

## A FRAMEWORK FOR DEFLATED AND AUGMENTED KRYLOV SUBSPACE METHODS\*

ANDRÉ GAUL<sup>†</sup>, MARTIN H. GUTKNECHT<sup>‡</sup>, JÖRG LIESEN<sup>†</sup>, AND REINHARD  
NABBEN<sup>†</sup>

**Abstract.** We consider deflation and augmentation techniques for accelerating the convergence of Krylov subspace methods for the solution of nonsingular linear algebraic systems. Despite some formal similarity, the two techniques are conceptually different from preconditioning. Deflation (in the sense the term is used here) “removes” certain parts from the operator making it singular, while augmentation adds a subspace to the Krylov subspace (often the one that is generated by the singular operator); in contrast, preconditioning changes the spectrum of the operator without making it singular. Deflation and augmentation have been used in a variety of methods and settings. Typically, deflation is combined with augmentation to compensate for the singularity of the operator, but both techniques can be applied separately. We introduce a framework of Krylov subspace methods that satisfy a Galerkin condition. It includes the families of orthogonal residual and minimal residual methods. We show that in this framework augmentation can be achieved either explicitly or, equivalently, implicitly by projecting the residuals appropriately and correcting the approximate solutions in a final step. We study conditions for a breakdown of the deflated methods, and we show several possibilities to avoid such breakdowns for the deflated minimum residual (MINRES) method. Numerical experiments illustrate properties of different variants of deflated MINRES analyzed in this paper.

**Key words.** Krylov subspace methods, augmentation, deflation, subspace recycling, CG, MINRES, GMRES, RMINRES

**AMS subject classifications.** 65F10, 65F08

**DOI.** 10.1137/110820713

**1. Introduction.** There are numerous techniques to accelerate the speed of convergence of Krylov subspace methods for solving large linear algebraic systems

$$(1.1) \quad \mathbf{Ax} = \mathbf{b},$$

where  $\mathbf{A} \in \mathbb{C}^{N \times N}$  is nonsingular and  $\mathbf{b} \in \mathbb{C}^N$ . The most widely used technique is *preconditioning*. Here the system (1.1) is modified using left- or right-multiplications with a *nonsingular* matrix (called the preconditioner). A typical goal of preconditioning is to obtain a modified matrix that is in some sense close to the identity matrix. For surveys of preconditioning techniques we refer to the books by Greenbaum [26, Part II] and Saad [46, Chapters 9–14] and the survey of Benzi [3].

Here we consider two approaches for convergence acceleration that are called *deflation* and *augmentation*. Let us briefly describe the main ideas of the two techniques. In deflation the system (1.1) is multiplied (at least implicitly) with a suitably chosen projection, and the general goal is to “eliminate” components that supposedly slow down convergence. Typically these are components that correspond to small

---

\*Received by the editors January 11, 2011; accepted for publication (in revised form) February 13, 2013; published electronically May 14, 2013.

<http://www.siam.org/journals/simax/34-2/82071.html>

<sup>†</sup>Institut für Mathematik, Technische Universität Berlin, D-10623 Berlin, Germany (gaul@math.tu-berlin.de, liesen@math.tu-berlin.de, nabben@math.tu-berlin.de). These authors were supported by the DFG Forschungszentrum MATHEON. The work of the third author was also supported by the Heisenberg Program of the DFG.

<sup>‡</sup>Seminar for Applied Mathematics, ETH Zurich, CH-8092 Zurich, Switzerland (mhg@math.ethz.ch). This author was supported by the DFG Forschungszentrum MATHEON and the Mercator Visiting Professorship Program of the DFG.

eigenvalues. Multiplication by the projection turns the system (1.1) into a consistent singular one, which is then solved by a Krylov subspace method. We need to mention, however, that techniques have been proposed that move small eigenvalues of  $\mathbf{A}$  to some large common value, say, to the value 1; see [1, 17, 31]. Some authors refer to these techniques as “deflation” too. In augmentation techniques the search space of the Krylov subspace method, which is at the same time the Galerkin test space, is “enlarged” by a suitably chosen subspace. A typical goal is to add information about the problem to the search space that is slowly revealed in the Krylov subspace itself, e.g., eigenvectors corresponding to small eigenvalues.

Deflation and augmentation techniques can be combined with conventional preconditioning techniques. Then the projection and augmentation parameters have to be adapted to the preconditioned matrix. In this paper, we assume that (1.1) is already in preconditioned form, i.e.,  $\mathbf{A}$  is the preconditioned matrix and  $\mathbf{b}$  the preconditioned right-hand side. Details of preconditioning techniques will thus not be addressed here.

We will now give a brief overview of existing deflation and augmentation strategies. For a more comprehensive presentation we refer to section 9 of the survey article by Simoncini and Szyld [49]. The first deflation and augmentation techniques in the context of Krylov subspace methods appeared in the papers of Nicolaides [41] and Dostál [12]. Both proposed deflated variants of the CG method [29] to accelerate the speed of convergence for symmetric positive definite (spd) matrices  $\mathbf{A}$  arising from discretized elliptic partial differential equations. Since these early works deflation and augmentation have become widely used tools. Several authors working in different fields of numerical analysis applied them to many Krylov subspace methods, and they use a variety of techniques to determine a deflation subspace. A review of all applications is well beyond this introduction. We concentrate in the following on some—but not all—key contributions.

For nonsymmetric systems Morgan [36] and also Chapman and Saad [6] extracted approximate eigenvectors of  $\mathbf{A}$  from the Krylov subspace generated by the GMRES method [47], and then they augmented the Krylov subspace with these vectors; for related references we refer to [22]. A comparable approach in the context of the CG method for spd matrices  $\mathbf{A}$  was described by Saad, Yeung, Erhel, and Guyomarc’h [48]. De Sturler [10] introduced the GCRO method, which involves an outer GCR iteration [15, 16] and an inner deflated GMRES method where the space used for deflation depends on the outer iteration. This method has been extended to GCROT in [11] to incorporate truncation strategies when restarts are necessary. In [33] Kolotilina used a twofold deflation technique for simultaneously deflating the  $r$  largest and the  $r$  smallest eigenvalues by an appropriate deflating subspace of dimension  $r$ . An analysis of acceleration strategies (including augmentation) for minimal residual methods was given by Saad [45] and for restarted methods by Eiermann, Ernst, and Schneider [14]. The latter work analyzes minimal residual (MR) and orthogonal residual (OR) methods in a general framework that allows approximations from arbitrary correction spaces. By using multiple correction spaces forming a direct sum, several cases of augmentation and deflation are discussed. The analysis concentrates on (nearly)  $\mathbf{A}$ -invariant augmentation spaces.

