r = K_\sigma^{-1}(f - (K - \omega^2 M)\widetilde{x}) = (\omega^2 - \sigma)v_{k+1}\beta_k e_k^T z \ .$$

We assume that the Lanczos method is able to compute $\widetilde{x}$ with a small residual norm for $\omega \in \Omega$. If not, a larger value of $k$ could be used. An alternative is to split $\Omega$ in smaller intervals and to perform a separate Lanczos run with a new $\sigma$ for each interval as for the eigenvalue problem [17]. This is outside the scope of this paper. Therefore, we assume that $\|r(\omega)\|_M$ is "small" for all $\omega \in \Omega$.

Note that the vectors $V_k$ need not be stored if $\widetilde{x}$ is updated in each Lanczos iteration. Often, only a few elements of $x$ are wanted, and so we only need to store the desired elements of $\widetilde{x}$ for all wanted $\omega$'s. In finite precision arithmetic, the columns of $V_k$ lose orthogonality. Reorthogonalization can be used to restore the orthogonality although, strictly speaking, this is not required for convergence. Reorthogonalization, however, may produce a smaller residual norm: See, e.g., the numerical examples in [22].

**2.4. Deflation with inexact eigenvectors.** In exact arithmetic, $\mathcal{U}_p^\perp$ is an invariant subspace of $K_\sigma^{-1}M$. With $P_\perp b \in \mathcal{U}_p^\perp$, the Lanczos method applied to $K_\sigma^{-1}M$ with starting vector $P_\perp b$ produces $V_k$ whose columns are in $\mathcal{U}_p^\perp$. Therefore, the convergence is determined only by the eigenvalues of $K_\sigma^{-1}A$ associated with $\mathcal{U}_p^\perp$. As we shall see in section 3.1, these eigenvalues are favorable for fast convergence.

In practical computations, the eigenpairs are not available to full accuracy. Alternatively, if we do not need accurate eigenvalue estimates, we may save computation time in the eigenvalue computation. Define $\widehat{U}_p = [\hat{u}_1, \ldots, \hat{u}_p]$ and $\widehat{\Theta}_p = \mathrm{diag}([\hat{\theta}_1, \ldots, \hat{\theta}_p])$, with $(\hat{\theta}_j, \hat{u}_j)$ for $j = 1, \ldots, p$ approximate eigenpairs of $K_\sigma^{-1}M$. We assume that $\hat{U}_p^T M \hat{U}_p = I$. We introduce the projectors $\widehat{P} = \widehat{U}_p \widehat{U}_p^T M$ and $\widehat{P}_\perp = I - \widehat{U}_p \widehat{U}_p^T M$.

Define the residual of the approximate eigenpairs by

(2.7) $$R_p = K_\sigma^{-1}M\widehat{U}_p - \widehat{U}_p\widehat{\Theta}_p \ .$$

**2.4.1. Deflated right-hand side Lanczos (DRHSL).** The Init-CG method [9, 28, 15, 31] is an iterative method that uses inexact eigenvectors as a preconditioner. Let us first compute

$$\widetilde{x}^{(1)} = \widehat{U}_p z^{(1)} \quad \text{with} \quad z^{(1)} = (I - (\omega^2 - \sigma)\widehat{\Theta}_p)^{-1}\widehat{U}_p^T Mb$$

as the solution of $K_\sigma^{-1}A\widetilde{x}^{(1)} = \widehat{P}b$, From (2.7), we deduce that

$$K_\sigma^{-1}A\widehat{U}_p = \widehat{U}_p(I - (\omega^2 - \sigma)\widehat{\Theta}_p) - (\omega^2 - \sigma)R_p,$$
$$\widehat{P}b - K_\sigma^{-1}Ax^{(1)} = (\omega^2 - \sigma)R_pz^{(1)} .$$

Then, compute $\widetilde{x} = \widetilde{x}^{(1)} + \widetilde{x}^{(2)}$ where $\widetilde{x}^{(2)}$ is the solution of

$$K_\sigma^{-1}A\tilde{x}^{(2)} = b - K_\sigma^{-1}A\tilde{x}^{(1)} .$$

So,

$$K_\sigma^{-1}A\widetilde{x}^{(2)} = \widehat{P}_\perp b + (\widehat{P}b - K_\sigma^{-1}A\widetilde{x}^{(1)})$$

(2.8)
$$= \widehat{P}_\perp b + (\omega^2 - \sigma)R_pz^{(1)} .$$

The difficulty is that the parameterized Lanczos method cannot be used since the right-hand side in (2.8) depends on $\omega$. However, if $R_pz^{(1)}$ is small, we can omit the second term in (2.8) and work with $P_\perp b$ as for exact deflation.

So, we solve the preconditioned equation with deflated right-hand side:

$$K_\sigma^{-1}Ax^{(2)} = \widehat{P}_\perp b$$

by the parameterized Lanczos procedure. We have the recurrence relation

$$K_\sigma^{-1}MV_k = V_kT_k + \beta_kv_{k+1}e_k^T,$$

where $V_k$ is not necessarily orthogonal to $\widehat{U}_p$. An approximate solution $\widetilde{x}$ to (1.1) takes the form

$$\widetilde{x} = \widehat{U}_pz^{(1)} + V_kz^{(2)}$$

with

$$z^{(2)} = (I - (\omega^2 - \sigma)T_k)^{-1}e_1\|\widehat{P}_\perp b\|_M .$$

The residual for this solution takes the form

(2.9)        $$r = K_\sigma^{-1}(f - A\widetilde{x}) = (\omega^2 - \sigma)R_pz^{(1)} + (\omega^2 - \sigma)v_{k+1}\beta_ke_k^Tz^{(2)},$$

where the second term can be made arbitrarily small by increasing $k$. The first term is small only if $R_p$ is small enough. We will comment on this in sections 3.3 and 4.

The Proj-CG method [15] is a variation on the Init-CG method, where the Lanczos method is applied to $\widehat{P}_\perp K_\sigma^{-1}A$. We have the same difficulty as with the Init-CG method; i.e., the right-hand side depends on $\omega$.

**2.4.2. Deflated Matrix Lanczos (DML).** A practical problem with DRHSL is that in finite precision arithmetic, the components in the deflated eigenvectors can grow in the Lanczos method due to rounding errors. It is therefore usually wise to explicitly orthogonalize. The Lanczos method is now applied to the deflated matrix

$$\widehat{P}_\perp K_\sigma^{-1}M\widehat{P}_\perp.$$

We then have the recurrence relation

$$(\widehat{P}_\perp K_\sigma^{-1}M\widehat{P}_\perp)V_k = V_kT_k + \beta_kv_{k+1}e_k^T,$$

where $\widehat{P}_\perp V_k = V_k$. We thus obtain

$$(I - \widehat{U}_p \widehat{U}_p^T M) K_\sigma^{-1} M V_k = V_k T_k + \beta_k v_{k+1} e_k^T,$$

which we can rewrite as

$$(2.10) \qquad K_\sigma^{-1} M V_k = V_k T_k + \widehat{U}_p C + \beta_k v_{k+1} e_k^T$$

with $C = \widehat{U}_p^T M K_\sigma^{-1} M V_k$. From (2.7), we have that $\widehat{U}_p^T M K_\sigma^{-1} = R_p^T + \widehat{\Theta}_p \widehat{U}_p^T$. Together with $\widehat{U}_p^T M V_k = 0$, we derive that $C = R_p^T M V_k$.