In [37] Morgan proposed a block-GMRES method for multiple right-hand sides that deflates approximated eigenvectors when GMRES is restarted. A similar method for solving systems with multiple shifts and multiple right-hand sides has been introduced by Darnell, Morgan, and Wilcox [9]. Giraud et al. [25] recently developed a flexible GMRES variant with deflated restarting where the preconditioner may vary

from one iteration to the next. In [42] Olshanskii and Simoncini studied spectral properties of saddle point matrices preconditioned with a block-diagonal preconditioner and applied a deflated minimum residual (MINRES) method to the resulting symmetric and indefinite matrix in order to alleviate the influence of a few small outlying eigenvalues. Theoretical results for deflated GMRES based on an exactly  $\mathbf{A}$ -invariant subspace have been presented in [61].

In addition to deflation/augmentation spaces based on approximative eigenvectors, other choices have been studied. Mansfield [34] showed how Schur complement-type domain decomposition methods can be seen as a series of deflations. Nicolaidis [41] constructed a deflation technique based on piecewise constant interpolation from a set of  $r$  subdomains, and he pointed out that deflation might be effectively used with a conventional preconditioner. In [35] Mansfield used the same “subdomain deflation” in combination with damped Jacobi smoothing and obtained a preconditioner that is related to the two-grid method. Baker, Jessup, and Manteuffel [2] proposed a GMRES method that is augmented upon restarts by approximations to the error.

In [38, 39, 40] Nabben and Vuik described similarities between the deflation approach and domain decomposition methods for arbitrary deflation spaces. This comparison was extended to multigrid methods in [54, 53].

This brief survey indicates that in principle deflation or augmentation can be incorporated into every Krylov subspace method. However, some methods may suffer from mathematical shortcomings like breakdowns or numerical problems due to round-off errors. The main goal of this paper is not to add further examples to the existing collection but to introduce first a suitable framework for a whole family of such augmented and deflated methods (section 2) and then to prove some results just assuming this framework (section 3). The framework focuses on Krylov subspace methods whose residuals satisfy a certain Galerkin condition with respect to a true or formal inner product. In section 3, we mathematically characterize the equivalence of two approaches for realizing such methods and discuss them along with potential pitfalls. We then discuss known approaches to deflate CG (section 4), GMRES (section 5), and MINRES (section 6) in the light of our general equivalence theorem. Among other results, this will show that a recent version of deflated MINRES, which is part of the “recycling” MINRES (RMINRES) method suggested by Wang, de Sturler, and Paulino [57], can break down and will show how these breakdowns can be avoided by either adapting the right-hand side or the initial guess. We do not focus on specific implementations or algorithmic details but on the mathematical theory of these methods. For the numerical application in section 6.3 we draw on the most robust MINRES implementation that is available.

**2. A framework for deflated and augmented Krylov methods.** In this section we describe a general framework for deflation and augmentation, which simultaneously covers several Krylov subspace methods whose residuals satisfy a Galerkin condition. Given an initial guess  $\mathbf{x}_0 \in \mathbb{C}^N$ , a positive integer  $n$ , an  $n$ -dimensional subspace  $\mathcal{S}_n$  of  $\mathbb{C}^N$ , and a nonsingular matrix  $\mathbf{B} \in \mathbb{C}^{N \times N}$ , let us first consider an approximation  $\mathbf{x}_n$  to the solution  $\mathbf{x}$  of the form

$$(2.1) \quad \mathbf{x}_n \in \mathbf{x}_0 + \mathcal{S}_n,$$

so that the corresponding residual

$$\mathbf{r}_n := \mathbf{b} - \mathbf{A}\mathbf{x}_n \in \mathbf{r}_0 + \mathbf{A}\mathcal{S}_n$$

satisfies

$$(2.2) \quad \mathbf{r}_n \perp \mathbf{B}\mathcal{S}_n.$$

If  $\mathbf{B}^H\mathbf{A}$  is Hermitian positive definite (Hpd), then  $\mathbf{B}^H\mathbf{A}$  induces an inner product  $\langle \cdot, \cdot \rangle_{\mathbf{B}^H\mathbf{A}}$ , a corresponding norm  $\|\cdot\|_{\mathbf{B}^H\mathbf{A}}$ , and an orthogonality  $\perp_{\mathbf{B}^H\mathbf{A}}$ . Imposing (2.1) and (2.2) can then be seen to be equivalent to solving the following minimization problem:

$$(2.3) \quad \text{find } \mathbf{x}_n \in \mathbf{x}_0 + \mathcal{S}_n \quad \text{s.t.} \quad \|\mathbf{x} - \mathbf{x}_n\|_{\mathbf{B}^H\mathbf{A}} = \min_{\mathbf{y} \in \mathbf{x}_0 + \mathcal{S}_n} \|\mathbf{x} - \mathbf{y}\|_{\mathbf{B}^H\mathbf{A}}.$$

Note that due to  $\mathbf{r}_n = \mathbf{A}(\mathbf{x} - \mathbf{x}_n)$  the condition (2.2) can be written as orthogonality condition for the error  $\mathbf{x} - \mathbf{x}_n$ :

$$(2.4) \quad (\mathbf{x} - \mathbf{x}_n) \perp_{\mathbf{B}^H\mathbf{A}} \mathcal{S}_n.$$

The following two cases where  $\mathbf{B}^H\mathbf{A}$  is Hpd are of particular interest:

- (1)  $\mathbf{B} = \mathbf{I}$  if  $\mathbf{A}$  itself is Hpd;
- (2)  $\mathbf{B} = \mathbf{A}$  for general nonsingular  $\mathbf{A}$ .

The case (1) is the one where (2.2) is a typical Galerkin condition:  $\mathbf{A}$  is Hpd and the residual  $\mathbf{r}_n$  is orthogonal to the linear search space  $\mathcal{S}_n$  for  $\mathbf{x}_n - \mathbf{x}_0$ . In (2.3) we then have

$$(2.5) \quad \|\mathbf{x} - \mathbf{x}_n\|_{\mathbf{A}} = \|\mathbf{r}_n\|_{\mathbf{A}^{-1}},$$

so while the error is minimal in the  $\mathbf{A}$ -norm, the residual is minimal in the  $\mathbf{A}^{-1}$ -norm.