We compute the solution by projecting (2.3) on the space spanned by the Krylov vectors and the Ritz vectors [18], i.e., compute $z$ from

$$(2.11) \qquad \left[\widehat{U}_p \; V_k\right]^T M K_\sigma^{-1} A \left[\widehat{U}_p \; V_k\right] z = \begin{bmatrix} \widehat{U}_p^T M b \\ V_k^T M b \end{bmatrix}.$$

If we use the spectral transformation Lanczos method as eigenvalue solver, with shift $\sigma$, we have that $\widehat{U}_p^T M \widehat{U}_p = I$, $\widehat{\Theta}_p = \widehat{U}_p^T M K_\sigma^{-1} M \widehat{U}_p$, and $\widehat{U}_p^T M R_p = 0$. As a result, (2.11) can be written as

$$(2.12) \qquad \begin{bmatrix} I + (\sigma - \omega^2)\widehat{\Theta}_p & (\sigma - \omega^2) C \\ (\sigma - \omega^2) C^T & I + (\sigma - \omega^2) T_k \end{bmatrix} \begin{bmatrix} z^{(1)} \\ z^{(2)} \end{bmatrix} = \begin{bmatrix} \widehat{U}_p^T M b \\ V_k^T M b \end{bmatrix}.$$

It is easy to see that the residual is

$$
\begin{aligned}
r &= K_\sigma^{-1}(f - A(\widehat{U}_p z^{(1)} + V_k z^{(2)})) \\
(2.13) \qquad &= (\omega^2 - \sigma)(I - V_k V_k^T M) R_p z^{(1)} + (\omega^2 - \sigma)\beta_k v_{k+1} e_k^T z^{(2)}.
\end{aligned}
$$

There is no guarantee that it is small unless $\|R_p z^{(1)}\|$ is small. We will discuss this in section 4.

Since we expect $\|R_p\|$ to be small, the off-diagonal blocks in the matrix in (2.12) can usually be omitted. We then obtain two decoupled equations, one for $z^{(1)}$ and one for $z^{(2)}$, as for DRHSL.

It should be noted that this method can be quite expensive due to the orthogonalization of $V_k$ against $\widehat{U}_p$, especially when $p$ is large. Without discussing the details, the orthogonalization cost has been successfully reduced for eigenvalue computations by selective reorthogonalization [17] and can be applied in the context of deflated iterative methods as well [31].

**3. Convergence analysis.** We analyze the spectral properties in order to understand the convergence behavior of the iterative process. We also show a connection with Padé approximation and make a statement about the required accuracy of the deflated eigenvalues.

**3.1. Spectral convergence analysis.** The spectrum of $B = K_\sigma^{-1} A$ is

$$(3.1) \qquad \phi_j = \frac{\lambda_j - \omega^2}{\lambda_j - \sigma}, \quad j = 1, \dots, n.$$

The spectral condition number of $K_\sigma^{-1} A$ is

$$\frac{\max_j\{|\phi_j|\}}{\min_j\{|\phi_j|\}}.$$

When $\lambda_j$ is far away from both $\sigma$ and $\omega^2$, then $\phi_j$ is close to one. If $\lambda_j$ is close to $\omega^2$, $|\phi_j|$ is small. When $\lambda_j$ is close to $\sigma$, then $|\phi_j|$ is large. We want to deflate the eigenvalues that make the condition number large.

LEMMA 3.1. *The matrix $B$ defined is self-adjoint with respect to the $M$ inner product, i.e., $x^T M B y = y^T M B x$. In addition, if $\mathbb{L}_p$ contains all eigenvalues of (1.2) between $\sigma$ and $\omega^2$, then the matrix $B$, restricted to $\mathcal{U}_p^\perp$, is positive definite.*

The convergence rate of the Lanczos method for the positive definite matrix $B$ (i.e., the conjugate gradients method) is then bounded from above by

$$(3.2) \qquad \left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}\right)^k,$$

where $\kappa$ is the condition number of $B$ restricted to $\mathcal{U}_p^\perp$; see, e.g., [16, Theorem 3.1.1.].

THEOREM 3.2. *Let $\mathbb{L}_p$ contain all eigenvalues in $\mathcal{I}_p = [\min(\mathbb{L}_p), \max(\mathbb{L}_p)]$. Let*

$$\lambda_m = \max\{\lambda : \lambda \in \mathbb{L}\backslash\mathbb{L}_p, \lambda \leq \min(\mathbb{L}_p)\},$$
$$\lambda_M = \min\{\lambda : \lambda \in \mathbb{L}\backslash\mathbb{L}_p, \lambda \geq \max(\mathbb{L}_p)\},$$

*and*

$$\gamma(\omega) = \frac{\lambda_M - \omega^2}{\lambda_M - \sigma} \bigg/ \frac{\omega^2 - \lambda_m}{\sigma - \lambda_m}.$$

*If $\sigma, \omega^2 \in \mathcal{I}_p$, then*

$$\kappa_M(B) \leq \max(\gamma(\omega), \gamma(\omega)^{-1}).$$

*If there are no eigenvalues on the left of $\min(\mathbb{L}_p)$, then*

$$\gamma = \frac{\lambda_M - \omega^2}{\lambda_M - \sigma}.$$

*A similar conclusion holds when there are no eigenvalues to the right of $\max(\mathbb{L}_p)$.*

*Proof.* The condition number of $B$ restricted to $\mathcal{U}_p^\perp$ is defined by

$$\kappa_M(B) = \frac{\max_{\det(B-\phi I)=0} |\phi|}{\min_{\det(B-\phi I)=0, \phi \neq 0} |\phi|}$$

$$= \max_{\lambda \in \mathbb{L}\backslash\mathbb{L}_p} \left|\frac{\lambda - \omega^2}{\lambda - \sigma}\right| \bigg/ \min_{\lambda \in \mathbb{L}\backslash\mathbb{L}_p} \left|\frac{\lambda - \omega^2}{\lambda - \sigma}\right|.$$

Figure 3.1 shows the situation where $\sigma$ and $\omega^2$ lie in $\mathcal{I}_p = [\min(\mathbb{L}_p), \max(\mathbb{L}_p)]$. The Figure also plots $|\phi| = |\lambda - \omega^2|/|\lambda - \sigma|$. If $\sigma < \omega^2$, $|\phi(\lambda)| > 1$ for $\lambda \leq \lambda_m$ and $|\phi(\lambda)| < 1$ for $\lambda \geq \lambda_M$. The maximum of $|\phi|$ outside $\mathcal{I}_p$ is attained at $\lambda_m$ and the minimum at $\lambda_M$. So, $\kappa = \gamma^{-1}$. If $\sigma > \omega^2$, $|\phi(\lambda)| < 1$ for $\lambda \leq \lambda_m$ and $|\phi(\lambda)| > 1$ for $\lambda \geq \lambda_M$. The maximum of $|\phi|$ outside $\mathcal{I}_p$ is attained at $\lambda_M$ and the minimum at $\lambda_m$. So, $\kappa = \gamma$. This proves the theorem. $\square$

Figure 3.1 shows the situation where $\sigma$ and $\omega^2$ lie in $\mathcal{I}_p = [\min(\mathbb{L}_p), \max(\mathbb{L}_p)]$. We see that if $\omega^2$ is somewhere in the middle of the interval $\mathcal{I}_p$, $|\lambda - \omega^2|$ and $|\lambda - \sigma|$ are both large so that their ratio is almost one, leading to a small $\kappa$. When $\omega^2$ lies close to $\min(\mathbb{L}_p)$ or $\max(\mathbb{L}_p)$, then some $|\phi_j|$ may be small.
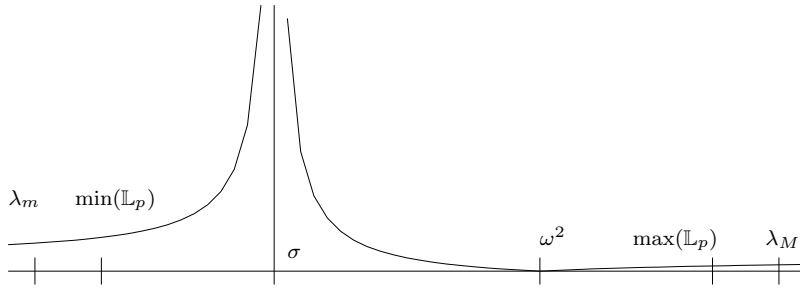
FIG. 3.1. $|\phi(\lambda)|$ where $\sigma$ and $\omega^2$ lie in $\mathbb{L}_p$.