In this paper we will refer to (2.2) also in case (2) as a Galerkin condition, because the search space and the test space are still essentially the same. However, in this case

$$(2.6) \quad \|\mathbf{x} - \mathbf{x}_n\|_{\mathbf{A}^H\mathbf{A}} = \|\mathbf{r}_n\|_2,$$

so (2.2) implies that the 2-norm of the residual is minimized. Consequently, in both cases a minimization property holds.

If the search space  $\mathcal{S}_n$  is the  $n$ th Krylov subspace generated by  $\mathbf{A}$  and the initial residual  $\mathbf{r}_0 := \mathbf{b} - \mathbf{A}\mathbf{x}_0$ , i.e., if

$$(2.7) \quad \mathcal{S}_n = \mathcal{K}_n(\mathbf{A}, \mathbf{r}_0) := \text{span}\{\mathbf{r}_0, \mathbf{A}\mathbf{r}_0, \dots, \mathbf{A}^{n-1}\mathbf{r}_0\},$$

then, in case (1), conditions (2.1)–(2.2) mathematically characterize the CG method [29]. It is the prototype of an OR method characterized by (2.1) and (2.2) with  $\mathbf{B} = \mathbf{I}$ .

In case (2), conditions (2.1)–(2.2) with  $\mathcal{S}_n = \mathcal{K}_n(\mathbf{A}, \mathbf{r}_0)$  mathematically characterize the GCR [15] and GMRES [47] methods and, for Hermitian  $\mathbf{A}$ , the MINRES [43] method. If  $\mathbf{A}$  is even Hpd, we can resort to Stiefel's conjugate residual (CR) method [52]. All these are prototype MR methods characterized by (2.1) and (2.2) with  $\mathbf{B} = \mathbf{A}$ .

OR and MR methods often come in pairs defined by the properties of  $\mathbf{A}$ , the Krylov search space, and, to some extent, the fundamental structure of the algorithms. Examples of such pairs are CG/CR, GCG/GCR, full orthogonalization method (FOM)/GMRES, and CGNE/CGNR. It has been pointed out many times, see, e.g., [4, 8, 13, 14, 27], that the residuals of these pairs of OR/OM methods and in particular the residual norms are related in a simple fashion.

A fact related to the OR/MR residual connection is that the iterates and residuals of an MR method can be found from those of the corresponding OR method by a smoothing process introduced by Schönauer; see [58, 60, 27, 28]. The reverse process also exists [27]. Again these processes hold for the residuals of the deflated system and, since they are identical, for those of the explicit augmentation approach.

If  $\mathbf{B}^H\mathbf{A}$  is not Hpd, the minimization property (2.3) no longer makes sense, but we may still request that the orthogonality condition (2.2), or, equivalently, (2.4), hold. Resulting algorithms may then break down since an approximate solution  $\mathbf{x}_n$  satisfying the conditions may not exist for some  $n$ . Nevertheless, such methods are occasionally applied in practice. In particular, the choice

$$(3) \quad \mathbf{B} = \mathbf{I} \text{ and } \mathbf{A} \text{ nonsingular}$$

covers the FOM of Saad [44, 46], which is sometimes also referred to as Arnoldi method for linear algebraic systems.

For minimizing the error  $\mathbf{x}_n - \mathbf{x}$  in the 2-norm one has to choose

$$(4) \quad \mathbf{B} = \mathbf{A}^{-H} \text{ and } \mathbf{A} \text{ nonsingular.}$$

Since multiplication by  $\mathbf{A}^{-H}$  is not feasible, these methods only work for particular search spaces; the simplest choice is

$$(2.8) \quad \mathcal{S}_n = \mathbf{A}^H \mathcal{K}_n(\mathbf{A}^H, \mathbf{r}_0).$$

Unlike the normal Krylov search space of (2.7), this one has the drawback that the (exact) solution of the system need not be in one of these spaces, i.e., even in exact arithmetic convergence is not guaranteed. One interesting example based on this choice is the generalized minimum error (GMERR) method of Weiss [59]. Earlier, for spd matrices, such a method was proposed by Fridman [24], and an alternative algorithm was mentioned by Fletcher [21]. Symmetric indefinite systems can be treated in this way with the SYMMLQ algorithm of Paige and Saunders [43]; see also Freund [23] for a review of methods featuring this optimality criterion and yet another algorithm called ME to achieve it.

Finally, we can easily incorporate the CGNR method [29] for solving overdetermined linear systems in the setting of (2.1) and (2.2) by choosing the appropriate Krylov search space. Given such a system

$$(2.9) \quad \mathbf{E}\mathbf{x} = \mathbf{f}$$

with a full-rank  $M \times N$ -matrix  $\mathbf{E}$  (where  $M \geq N$ ), the corresponding normal equations are  $\mathbf{E}^H\mathbf{E}\mathbf{x} = \mathbf{E}^H\mathbf{f}$ , i.e.,  $\mathbf{A}\mathbf{x} = \mathbf{b}$  with  $\mathbf{A} := \mathbf{E}^H\mathbf{E}$  and  $\mathbf{b} := \mathbf{E}^H\mathbf{f}$ . Since  $\mathbf{A}$  is Hpd, we can apply the CG method which corresponds to case (1) and

$$(2.10) \quad \mathcal{S}_n = \mathcal{K}_n(\mathbf{E}^H\mathbf{E}, \mathbf{E}^H\mathbf{s}_0)$$

with  $\mathbf{s}_0 := \mathbf{f} - \mathbf{E}\mathbf{x}_0$ . In this situation we have to distinguish between the residuals  $\mathbf{r}_n := \mathbf{b} - \mathbf{A}\mathbf{x}_n = \mathbf{E}^H\mathbf{f} - \mathbf{E}^H\mathbf{E}\mathbf{x}_n$  of the normal equations and the residuals  $\mathbf{s}_n := \mathbf{f} - \mathbf{E}\mathbf{x}_n$  of the given system (2.9). The CGNR method allows one to keep track of both. The latter residuals satisfy

$$(2.11) \quad \mathbf{s}_n \in \mathbf{s}_0 + \mathbf{E}\mathcal{K}_n(\mathbf{E}^H\mathbf{E}, \mathbf{E}^H\mathbf{f}), \quad \mathbf{s}_n \perp \mathbf{E}\mathcal{K}_n(\mathbf{E}^H\mathbf{E}, \mathbf{E}^H\mathbf{f}),$$

and they can be seen to minimize the 2-norm of  $\mathbf{s}_n$ . Note that it can be viewed as an MR method with a possibly nonsquare  $\mathbf{B} = \mathbf{E}$ ; see [27].