Consider as an example, the interval $\mathcal{I}_p = [0, 100]$ and $\sigma = 80$. Let the eigenvalues be $\lambda_j = 5 + 10j$ for $j = 1, \ldots, n = 1000$. Note that $\lambda_1, \ldots, \lambda_9$ lie in $\mathcal{I}_p$. Then for $\omega^2 = 10$,

$$\phi_j = \frac{10j - 5}{10j - 75} = \frac{2j - 1}{2j - 15} \ .$$

We find that $\max_j\{|\phi_j|\} = 15$ and $\min_j\{|\phi_j|\} = 1/13$, so $\kappa = 325$. When we restrict to $\mathbf{R}\backslash\mathcal{I}_p$, we have $\max_j\{|\phi_j|\} = 19/5$ and $\min_j\{|\phi_j|\} = 9995/9925 \simeq 1$, so that $\kappa \simeq 4$.

When $\mathbb{L}_p$ includes eigenvalues outside $\Omega^2$, $\lambda_M$ and $\lambda_m$ move farther away from $\sigma$, $\omega^2_{\min}$, and $\omega^2_{\max}$, and $\gamma$ gets closer to one.

**3.2. Vector Padé connection.** As we know from [13, 2, 3, 22], $k$ steps of the Lanczos method produce a rational function that matches the first $k$ derivatives of $x^{(2)}$ in $\sigma$. In the case of exact deflation, the solution $\tilde{x} = x^{(1)} + x^{(2)}$ has two terms:

$$x^{(1)} = \sum_{j=1}^{p} u_j \frac{u_j^T f}{\lambda_j - \omega^2} \quad \text{and} \quad x^{(2)} = V_{k-p}z \ .$$

The first term is the exact solution for the right-hand side $Pb$, and the (higher order) derivatives are also exact.

The second term is computed by the Lanczos method, which means that the $k - p$ first derivatives of $x^{(2)}$ match with the exact solution for the right-hand side $P_\perp b$. As a conclusion, $\tilde{x}$ matches the first $k - p$ derivatives of the exact solution of (1.1) and interpolates the function value and the derivatives in the deflated eigenvalues $\lambda_1, \ldots, \lambda_p$ exactly.

**3.3. Residual terms.** Let

$$(3.3) \qquad\qquad\qquad T_k Y = Y \widehat{\Theta}_k$$

be the eigendecomposition of $T_k$, where $Y = [y_1, \ldots, y_k]$ and $\widehat{\Theta}_k = \mathrm{diag}(\widehat{\theta}_1, \ldots, \widehat{\theta}_k)$. The $\widehat{\theta}_j$'s are called Ritz values and the columns of $\widehat{U}_k = V_k Y = [\hat{u}_1, \ldots, \hat{u}_k]$ Ritz vectors of $K_\sigma^{-1}M$.

In practical algorithms, $k$ is either fixed beforehand or determined dynamically within the iterative process so that the solution is accurate. See [30, 22] for more details. The residuals in (2.9) and (2.13) have two terms. The second term is small when $k$ is selected high enough. The first term in (2.9) cannot be controlled by $k$. In

(2.9), we have the term

$$(\omega^2 - \sigma)R_p e_j \frac{\widehat{u}_j^T M b}{1 - (\omega^2 - \sigma)\widehat{\theta}_j} \simeq R_p e_j \frac{\sigma - \omega^2}{\widehat{\lambda}_j - \omega^2} \widehat{u}_j^T f,$$

which can be large if $\widehat{\lambda}_j \in \Omega^2$, but which is small otherwise. If we want $\|R_p z^{(1)}\|$ smaller than some tolerance, we have to require that $\|R_p e_j\|$ is small when $\widehat{\lambda}_j \in \Omega^2$. We do not have to put a strong condition on $\|R_p e_j\|$ for $\widehat{\lambda}_j \notin \Omega^2$. This also holds to some extent for (2.13).

**4. Lanczos Ritz values.** Ritz pairs are often computed by the Lanczos method. In this section, we make a connection between the solution of (1.1) by the parameterized Lanczos method and the solution of (1.2).

Recall (3.3). From (2.5), we find that

$$\begin{aligned}
\rho_j &= \|K_\sigma^{-1} M \widehat{u}_j - \widehat{\theta}_j \widehat{u}_j\|_M \\
&= \|v_{k+1} \beta_k e_k^T y_j\|_M \\
&= \beta_k |e_k^T y_j| \ .
\end{aligned}$$

Then define $\widehat{\lambda}_j = \sigma + \widehat{\theta}_j^{-1}$. If $\rho_j$ is small, we have that

$$K \widehat{u}_j \simeq \widehat{\lambda}_j M \widehat{u}_j \ .$$

The Ritz vectors form a basis of the Krylov space. Following (2.6), the solution of (1.1) can thus be expressed in terms of Ritz vectors as follows :

$$(4.1) \qquad x = \sum_{j=1}^k \hat{u}_j \frac{\hat{u}_j^T M b}{1 - (\omega^2 - \sigma)\hat{\theta}_j} = \sum_{j=1}^k \hat{u}_j \frac{\hat{w}_j^T f}{\hat{\lambda}_j - \omega^2},$$

with

$$\hat{w}_j = \hat{\theta}_j^{-1} K_\sigma^{-1} M \hat{u}_j = \hat{u}_j + v_{k+1} \frac{\beta_k e_k^T y_j}{\hat{\theta}_j}$$

the purified Ritz vector [25]. The right-hand side in (4.1) follows from $\hat{u}^T M b = (\hat{u}^T K_\sigma^{-1} M) f$.

In this section, we analyze how close these Ritz values are to the eigenvalues of $K_\sigma^{-1} A$ when the Lanczos method is applied to solve (1.1).

Usually, the stopping criterion for solving (1.1) takes the form

$$(4.2) \qquad \|r\|_M \le \tau(\|b\|_M + \|K_\sigma^{-1} A\| \|x\|_M),$$

where $\tau$ is a prescribed tolerance. The following theorem shows that the residual norms of the Ritz pairs corresponding to $\hat{\lambda}_j$ in $[\omega_{\min}^2, \omega_{\max}^2]$ are proportional to the residual tolerance for the linear system.

THEOREM 4.1. *Let* $(\hat{\theta}_j, \hat{u}_j)$, $j = 1, \ldots, k$ *be the Ritz pairs from the Lanczos method. If (4.2) holds for all* $\omega^2 \in \Omega^2$, *then*

$$\|K_\sigma^{-1} M \hat{u}_j - \hat{\theta}_j \hat{u}_j\|_M \le 6\tau |\hat{\theta}_j| \|K_\sigma^{-1} A\|_M$$

*when* $\hat{\lambda}_j = \sigma + \hat{\theta}_j^{-1} \in \Omega^2$.

*Proof.* The proof is similar to the proof of Lemma 4.1 in [22]. Let $\alpha = \omega^2 - \sigma$. From

$$r = b - K_\sigma^{-1} A \widetilde{x},$$
$$K_\sigma^{-1} A = I - \alpha K_\sigma^{-1} M,$$

$b = v_1 \|b\|_M$, and (2.6), we have

$$r = \|b\|_M v_1 - V_k (I - \alpha T_k) z + \alpha \beta_k v_{k+1} e_k^T z$$
$$= \alpha \beta_k v_{k+1} e_k^T z .$$

Next, from (3.3) and (4.1), we have that

$$z = \sum_{j=1}^k y_j \frac{y_j^T e_1 \|b\|_M}{1 - \alpha \hat{\theta}_j} = \sum_{j=1}^k y_j \frac{\hat{u}_j^T M b}{1 - \alpha \hat{\theta}_j} = \sum_{j=1}^k y_j \frac{\hat{\lambda}_j - \sigma}{\hat{\lambda}_j - \omega^2} (\hat{u}_j^T M b) .$$

With $\rho_j = \beta_k e_k^T y_j$, we have

$$r = \sum_{j=1}^k \alpha \beta_k v_{k+1} e_k^T y_j \frac{\hat{\lambda}_j - \sigma}{\hat{\lambda}_j - \omega^2} (\hat{u}_j^T M b)$$

$$= \alpha v_{k+1} \sum_{j=1}^k \rho_j \frac{\hat{\lambda}_j - \sigma}{\hat{\lambda}_j - \omega^2} (\hat{u}_j^T M b) .$$

For each $i = 1, \ldots, k$, for which $\hat{\lambda}_i \in \mathcal{I}_p$, we can determine $\omega^2 \in \mathcal{I}_p$ so that the following four statements hold:

$$(4.3) \qquad \left| \rho_i \frac{\hat{\lambda}_i - \sigma}{\hat{\lambda}_i - \omega^2} (\hat{u}_i^T M b) \right| \geq 2 \left| \sum_{j \neq i} \rho_j \frac{\hat{\lambda}_j - \sigma}{\hat{\lambda}_j - \omega^2} (\hat{u}_j^T M b) \right|,$$

$$(4.4) \qquad \left| \frac{\hat{\lambda}_i - \sigma}{\hat{\lambda}_i - \omega^2} (\hat{u}_i^T M b) \right|^2 \geq \sum_{j \neq i} \left| \frac{\hat{\lambda}_j - \sigma}{\hat{\lambda}_j - \omega^2} (\hat{u}_j^T M b) \right|^2,$$

$$(4.5) \qquad \left| \frac{\hat{\lambda}_i - \omega^2}{(\hat{\lambda}_i - \sigma)(\hat{u}_i^T M b)} \right| \|b\|_M \leq \|K_\sigma^{-1} A\|,$$

$$(4.6) \qquad |\omega^2 - \sigma| \geq |\hat{\lambda}_i - \sigma|,$$

since $|\hat{\lambda}_i - \omega^2|$ can be made arbitrarily small, by picking $\omega^2$ close to $\hat{\lambda}_i$. Note that in the Lanczos process all Ritz values are simple, so the terms in the summation for $x$ and $r$ with $j \neq i$ remain small.

We then have from (4.3) that

$$\|r\|_M / |\alpha| \geq \left| \rho_i \frac{\hat{\lambda}_i - \sigma}{\hat{\lambda}_i - \omega^2} (\hat{u}_i^T M b) \right| - \left| \sum_{j \neq i} \rho_j \frac{\hat{\lambda}_j - \sigma}{\hat{\lambda}_j - \omega^2} (\hat{u}_j^T M b) \right|$$

$$\geq \frac{1}{2} \left| \rho_i \frac{\hat{\lambda}_i - \sigma}{\hat{\lambda}_i - \omega^2} (\hat{u}_i^T M b) \right|,$$

$$(4.7) \qquad \left| \rho_i \frac{\hat{\lambda}_i - \sigma}{\hat{\lambda}_i - \omega^2} (\hat{u}_i^T M b) \right| \leq 2 \|r\|_M / |\alpha|$$

and from (4.4) that

$$\|x\|_M \le 2 \left| \frac{\hat{\lambda}_i - \sigma}{\hat{\lambda}_i - \omega^2} (\hat{u}_i^T M b) \right| .$$

From (4.2) and (4.7), we have that

$$\left| \rho_i \frac{\hat{\lambda}_i - \sigma}{\hat{\lambda}_i - \omega^2} (\hat{u}_i^T M b) \right| \le 2 \frac{\tau}{|\alpha|} \left( \|b\|_M + 2\|K_\sigma^{-1} A\| \left| \frac{\hat{\lambda}_i - \sigma}{\hat{\lambda}_i - \omega^2} (\hat{u}_i^T M b) \right| \right) ,$$

$$|\rho_i| \le 2 \frac{\tau}{|\alpha|} \left( \left| \frac{\hat{\lambda}_i - \omega^2}{(\hat{\lambda}_i - \sigma)(\hat{u}_i^T M b)} \right| \|b\|_M + 2\|K_\sigma^{-1} A\| \right) .$$

By applying (4.5) and (4.6), we finally have

$$|\rho_i| \le \frac{6\tau}{|\hat{\lambda}_i - \sigma|} \|K_\sigma^{-1} A\| .$$

Note that

$$\|K_\sigma^{-1} M \hat{u}_j - \hat{\theta}_j \hat{u}_j\|_M = |\rho_j| .$$

This proves the theorem. □

So, the backward error for the linear solves determines the backward error on the Ritz pairs. A reasonable precision for $x$ for all $\omega \in \mathcal{I}_p$ does provide accurate enough Ritz values near $\omega^2 \in \Omega^2$.

**5. Multiple right-hand sides.** The goal is to solve

(5.1) $$(K - \omega^2 M)[x_1, \dots, x_s] = [f_1, \dots, f_s] .$$

We can use a block Krylov method for solving all right-hand sides at once. The alternative is to solve each system independently, which can be useful for saving memory (less vectors to store), or is the only option when the right-hand sides are not available at once, or when $s$ is large. (The latter is actually the case for the application in example 6.2.) Probably the best alternative is to solve (5.1) block by block, i.e., using a block method with $n_b$ right-hand sides (RHS) at a time, where $n_b$ is not too large, e.g., 5. This is conceptually similar to solving (5.1) right-hand side, by right-hand side and therefore, we do not consider this approach.

The algorithms exploit simplifications by ignoring the second term in (2.8) and dropping the diagonal elements in (2.12). Those will be used by the numerical experiments.

We use the following algorithms.

ALGORITHM 5.1 (Multiple RHS solution using DRHSL).
1. Solve $K_\sigma^{-1}(K - (\omega^2 - \sigma)M)x_1 = K_\sigma^{-1} f_1$ using the parameterized Lanczos method.
2. Compute Ritz pairs $(\widehat{\Theta}_p, \widehat{U}_p)$ of $K_\sigma^{-1} M$.
3. For $j = 2, 3, \dots, s$:
   3.1. Solve $K_\sigma b_j = f_j$.
   3.2. Solve the diagonal system $(I - (\omega^2 - \sigma)\widehat{\Theta}_p)z_j = \widehat{U}_p^T M b_j$.
   3.3. Solve $K_\sigma^{-1}(K - (\omega^2 - \sigma)M)\widetilde{x}_j^{(2)} = \widehat{P}_\perp b_j$ using the parameterized Lanczos method.
   3.4. Let the solution be $\widetilde{x}_j = U_p z_j + \widetilde{x}_j^{(2)}$.