A method that also fits into our framework, though with some modifications, is the CGNE method, also called Craig’s method [7], which can also be used for solving

underdetermined linear algebraic systems (2.9) with a full-rank  $M \times N$ -matrix  $\mathbf{E}$ . The search space for  $\mathbf{x}_n$  in this case is (2.10), but the Galerkin condition becomes  $\mathbf{s}_n \perp \mathcal{K}_n(\mathbf{E}\mathbf{E}^H, \mathbf{s}_0)$ .

Since we are aiming at a *general framework*, let us for the moment consider an arbitrary, possibly singular matrix  $\widehat{\mathbf{A}} \in \mathbb{C}^{N \times N}$  and an arbitrary vector  $\widehat{\mathbf{v}} \in \mathbb{C}^N$ , such that the Krylov subspace  $\mathcal{K}_n(\widehat{\mathbf{A}}, \widehat{\mathbf{v}})$  has dimension  $n$ .

Instead of a search space of the form  $\mathcal{S}_n = \mathcal{K}_n(\mathbf{A}, \mathbf{r}_0)$  we focus from now on augmented Krylov subspaces of the form

$$(2.12) \quad \mathcal{S}_n := \mathcal{K}_n(\widehat{\mathbf{A}}, \widehat{\mathbf{v}}) + \mathcal{U}.$$

We suppose that  $\mathcal{U}$  has dimension  $k$ ,  $0 < k < N$ , and denote by  $\mathbf{U} \in \mathbb{C}^{N \times k}$  a matrix whose columns form a basis of  $\mathcal{U}$ , and by  $\mathbf{V}_n \in \mathbb{C}^{N \times n}$  one whose columns form a basis of  $\mathcal{K}_n(\widehat{\mathbf{A}}, \widehat{\mathbf{v}})$ , so that (2.1) can be written as

$$(2.13) \quad \mathbf{x}_n = \mathbf{x}_0 + \mathbf{V}_n \mathbf{y}_n + \mathbf{U} \mathbf{u}_n$$

for some vectors  $\mathbf{y}_n \in \mathbb{C}^n$  and  $\mathbf{u}_n \in \mathbb{C}^k$ . Of course,  $\mathbf{U}$  may be redefined when an algorithm like GMRES is restarted, but we will not account for that in our notation.

Assuming the general structure of the search space  $\mathcal{S}_n$  in (2.12) we will now investigate augmented Galerkin-type methods that still satisfy (2.1) and (2.2).

**3. A general equivalence theorem.** Our goal in this section is to show that augmentation can be achieved either explicitly as in (2.12), or implicitly, namely, by projecting the residuals appropriately and correcting the approximate solutions in a final step. Our main result is stated in Theorem 3.2 below.

In order to satisfy (2.2), the residual  $\mathbf{r}_n = \mathbf{b} - \mathbf{A}\mathbf{x}_0$  must be orthogonal to both  $\mathbf{B}\mathcal{K}_n(\widehat{\mathbf{A}}, \widehat{\mathbf{v}})$  and  $\mathbf{B}\mathcal{U}$ , hence it must satisfy the pair of orthogonality conditions

$$(3.1) \quad \mathbf{r}_n \perp \mathbf{B}\mathcal{K}_n(\widehat{\mathbf{A}}, \widehat{\mathbf{v}}) \quad \text{and} \quad \mathbf{r}_n \perp \mathbf{B}\mathcal{U}.$$

Let us concentrate on the second condition of (3.1), which can be written as

$$\mathbf{0} = \mathbf{U}^H \mathbf{B}^H \mathbf{r}_n = \mathbf{U}^H \mathbf{B}^H (\mathbf{r}_0 - \mathbf{A}\mathbf{V}_n \mathbf{y}_n - \mathbf{A}\mathbf{U} \mathbf{u}_n) = \mathbf{U}^H \mathbf{B}^H (\mathbf{r}_0 - \mathbf{A}\mathbf{V}_n \mathbf{y}_n) - \mathbf{E}_B \mathbf{u}_n,$$

where

$$(3.2) \quad \mathbf{E}_B := \mathbf{U}^H \mathbf{B}^H \mathbf{A} \mathbf{U} \in \mathbb{C}^{k \times k}.$$

Clearly, if  $\mathbf{B}^H \mathbf{A}$  is Hpd, then  $\mathbf{E}_B$  is Hpd too—though in a smaller space—and thus nonsingular. In the following derivation,  $\mathbf{B}^H \mathbf{A}$  need not be Hpd, but we must then assume that  $\mathbf{E}_B$  is nonsingular. Then the second orthogonality condition is equivalent to

$$(3.3) \quad \mathbf{u}_n = \mathbf{E}_B^{-1} \mathbf{U}^H \mathbf{B}^H (\mathbf{r}_0 - \mathbf{A}\mathbf{V}_n \mathbf{y}_n).$$

Substituting this into (2.13) gives

$$(3.4) \quad \begin{aligned} \mathbf{x}_n &= \mathbf{x}_0 + \mathbf{V}_n \mathbf{y}_n + \mathbf{U} (\mathbf{E}_B^{-1} \mathbf{U}^H \mathbf{B}^H (\mathbf{r}_0 - \mathbf{A}\mathbf{V}_n \mathbf{y}_n)) \\ &= (\mathbf{I} - \mathbf{U} \mathbf{E}_B^{-1} \mathbf{U}^H \mathbf{B}^H \mathbf{A}) (\mathbf{x}_0 + \mathbf{V}_n \mathbf{y}_n) + \mathbf{U} \mathbf{E}_B^{-1} \mathbf{U}^H \mathbf{B}^H \mathbf{b}, \end{aligned}$$

$$(3.5) \quad \begin{aligned} \mathbf{r}_n &= \mathbf{r}_0 - \mathbf{A}\mathbf{V}_n \mathbf{y}_n - \mathbf{A}\mathbf{U} (\mathbf{E}_B^{-1} \mathbf{U}^H \mathbf{B}^H (\mathbf{r}_0 - \mathbf{A}\mathbf{V}_n \mathbf{y}_n)) \\ &= (\mathbf{I} - \mathbf{A}\mathbf{U} \mathbf{E}_B^{-1} \mathbf{U}^H \mathbf{B}^H) (\mathbf{r}_0 - \mathbf{A}\mathbf{V}_n \mathbf{y}_n). \end{aligned}$$