The following algorithm is conceptually close to [18] and is quite similar to [31, 32].

ALGORITHM 5.2 (Multiple RHS solution using DML).

1. Solve $K_\sigma^{-1}(K - (\omega^2 - \sigma)M)x_1 = K_\sigma^{-1}f_1$ using the parameterized Lanczos method.
2. Compute Ritz pairs $(\widehat{\Theta}_p, \widehat{U}_p)$ of $K_\sigma^{-1}M$.
3. For $j = 2, 3, \ldots, s$ :
   3.1. Solve $K_\sigma b_j = f_j$.
   3.2. Compute the Krylov space for $\hat{P}_\perp K_\sigma^{-1} M \hat{P}_\perp$ with starting vector $\hat{P}_\perp K_\sigma^{-1} f_j$.
   3.3. Solve $[\widehat{U}_p\ V_k]^T M K_\sigma^{-1} A[\widehat{U}_p\ V_k]z_j = [\widehat{U}_p\ V_k]^T M b_j$ for $z_j$.
   3.4. Let the solution be $\widetilde{x}_j = [\widehat{U}_p\ V_k]z_j$.

Algorithm 5.1 does not require the storage of the Lanczos vectors since the solution can be updated at each iteration step. However, Lanczos vectors cannot be reorthogonalized, and this leads to a loss of precision; see [22]. Algorithm 5.2 requires the storage of the Lanczos vectors, which may be undesirable when $k$ is large.

If the Ritz pairs are computed by the Lanczos method for the first right-hand side as in Algorithms 5.1 and 5.2, the residual terms related to Ritz values in $\Omega^2$ are small following the analysis in section 4. If the $p$ Ritz values contain all eigenvalues in $\Omega^2$, all linear systems in step 2.3 are positive definite and condition numbers are most likely good.

We have shown that $k$ is usually not large, since the condition number $\kappa(B)$ is small. There is one situation where $\kappa(B)$ can be large, i.e., when $\mathbb{L}_p$ does not contain approximations to all eigenvalues in $\Omega^2$. Although this is impossible to happen in theory when $u_j^T M b \neq 0$, it may happen in practice when eigenvalues are clustered or multiple.

A problem may arise also when Ritz values outside $\Omega^2$ are deflated, since their residual norms are not bounded by Theorem 4.1. However, from [4], the error on $x$ usually increases more or less monotonically when $\omega^2$ goes away from $\sigma$. From section 4, we may conclude that the Ritz residual norms also increase more or less monotonically. Moreover, in section 3.3, we argued that the Ritz values outside $\Omega^2$ do not have to have small residual norms.

**6. Numerical examples.** We now illustrate the algorithms for a number of examples. In the first example, we compare Algorithms 5.1 and 5.2. Since the simplest of both, DRHSL, performs as well as DML, we use this method for the remaining examples. The second example is related to an industrial application. The third example is constructed to make recycling fail.

The direct solver MUMPS [24] was used for solving the linear systems with $K_\sigma$ in the Lanczos procedure and Algorithms 5.1 and 5.2. The computations are dominated by the sparse factorization of $K_\sigma$ and the construction of the Krylov space. The backward solves are the dominant cost in the Krylov methods. From experience in collaboration with industry (Free Field Technologies), the cost of the factorization often corresponds to 20 to 100 times the cost of the backtransformation, depending on the problem and the linear solver used. To illustrate this, we report timings for the largest problem (second example).

**6.1. Windscreen problem.** In this section, we show the numerical performance of Algorithms 5.1 and 5.2 for a test problem arising from a structural model of a car windscreen. This is a three-dimensional (3D) problem discretized with 7564 nodes and 5400 linear hexahedral elements (3 layers of $60 \times 30$ elements). The mesh is shown in Figure 6.1(a). The material is glass with the following properties: the Young
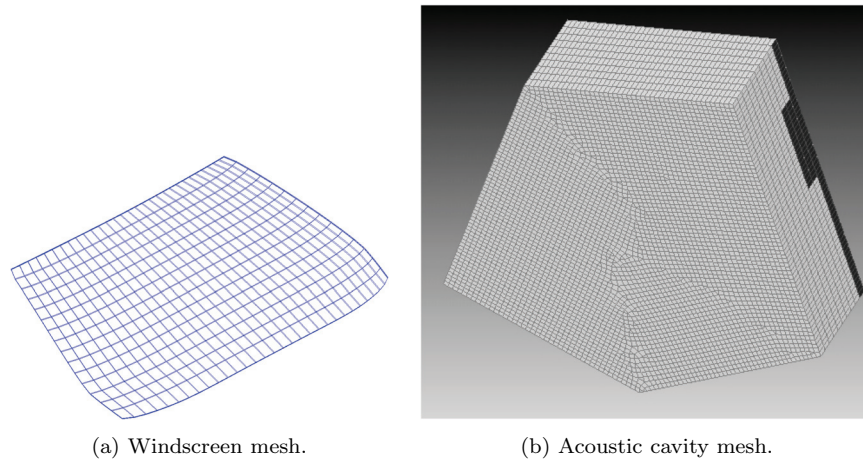
(a) Windscreen mesh.                    (b) Acoustic cavity mesh.

FIG. 6.1. *Meshes of test examples.*

modulus is $7 \times 10^{10} \text{N/m}^2$, the density is $2490 \text{kg/m}^3$, and the Poisson ratio is 0.23. The structural boundaries are free (free-free boundary conditions). The plate is subjected to a point force applied on a corner node [22, 1].

The discretized problem has dimension $n = 22,692$. The goal is to compute $x(\omega)$ with $\omega \in [0, 200]$. In order to generate the plots, the frequency range was discretized as $\{\omega_1, \ldots, \omega_m\} = \{0.5j, j = 1, \ldots, m\}$ with $m = 200$. We used shift $\sigma = 1$.

We performed a first run with right-hand side $f$ with $f_j = 0$ for $j \neq 5673$ and $f_{5673} = 1$, which corresponds to a point load on a corner of the windscreen. We used $k = 20$ Lanczos iterations. Then we kept the $p = 14$ Ritz values below $2 \times 100^2$. The residual norms of the eigenvalues in $[0, 2 \times 100^2]$ were all below $3 \times 10^{-7}$. The largest residual norm is for the Ritz value corresponding to $\omega = 104$, which is just outside the interval $\Omega$. The Ritz values in $\Omega^2 = [0, 100^2]$ have residual norms below $6 \times 10^{-9}$.

Next, we performed a second run with the right-hand side $f$ with $f_j = 0$ for $j > 1$ and $f_1 = 1$. We used 6 (additional) Lanczos iterations to make a total of $p + k = 20$ vectors. For given $\omega$,

$$\kappa = \max_{j>p} \frac{\lambda_j - \omega^2}{\lambda_j - \sigma} \Big/ \min_{j>p} \frac{\lambda_j - \omega^2}{\lambda_j - \sigma} < \widetilde{\kappa} := 1 \cdot \max_{j>p} \frac{\lambda_j - \sigma}{\lambda_j - \omega^2} = \frac{\lambda_{p+1} - \sigma}{\lambda_{p+1} - \omega^2}.$$

The largest $\widetilde{\kappa}$ is for $\omega = \max(\Omega)$. In this example, the maximum $\widetilde{\kappa}$ is $(142.089^2 - 1)/(142.089^2 - 100^2) = 1.9813$, so the convergence ratio in (3.2) is 0.1693. After six iterations, the error norm is reduced by approximately $2 \times 10^{-5}$. Figure 6.2 shows the results. Both Algorithms 5.1 and 5.2 produce the same results: Solution and error curves cannot be distinguished in the figures. The additional iterations (only six) is low so that loss of orthogonality in the Lanczos vectors is most likely not having an impact.

To see the effect of ignoring that $\|R_p\|$ is not zero, we compared with 20 Lanczos iterations without recycling. We obtained the same number of vectors as with recycling of 14 vectors and 6 additional iterations. We observed no visual difference in the error curves for $\omega$'s below 80. For $\omega$ between 80 and 100, 20 Lanczos iterations gained one digit of accuracy. We also compared with 20 instead of 6 additional iterations after recycling 14 Ritz vectors. The error curves did not show any visual difference with 6 additional iterations. In other words, performing 20 iterations without
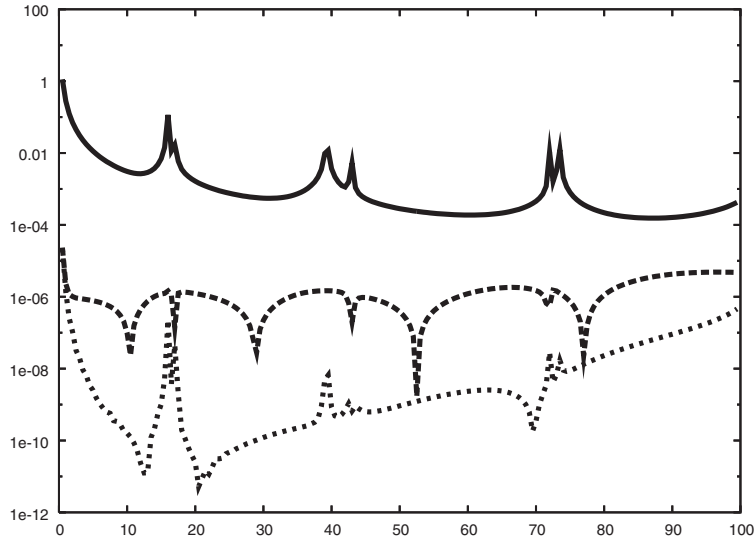
FIG. 6.2. *Windscreen problem: solution norm $\|x\|_2$ (vertical axis) in function of $\omega$ (horizontal axis) in solid line, error on $\|x\|_2$ for superposition on 14 modes as a dashed line, and error on $\|x\|_2$ for superposition on 14 modes and 6 additional Lanczos iteration as a dotted line.*

recycling is slightly more accurate for the higher frequencies (i.e., $\omega$'s farther away from the shift) than 20 iterations with recycling. The observations confirm the analysis from section 3.3.

We performed a third run with $f$ being 1 everywhere. The conclusions are similar as for the previous situation. Figure 6.3 shows the results.
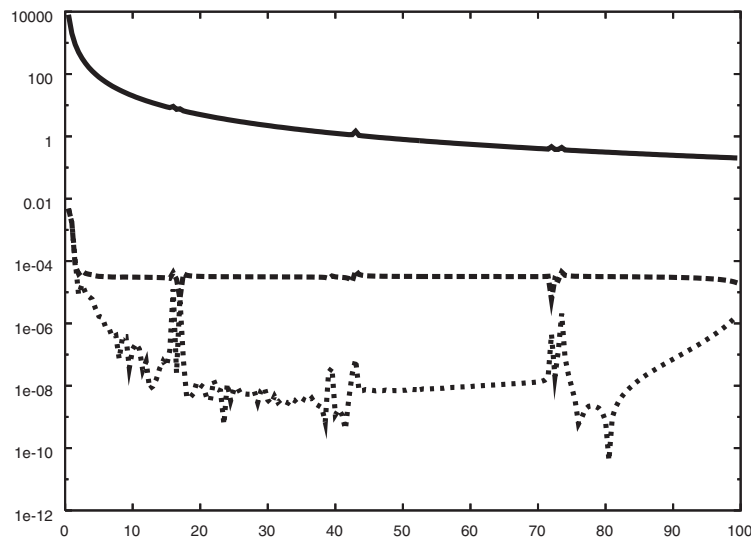


FIG. 6.3. *Windscreen problem: solution norm $\|x\|_2$ (vertical axis) in function of $\omega$ (horizontal axis) in solid line, error on $\|x\|_2$ for superposition on 14 modes as a dashed line, and error on $\|x\|_2$ for superposition on 14 modes and 6 additional Lanczos iteration as a dotted line.*

**6.2. An acoustic cavity problem.** We consider a model of an acoustic cavity discretized by 42880 finite elements and 48158 nodes, whose mesh can be seen in Figure 6.1(b). The excitation (i.e., the right-hand side) consists of 202 columns, each of them corresponding to a velocity excitation applied to one element. (See the darker elements in the mesh.) The frequency response function (FRF) is computed in the frequency range $\Omega = [0, 10000]$. Computations were performed on a Dell Precision 390 with 2GB RAM. We applied the Lanczos method with 50 vectors for the first right-hand side. The computational cost for the Krylov method was 8 seconds for the sparse matrix factorization of $K - \sigma M$ and 6 seconds for the construction of the Krylov basis.

For the second right-hand side, we kept all 31 Ritz vectors associated with the Ritz values in $[0, 2 \times 10.000^2]$. We used the DRHSL method. We then performed the following three experiments:

1. We performed 19 additional Lanczos iterations in order to obtain a basis of dimension 50. The computational cost was approximately 2 seconds, i.e., the computational cost was divided by approximately three. For 202 right-hand sides the total cost with the original Lanczos method would be of the order of 1220 seconds, whereas with recycling this would be only 416 seconds. The results are shown in Figure 6.4(a): There is no visual difference between the exact solution and the computed solution.
2. We performed $k = 50$ iterations of the Lanczos method (without recycling) as a first reference. The results are shown in Figure 6.4(b): There is no visual difference between the exact solution and the computed solution.
3. We performed $k = 19$ iterations of the Lanczos method (without recycling) as a second reference. The results are shown in Figure 6.4(c): There is a clear visual difference between the exact solution and the computed solution at a relatively low cost.