To simplify the notation we define the  $(N \times N)$ -matrices

$$\begin{aligned}
 \mathbf{M}_B &:= \mathbf{U}\mathbf{E}_B^{-1}\mathbf{U}^H = \mathbf{U}(\mathbf{U}^H\mathbf{B}^H\mathbf{A}\mathbf{U})^{-1}\mathbf{U}^H, \\
 \mathbf{P}_B &:= \mathbf{I} - \mathbf{A}\mathbf{M}_B\mathbf{B}^H, \\
 \mathbf{Q}_B &:= \mathbf{I} - \mathbf{M}_B\mathbf{B}^H\mathbf{A}.
 \end{aligned}
 \tag{3.6}$$

Using these matrices (3.4) and (3.5) take the form

$$\mathbf{x}_n = \mathbf{Q}_B(\mathbf{x}_0 + \mathbf{V}_n\mathbf{y}_n) + \mathbf{M}_B\mathbf{B}^H\mathbf{b},
 \tag{3.7}$$

$$\mathbf{r}_n = \mathbf{P}_B(\mathbf{r}_0 - \mathbf{A}\mathbf{V}_n\mathbf{y}_n).
 \tag{3.8}$$

Note that imposing the second orthogonality condition in (3.1) on the residual  $\mathbf{r}_n$  has determined the vector  $\mathbf{u}_n$ , which has therefore “disappeared” in (3.7)–(3.8). We next state some basic properties of the matrices  $\mathbf{P}_B$  and  $\mathbf{Q}_B$ . The proof of these properties is straightforward and is therefore omitted.

LEMMA 3.1. *Let  $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{N \times N}$  and  $\mathbf{U} \in \mathbb{C}^{N \times k}$  be such that  $\mathbf{E}_B := \mathbf{U}^H\mathbf{B}^H\mathbf{A}\mathbf{U}$  is nonsingular (which implies that  $\text{rank } \mathbf{U} = k$ ). Then the matrices in (3.6) are well defined and the following statements hold:*

1.  $\mathbf{P}_B^2 = \mathbf{P}_B$ ,  $\mathbf{P}_B\mathbf{A}\mathbf{U} = \mathbf{0}$ , and  $\mathbf{U}^H\mathbf{B}^H\mathbf{P}_B = \mathbf{0}$ , i.e.,  $\mathbf{P}_B$  is the projection onto  $(\mathbf{B}\mathbf{U})^\perp$  along  $\mathbf{A}\mathbf{U}$ .
2.  $\mathbf{Q}_B^2 = \mathbf{Q}_B$ ,  $\mathbf{Q}_B\mathbf{U} = \mathbf{0}$ , and  $\mathbf{U}^H\mathbf{B}^H\mathbf{A}\mathbf{Q}_B = \mathbf{0}$ , i.e.,  $\mathbf{Q}_B$  is the projection onto  $(\mathbf{A}^H\mathbf{B}\mathbf{U})^\perp$  along  $\mathbf{U}$ .
3.  $\mathbf{P}_B\mathbf{A} = \mathbf{P}_B\mathbf{A}\mathbf{Q}_B = \mathbf{A}\mathbf{Q}_B$ .
4.  $\mathbf{P}_A = \mathbf{P}_A^H$ , i.e.,  $\mathbf{P}_A$  is an orthogonal projection.

It remains to impose the first orthogonality condition in (3.1), which will determine the vector  $\mathbf{y}_n$ . To this end, let

$$\hat{\mathbf{x}}_n := \mathbf{x}_0 + \mathbf{V}_n\mathbf{y}_n \in \mathbf{x}_0 + \mathcal{K}_n(\hat{\mathbf{A}}, \hat{\mathbf{v}}),$$

so that by (3.7)  $\mathbf{x}_n = \mathbf{Q}_B\hat{\mathbf{x}}_n + \mathbf{M}_B\mathbf{B}^H\mathbf{b}$ . Using the definition of  $\mathbf{P}_B$  in (3.6) and statement 3 of Lemma 3.1, this orthogonality condition reads

$$\mathbf{r}_n = \mathbf{b} - \mathbf{A}\mathbf{x}_n = \mathbf{b} - \mathbf{A}\mathbf{Q}_B\hat{\mathbf{x}}_n - \mathbf{A}\mathbf{M}_B\mathbf{B}^H\mathbf{b} = \mathbf{P}_B(\mathbf{b} - \mathbf{A}\hat{\mathbf{x}}_n) \perp \mathbf{B}\mathcal{K}_n(\hat{\mathbf{A}}, \hat{\mathbf{v}}).$$

We summarize these considerations in the following theorem.

THEOREM 3.2. *Let the assumptions of Lemma 3.1 hold and let  $\hat{\mathbf{A}} \in \mathbb{C}^{N \times N}$ ,  $\hat{\mathbf{v}} \in \mathbb{C}^N$ , and  $n \in \mathbb{N}$  be such that the Krylov subspace  $\mathcal{K}_n(\hat{\mathbf{A}}, \hat{\mathbf{v}})$  has dimension  $n$ . Furthermore, let  $\mathbf{b}, \mathbf{x}_0 \in \mathbb{C}^N$  be arbitrary.*

*Then, with  $\mathcal{U} := \text{im}(\mathbf{U})$  and the definitions from (3.6) the two pairs of conditions,*

$$\begin{aligned}
 \mathbf{x}_n &\in \mathbf{x}_0 + \mathcal{K}_n(\hat{\mathbf{A}}, \hat{\mathbf{v}}) + \mathcal{U}, \\
 \mathbf{r}_n &:= \mathbf{b} - \mathbf{A}\mathbf{x}_n \perp \mathbf{B}\mathcal{K}_n(\hat{\mathbf{A}}, \hat{\mathbf{v}}) + \mathbf{B}\mathcal{U}
 \end{aligned}
 \tag{3.9}$$

and

$$\begin{aligned}
 \hat{\mathbf{x}}_n &\in \mathbf{x}_0 + \mathcal{K}_n(\hat{\mathbf{A}}, \hat{\mathbf{v}}), \\
 \hat{\mathbf{r}}_n &:= \mathbf{P}_B(\mathbf{b} - \mathbf{A}\hat{\mathbf{x}}_n) \perp \mathbf{B}\mathcal{K}_n(\hat{\mathbf{A}}, \hat{\mathbf{v}})
 \end{aligned}
 \tag{3.10}$$

are equivalent for  $n \geq 1$  in the sense that

$$\mathbf{x}_n = \mathbf{Q}_B\hat{\mathbf{x}}_n + \mathbf{M}_B\mathbf{B}^H\mathbf{b} \quad \text{and} \quad \mathbf{r}_n = \hat{\mathbf{r}}_n.
 \tag{3.11}$$