These results show that recycling indeed has a positive impact on the accuracy of the results.
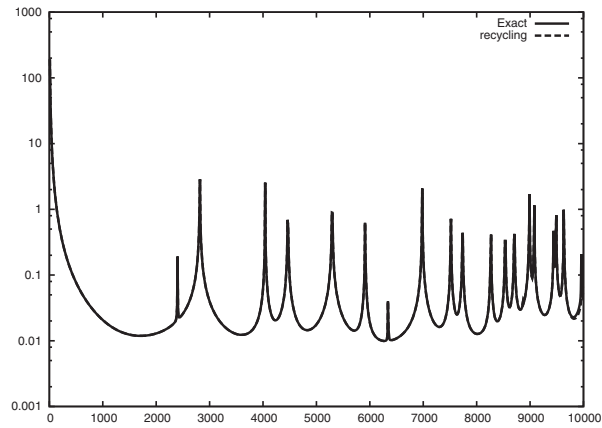
**6.3. Multiple eigenvalues.** The following example illustrates the presence of multiple eigenvalues. The matrix $K$ is the discretization of the 3D Laplacian on a unit cube and $M$ is the identity matrix. The matrix pair has multiple eigenvalues. Decompose the first right-hand side into
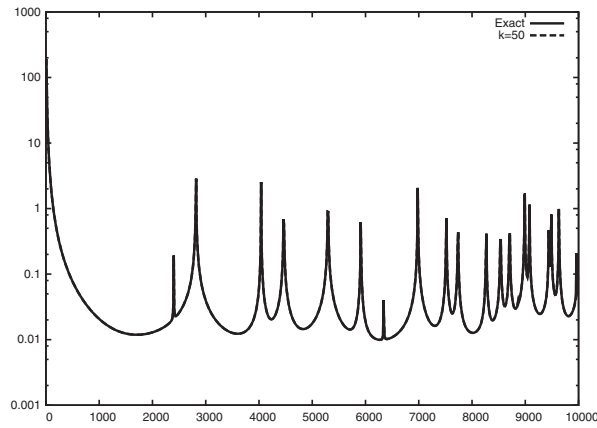
$$f_1 = \sum_{i=1}^{n} \alpha_{i,1} u_i .$$

Let $u_1$ and $u_2$ be (linearly independent) eigenvectors associated with a multiple eigenvalue, $\lambda_1$, say, so that $u_1^* M f_1 \neq 0$ and $u_2^T M f_1 = 0$. Consider a second right-hand side so that $u_2^T M f_2 \neq 0$. If a good approximation of $u_1$ is obtained by the Lanczos method for $f_1$, its recycling does not help the solution with right-hand side $f_2$, since $u_2$ is not computed. The Lanczos method cannot compute linearly independent eigenvectors of multiple eigenvalues. A similar situation arises when $f_1$ has no components in eigenvectors. In both cases, we may expect recycling not to be effective. We have chosen $f_1$ all ones and $f_2 = e_1$.

We first performed 30 steps of the Lanczos method for $f_1$, leading to the result shown in Figure 6.5(a). Note that the frequency response function does not show many peaks. Recycling the 22 Ritz pairs associated with $[0, 21^2]$, eight additional Lanczos steps produce the results from Figure 6.5(b). Figure 6.5(c) shows the results for a run with eight Lanczos iterations without recycling. The results are a little worse

(a) With recycling.
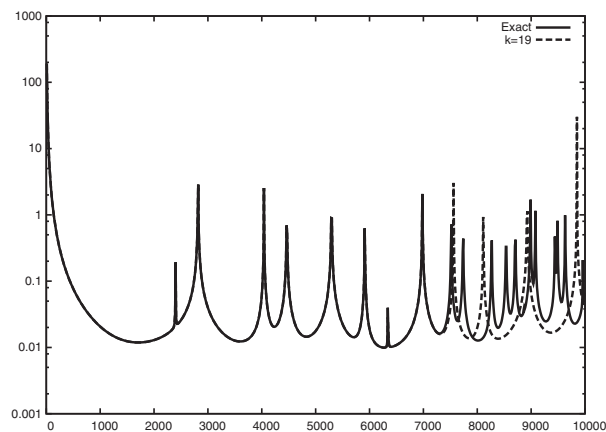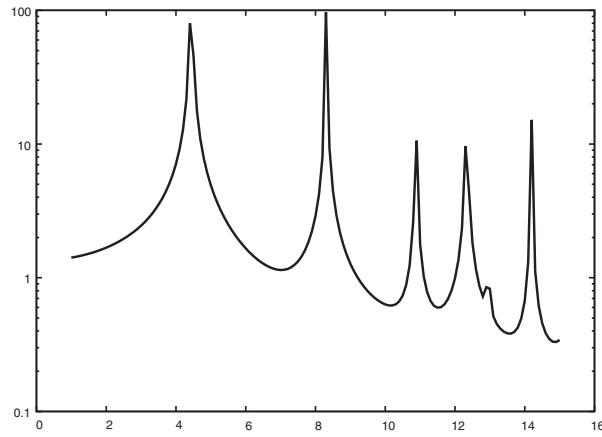
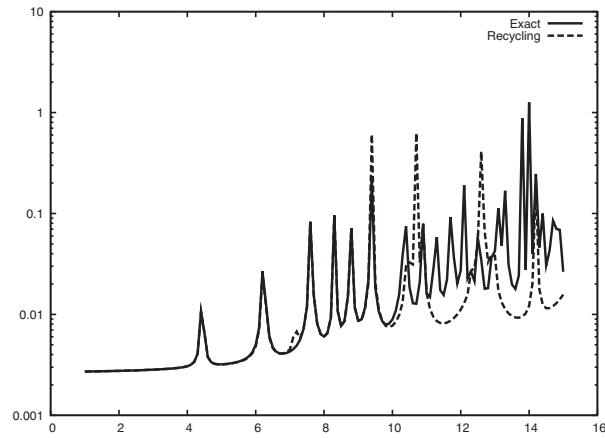

(b) Lanczos $k = 50$.



(c) Lanczos $k = 19$.



FIG. 6.4. *Exact and computed FRFs for the acoustic cavity problem. The horizontal axis corresponds to $\omega$ and the vertical axis corresponds to $\|x(\omega)\|_2$.*

(a) Lanczos for $f_1$.



(b) Recycling for $f_2$.
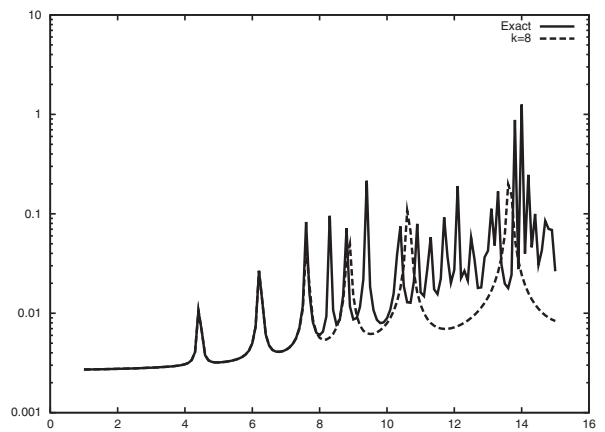


(c) Lanczos $k = 8$ for $f_2$.



FIG. 6.5. *Exact and computed results for the 3D Laplacian on a unit cube. The horizontal axis corresponds to $\omega$ and the vertical axis corresponds to $\|x(\omega)\|_2$.*

TABLE 6.1
*Ritz values computed for the 3D Laplacian. The Ritz values printed in italics are multiple Ritz.*

| Kept Ritz values from first run: | | | | |
|---|---|---|---|---|
| 20.87 | 20.0801 | 19.3112 | 18.9046 | 18.329 |
| 17.6468 | 17.0813 | 16.8657 | 15.8357 | 15.5699 |
| 14.1922 | 12.9514 | 12.3316 | 10.8806 | 10.3639 |
| 9.39104 | 8.78355 | 8.30884 | 8.30884 | *7.61551* |
| *6.23227* | 4.43694 | | | |
| Additional Ritz values from second run: | | | | |
| 24.914 | 19.2669 | 14.7007 | 11.7212 | 9.48057 |
| *7.61745* | 6.92621 | *6.23791* | | |

than with recycling. We notice that recycling does not help much for this example. The reason can be seen from Table 6.1: The second run computes a relatively large number of Ritz values in $\Omega^2$ that were not computed with the first run. These newly computed eigenvalues slow down the convergence.