We call (3.9) the *explicit* deflation and augmentation approach because the augmentation space  $\mathcal{U}$  is explicitly included in the search space. The equivalent conditions (3.10) show that the explicit inclusion of  $\mathcal{U}$  can be omitted when instead we first construct the iterate  $\widehat{\mathbf{x}}_n \in \mathbf{x}_0 + \mathcal{K}_n(\widehat{\mathbf{A}}, \widehat{\mathbf{v}})$  so that the projected residual  $\widehat{\mathbf{r}}_n = \mathbf{P}_B(\mathbf{b} - \widehat{\mathbf{A}}\widehat{\mathbf{x}}_n)$  satisfies the given orthogonality condition and then apply the affine correction (3.11) to  $\widehat{\mathbf{x}}_n$ , whose projected residual equals the one of  $\mathbf{x}_n$ . We call this second option the *implicit* deflation and augmentation approach.

Note that the theorem makes no assumption on relations between  $\widehat{\mathbf{A}}$ ,  $\widehat{\mathbf{v}}$ , and  $\mathcal{U}$ . The only assumption on the augmentation space  $\mathcal{U}$  is that the matrix  $\mathbf{U}^H \mathbf{B}^H \mathbf{A} \mathbf{U}$  is nonsingular. (Clearly, if this holds for one basis of  $\mathcal{U}$  it holds for all.) Moreover, in the theorem  $\widehat{\mathbf{A}}$  and  $\widehat{\mathbf{v}}$  are arbitrary except for the assumption that  $\mathcal{K}_n(\widehat{\mathbf{A}}, \widehat{\mathbf{v}})$  has dimension  $n$ .

In practice,  $\widehat{\mathbf{A}}$  and  $\widehat{\mathbf{v}}$  should be somehow related to  $\mathbf{A}$ , however. One specific choice is suggested by Theorem 3.2, in particular (3.10). If

$$\widehat{\mathbf{A}} := \mathbf{P}_B \mathbf{A}, \quad \widehat{\mathbf{v}} := \widehat{\mathbf{r}}_0 := \mathbf{P}_B \mathbf{r}_0 = \mathbf{P}_B(\mathbf{b} - \mathbf{A}\mathbf{x}_0), \quad \text{and} \quad \widehat{\mathbf{b}} := \mathbf{P}_B \mathbf{b},$$

then (3.10) becomes

$$(3.12) \quad \begin{aligned} \widehat{\mathbf{x}}_n &\in \mathbf{x}_0 + \mathcal{K}_n(\widehat{\mathbf{A}}, \widehat{\mathbf{r}}_0), \\ \widehat{\mathbf{r}}_n &:= \widehat{\mathbf{b}} - \widehat{\mathbf{A}}\widehat{\mathbf{x}}_n \perp \mathbf{B}\mathcal{K}_n(\widehat{\mathbf{A}}, \widehat{\mathbf{r}}_0), \end{aligned}$$

which is a formal Galerkin condition for the (consistent and singular) deflated system  $\widehat{\mathbf{A}}\widehat{\mathbf{x}} = \widehat{\mathbf{b}}$ . Based on the Jordan form of  $\mathbf{A}$  we show in the following theorem what the Jordan form of  $\widehat{\mathbf{A}} = \mathbf{P}_B \mathbf{A}$  looks like when (1)  $\mathcal{U}$  is a right invariant subspace or (2)  $\mathbf{B}\mathcal{U}$  is a left invariant subspace of  $\mathbf{A}$ .

**THEOREM 3.3.** *Suppose that the matrix  $\mathbf{A} \in \mathbb{C}^{N \times N}$  has a partitioned Jordan decomposition of the form*

$$(3.13) \quad \mathbf{A} = \mathbf{S} \mathbf{J} \mathbf{S}^{-1} = \begin{bmatrix} \mathbf{S}_1 & \mathbf{S}_2 \end{bmatrix} \begin{bmatrix} \mathbf{J}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{J}_2 \end{bmatrix} \begin{bmatrix} \widehat{\mathbf{S}}_1^H \\ \widehat{\mathbf{S}}_2^H \end{bmatrix},$$

where  $\mathbf{S}_1, \widehat{\mathbf{S}}_1 \in \mathbb{C}^{N \times k}$ ,  $\mathbf{S}_2, \widehat{\mathbf{S}}_2 \in \mathbb{C}^{N \times (N-k)}$ ,  $\mathbf{J}_1 \in \mathbb{C}^{k \times k}$ , and  $\mathbf{J}_2 \in \mathbb{C}^{(N-k) \times (N-k)}$ . Then the following assertions hold:

- (1) *If  $\mathcal{U} = \text{im}(\mathbf{S}_1)$ ,  $\mathbf{U} \in \mathbb{C}^{N \times k}$  is any matrix satisfying  $\text{im}(\mathbf{U}) = \mathcal{U}$ , and  $\mathbf{U}^H \mathbf{B}^H \mathbf{A} \mathbf{U}$  is nonsingular, then*

$$(3.14) \quad \begin{aligned} \widehat{\mathbf{A}} = \mathbf{P}_B \mathbf{A} &= \begin{bmatrix} \mathbf{U} & \mathbf{P}_B \mathbf{S}_2 \end{bmatrix} \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{J}_2 \end{bmatrix} \begin{bmatrix} \mathbf{U} & \mathbf{P}_B \mathbf{S}_2 \end{bmatrix}^{-1} \\ &\text{with} \quad \begin{bmatrix} \mathbf{U} & \mathbf{P}_B \mathbf{S}_2 \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{B} \mathbf{U} (\mathbf{U}^H \mathbf{B} \mathbf{U})^{-1} & \widehat{\mathbf{S}}_2 \end{bmatrix}^H. \end{aligned}$$

- (2) *If  $\mathbf{B}\mathcal{U} = \text{im}(\widehat{\mathbf{S}}_1)$ ,  $\mathbf{U} \in \mathbb{C}^{N \times k}$  is any matrix satisfying  $\text{im}(\mathbf{U}) = \mathcal{U}$ , and  $\mathbf{U}^H \mathbf{B}^H \mathbf{A} \mathbf{U}$  is nonsingular, then*

$$(3.15) \quad \begin{aligned} \widehat{\mathbf{A}} = \mathbf{P}_B \mathbf{A} &= \begin{bmatrix} \mathbf{U} & \mathbf{S}_2 \end{bmatrix} \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{J}_2 \end{bmatrix} \begin{bmatrix} \mathbf{U} & \mathbf{S}_2 \end{bmatrix}^{-1} \\ &\text{with} \quad \begin{bmatrix} \mathbf{U} & \mathbf{S}_2 \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{B} \mathbf{U} (\mathbf{U}^H \mathbf{B} \mathbf{U})^{-1} & \mathbf{Q}_B^H \widehat{\mathbf{S}}_2 \end{bmatrix}^H. \end{aligned}$$

In particular, in both cases the spectrum  $\Lambda(\widehat{\mathbf{A}})$  of  $\widehat{\mathbf{A}}$  is given by  $\Lambda(\widehat{\mathbf{A}}) = \{0\} \cup \Lambda(\mathbf{J}_2)$ .