**7. Conclusions.** We have used recycling Ritz vectors for solving linear systems with a parameter in frequency sweeping. We have presented an algorithm and theory for the symmetric case. We have given various arguments and a numerical example that show that recycling may significantly reduce the number of iterations in the Lanczos method.

Further extensions of this work lie in the application to proportional damping as in [23], and in using the Lanczos method as iterative method in the AMLS frequency sweeping method [19].

A comparison with the block Lanczos method would be interesting. We think that fundamental work on the reduction of the block size of the block Lanczos method should be carried out first in this context. See [27] for related work. The danger in block methods is larger memory consumption, but there may be a gain in computation time, especially when direct linear solvers are used to solve the linear systems with $K_\sigma$.

REFERENCES

[1] *Oberwolfach Model Reduction Benchmark Collection*, http://www.imtek.uni-freiburg.de/ simulation/benchmark/, 2004.
[2] Z. BAI AND R. W. FREUND, *A symmetric band Lanczos process based on coupled recurrences and some applications*, Numerical Analysis manuscript 00-8-04, Bell Laboratories, Murray Hill, NJ, 2000.
[3] Z. BAI AND R. W. FREUND, *A partial Padé-via-Lanczos method for reduced-order modeling*, Linear Algebra Appl., 332–334 (2001), pp. 139–164.
[4] Z. BAI AND Q. YE, *Error estimation of the Padé approximation of transfer functions via the Lanczos process*, Electron. Trans. Nnumer. Anal., 7 (1998), pp. 1–17.
[5] J. K. BENNIGHOF AND M. F. KAPLAN, *Frequency sweep analysis using multilevel substructuring, global modes and iteration*, in Proceedings of the 39th AIAA/ASME/ASCE/AHS Structures, Structural Dynamics and Materials Conference, 1998.
[6] J. K. BENNIGHOF AND C. K. KIM, *An adaptive multilevel substructuring method for efficient modeling of complex structures*, in Proceedings of the AIAA 33rd SDM Conference, Dallas, TX, 1992, pp. 1631–1639.
[7] J. K. BENNIGHOF AND R. B. LEHOUCQ, *An automated multilevel substructuring method for eigenspace computation in linear elastodynamics*, SIAM J. Sci. Comput., 25 (2004), pp. 2084–2106.

[8] T. F. Chan and M. K. Ng, *Galerkin projection methods for solving multiple linear systems*, SIAM J. Sci. Comput., 21 (1999), pp. 836–850.

[9] T. F. Chan and W. L. Wan, *Analysis of projection methods for solving linear systems with multiple right-hand sides*, SIAM J. Sci. Comput., 18 (1997), pp. 1698–1721.

[10] R. K. Clough and J. Penzien, *Dynamics of Structures*, McGraw-Hill, New York, 1975.

[11] D. Darnell, R. B. Morgan, and W. Wilcox, *Deflated GMRES for systems with multiple shifts and multiple right-hand sides*, Linear Algebra Appl., 429 (2008), pp. 2415–2434.

[12] T. Ericsson and A. Ruhe, *The spectral transformation Lanczos method for the numerical solution of large sparse generalized symmetric eigenvalue problems*, Math. Comp., 35 (1980), pp. 1251–1268.

[13] P. Feldman and R. W. Freund, *Efficient linear circuit analysis by Padé approximation via the Lanczos process*, IEEE Trans. Computer-Aided Design, CAD-14 (1995), pp. 639–649.

[14] A. Feriani, F. Perotti, and V. Simoncini, *Iterative system solvers for the frequency analysis of linear mechanical systems*, Technical report 1077, Instituto di Analysi Numerica–CNR, Via Addiategrasso, 20, 27100 Pavia, Italy, 1998.

[15] L. Giraud, D. Ruiz, and A. Touhami, *A comparative study of iterative solvers exploiting spectral information for SPD systems*, SIAM J. Sci. Comput., 27 (2006), pp. 1760–1786.

[16] A. Greenbaum, *Iterative Methods for Solving Linear Systems*, SIAM, Philadelphia, PA, 1997.

[17] R. G. Grimes, J. G. Lewis, and H. D. Simon, *A shifted block Lanczos algorithm for solving sparse symmetric generalized eigenproblems*, SIAM J. Matrix Anal. Appl., 15 (1994), pp. 228–272.

[18] M. E. Kilmer and E. de Sturler, *Recycling subspace information for diffuse optical tomography*, SIAM J. Sci. Comput., 27 (2006), pp. 2140–2166.

[19] J.-H. Ko and Z. Bai, *High-frequency response analysis via algebraic substructuring*, Technical report CSE-2007-18, University of California Davis, CA, 2007.

[20] C. Lanczos, *An iteration method for the solution of the eigenvalue problem of linear differential and integral operators*, J. Res. Natl. Bur. Stand., 45 (1950), pp. 255–282.

[21] C. Lanczos, *Solution of systems of linear equations by minimized iterations*, J. Res. Natl. Bur. Stand., 49 (1952), pp. 33–53.

[22] K. Meerbergen, *The solution of parametrized symmetric linear systems*, SIAM J. Matrix Anal. Appl., 24 (2003), pp. 1038–1059.

[23] K. Meerbergen, *Fast frequency response computation for Rayleigh damping*, Internat. J. Numer. Methods Engrg., 73 (2008), pp. 96–106.

[24] Mumps MUlfrontal Massively Parallel Solver, 2001, `http://graal.ens-lyon.fr/MUMPS/`.

[25] B. N.-Omid, B. N. Parlett, T. Ericsson, and P. S. Jensen, *How to implement the spectral transformation*, Math. Comp., 48 (1987), pp. 663–673.

[26] M. L. Parks, E. de Sturler, G. Mackey, D. D. Johnson, and S. Maiti *Recycling Krylov subspaces for sequences of linear systems*, SIAM J. Sci. Comput., 28 (2006), pp. 1651–1674.

[27] M. Robbé and M. Sadkane, *Use of near-breakdowns in the block Arnoldi method for solving large Sylvester equations*, Appl. Numer. Math., 58 (2008), pp. 486–498.

[28] Y. Saad, M. Yeung, J. Erhel, and F. Guyomarc'h, *A deflated version of the conjugate gradient algorithm*, SIAM J. Sci. Comput., 21 (2000), pp. 1909–1926.

[29] V. Simoncini, *Linear systems with a quadratic parameter and application to structural dynamics*, in D. R. Kincaid and A. Elster, editor, Iterative Methods in Scientific Computation IV, IMACS Series in Computational and Applied Mathematics, Vol. 5, D. R. Kincaid and A. Elster, eds., Elsevier, Amsterdam, 1999, pp. 451–461.

[30] V. Simoncini and F. Perotti, *On the numerical solution of $(\lambda^2 A + \lambda B + C)$, $x = b$ and application to structural dynamics*, SIAM J. Sci. Comput., 23 (2002), pp. 1875–1897.

[31] A. Stathopoulos and K. Orginos, *Computing and deflating eigenvalues while solving multiple right hand side linear systems with an application to quantum chromodynamics*, Technical report arXiv.org 0707.0131v2, College of William and Mary, 2008. Revised version 2009.

[32] S. Wang, E. d. Sturler, and G. H. Paulino, *Large-scale topology optimization using preconditioned Krylov subspace methods with recycling*, Internat. J. Numer. Methods Engrg., 69 (2007), pp. 2441–2468.

[33] E. L. Wilson, M.-W. Yuan, and J. M. Dickens, *Dynamic analysis by direct superposition of Ritz vectors*, Earthquake Engrg. Struct. Dynamics, 10 (1982), pp. 813–821.