*Proof.* (1) From Lemma 3.1 we can see that

$$\mathbf{P}_B \mathbf{A} \mathbf{U} = \mathbf{0} \quad \text{and} \quad (\mathbf{B} \mathbf{U} (\mathbf{U}^H \mathbf{B} \mathbf{U})^{-1})^H \mathbf{P}_B \mathbf{A} = \mathbf{0}.$$

By construction, there exists a nonsingular matrix  $\mathbf{R} \in \mathbb{C}^{k \times k}$  with  $\mathbf{A} \mathbf{U} = \mathbf{U} \mathbf{R}$ . Hence  $\mathbf{P}_B = \mathbf{I} - \mathbf{U} (\mathbf{U}^H \mathbf{B}^H \mathbf{U})^{-1} \mathbf{U}^H \mathbf{B}^H$  and from  $\widehat{\mathbf{S}}_2^H \mathbf{S}_1 = \mathbf{0}$  we conclude that  $\widehat{\mathbf{S}}_2^H \mathbf{U} = \mathbf{0}$  and thus

$$\widehat{\mathbf{S}}_2^H \mathbf{P}_B \mathbf{A} = \widehat{\mathbf{S}}_2^H \mathbf{A} = \mathbf{J}_2 \widehat{\mathbf{S}}_2^H.$$

Furthermore,

$$\mathbf{P}_B \mathbf{A} (\mathbf{P}_B \mathbf{S}_2) = \mathbf{P}_B \mathbf{A} \mathbf{S}_2 - \mathbf{P}_B \mathbf{A} \mathbf{U} (\mathbf{U}^H \mathbf{B}^H \mathbf{U})^{-1} \mathbf{U}^H \mathbf{B}^H \mathbf{S}_2 = (\mathbf{P}_B \mathbf{S}_2) \mathbf{J}_2$$

and the proof of (1) is complete after recognizing that

$$\begin{bmatrix} (\mathbf{U}^H \mathbf{B}^H \mathbf{U})^{-1} \mathbf{U}^H \mathbf{B}^H \\ \widehat{\mathbf{S}}_2^H \end{bmatrix} \begin{bmatrix} \mathbf{U} & \mathbf{P}_B \mathbf{S}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}.$$

The proof of (2) is analogous to (1).  $\square$

Results like the previous theorem motivate the term deflation, which means “making something smaller,” since the multiplication with the operator  $\mathbf{P}_B$  “removes” certain eigenvalues from the operator  $\mathbf{A}$  by “moving them to zero.” Special cases of the results shown in Theorem 3.3 have appeared in the literature: in particular, for spd matrices  $\mathbf{A}$  and  $\mathbf{B} = \mathbf{I}$  in the works of Frank and Vuik [22] and Nabben and Vuik [38, 39], and for nonsymmetric  $\mathbf{A}$  and  $\mathbf{B} = \mathbf{I}$  in the articles by Erlangga and Nabben [19, 20] and Yeung, Tang, and Vuik [61].

**4. Hpd matrices and CG.** This section presents some well-known results for deflation and augmentation techniques within the framework described in section 2 in the case where  $\mathbf{A}$  is Hpd. The first proposed deflated Krylov subspace methods for Hpd matrices are the deflated CG variants of Nicolaides [41] and Dostál [12]. With a full-rank matrix  $\mathbf{U} \in \mathbb{C}^{N \times k}$ ,  $\mathcal{U} = \text{im}(\mathbf{U})$  and  $\mathbf{B} = \mathbf{I}$  both essentially apply the CG method to the deflated system

$$(4.1) \quad \widehat{\mathbf{A}} \widehat{\mathbf{x}} = \widehat{\mathbf{b}}, \quad \text{where} \quad \widehat{\mathbf{A}} := \mathbf{P}_I \mathbf{A}, \quad \widehat{\mathbf{b}} := \mathbf{P}_I \mathbf{b}.$$

Here,  $\mathbf{P}_I = \mathbf{I} - \mathbf{A} \mathbf{U} (\mathbf{U}^H \mathbf{A} \mathbf{U})^{-1} \mathbf{U}^H$  is the projection onto  $\mathcal{U}^\perp$  along  $\mathbf{A} \mathcal{U}$  as defined in (3.6) when  $\mathbf{B} = \mathbf{I}$ . Moreover,  $\mathbf{Q}_I = \mathbf{P}_I^H$  is then the projection onto  $(\mathbf{A} \mathcal{U})^\perp$  along  $\mathcal{U}$ . Note that all the matrices in (3.6) are well defined because  $\mathbf{E}_I = \mathbf{U}^H \mathbf{A} \mathbf{U}$  is Hpd if  $\mathbf{A}$  is Hpd. Clearly, the deflated matrix  $\widehat{\mathbf{A}}$  is Hermitian but singular, since  $\mathbf{P}_I$  is a nontrivial projection if  $0 < k < N$ . In fact, this matrix  $\widehat{\mathbf{A}}$  is positive semidefinite, since

$$\mathbf{v}^H \widehat{\mathbf{A}} \mathbf{v} = \mathbf{v}^H \mathbf{P}_I \mathbf{A} \mathbf{v} = \mathbf{v}^H \mathbf{P}_I^2 \mathbf{A} \mathbf{v} = \mathbf{v}^H \mathbf{P}_I (\mathbf{P}_I \mathbf{A}) \mathbf{v} = \mathbf{v}^H \mathbf{P}_I (\mathbf{P}_I \mathbf{A})^H \mathbf{v} = \mathbf{v}^H \mathbf{P}_I \mathbf{A} \mathbf{P}_I^H \mathbf{v} \geq 0$$

holds for any  $\mathbf{v} \in \mathbb{C}^N$ . The system (4.1) is consistent since it results from a left-multiplication of the nonsingular system  $\mathbf{A} \mathbf{x} = \mathbf{b}$  by  $\mathbf{P}_I$ . (We note that in [41, 12] the application of the projection  $\mathbf{P}_I$  to  $\mathbf{b}$  is carried out implicitly by adapting the initial guess such that the initial residual is orthogonal to  $\mathcal{U} = \text{im}(\mathbf{U})$ .) The solution  $\mathbf{x}$  of  $\mathbf{A} \mathbf{x} = \mathbf{b}$  thus also solves  $\widehat{\mathbf{A}} \widehat{\mathbf{x}} = \widehat{\mathbf{b}}$ , but in (4.1) we replaced  $\mathbf{x}$  by  $\widehat{\mathbf{x}}$  to indicate the nonuniqueness of the solution. In fact, the general solution is  $\widehat{\mathbf{x}} = \mathbf{x} + \mathbf{h}$  with  $\mathbf{h} \in \mathcal{U}$

since  $\mathbf{P_I A h} = \mathbf{A Q_I h} = \mathbf{0}$  if and only if  $\mathbf{Q_I h} = \mathbf{0}$ , that is,  $\mathbf{h} \in \mathcal{U}$ ; see statement 2 of Lemma 3.1. The application of  $\mathbf{Q_I}$  in the final correction (3.11) will annihilate  $\mathbf{h}$ . Note that a deflation version including the final correction (3.11) is used by Frank and Vuik [22] and Nabben and Vuik [38, 39].

In the context of Hpd matrices the application of the CG method to a deflated system like (4.1) is a commonly used technique; see, e.g., [54] for a survey of results. Finally, we point out that  $\widehat{\mathbf{A}}$  as defined in (4.1) is completely determined by  $\mathbf{A}$  and the choice of the space  $\mathcal{U}$ .

According to Nicolaidis [41, section 3] and Kaasschieter [30, section 2], the CG method is well defined (in exact arithmetic) for each step  $n$  until it terminates with an exact solution when it is applied to a consistent linear algebraic system with a real and symmetric positive semidefinite matrix. This result easily generalizes to complex and Hermitian positive semidefinite matrices.

Mathematically, the  $n$ th step of the CG method applied to the deflated system (4.1) with the initial guess  $\mathbf{x}_0$  and the corresponding initial residual  $\widehat{\mathbf{r}}_0 = \widehat{\mathbf{b}} - \widehat{\mathbf{A}}\mathbf{x}_0$  is characterized by the two conditions

$$\begin{aligned}\widehat{\mathbf{x}}_n &\in \mathbf{x}_0 + \mathcal{K}_n(\widehat{\mathbf{A}}, \widehat{\mathbf{r}}_0), \\ \widehat{\mathbf{r}}_n &= \widehat{\mathbf{b}} - \widehat{\mathbf{A}}\widehat{\mathbf{x}}_n = \mathbf{P_I}(\mathbf{b} - \mathbf{A}\widehat{\mathbf{x}}_n) \perp \mathcal{K}_n(\widehat{\mathbf{A}}, \widehat{\mathbf{r}}_0).\end{aligned}$$

This is nothing but the set of conditions (3.10) in Theorem 3.2 with  $\mathbf{B} = \mathbf{I}$ . In the sense of relation (3.11) these conditions have been shown to be equivalent to (3.9), namely,

$$\begin{aligned}\mathbf{x}_n &\in \mathbf{x}_0 + \mathcal{K}_n(\widehat{\mathbf{A}}, \widehat{\mathbf{r}}_0) + \mathcal{U}, \\ \mathbf{r}_n &= \mathbf{b} - \mathbf{A}\mathbf{x}_n \perp \mathcal{K}_n(\widehat{\mathbf{A}}, \widehat{\mathbf{r}}_0) + \mathcal{U},\end{aligned}$$

which is the starting point of the theory for the deflated CG method developed in [48], where the authors also showed the equivalence between CG with explicit augmentation and CG applied to the deflated system (4.1); see section 4 in [48], in particular Theorem 4.6. In a partly similar treatment, Erhel and Guyomarc'h [18] considered an augmented and deflated CG method where the augmentation space  $\mathcal{U}$  is itself a Krylov space. It is worth mentioning that both Saad et al. [48, equation (3.12)] and Erhel and Guyomarc'h [18, equation (3.2)] use the initial correction

$$\mathbf{x}_0 := \mathbf{x}_{-1} + \mathbf{M_I r}_{-1} \quad \text{with} \quad \mathbf{r}_{-1} := \mathbf{b} - \mathbf{A}\mathbf{x}_{-1}$$

to replace a given initial approximation  $\mathbf{x}_{-1}$  by one with  $\mathbf{r}_0 \perp \mathcal{U}$ ; in fact, it is easily seen that  $\mathbf{r}_0 = \mathbf{P_I r}_{-1}$ .

The goal of deflation is to obtain a deflated matrix  $\widehat{\mathbf{A}}$  whose effective condition number is smaller than the one of  $\mathbf{A}$ , for example, by “eliminating” the smallest eigenvalues of  $\mathbf{A}$ . A detailed analysis of spectral properties of  $\mathbf{P_I A}$  and other projection-type preconditioners arising from domain decomposition and multigrid methods was carried out in [39] and [54]. In particular, it was shown in these papers that the effective condition number of  $\widehat{\mathbf{A}}$  is less than or equal to the condition number of  $\mathbf{A}$  for any augmentation space  $\mathcal{U}$ . Moreover, if  $\Lambda = \Lambda(\mathbf{A})$  is the spectrum of  $\mathbf{A}$  and  $\mathcal{U}$  is an  $\mathbf{A}$ -invariant subspace associated with the eigenvalues  $\Theta = \{\theta_1, \dots, \theta_k\} \subset \Lambda$ , then the effective 2-norm condition number is

$$\kappa_2(\widehat{\mathbf{A}}) = \frac{\max_{\lambda \in \Lambda \setminus \Theta} \lambda}{\min_{\lambda \in \Lambda \setminus \Theta} \lambda}.$$

In summary, for any augmentation space  $\mathcal{U}$ , the CG method applied to the (singular) deflated system (4.1) is well defined for any iteration step  $n$ , and it terminates with an exact solution  $\hat{\mathbf{x}}$  (in exact arithmetic). Once CG has terminated with a solution  $\hat{\mathbf{x}}$  of the deflated system, we can obtain the uniquely defined solution of the original system using the final correction step

$$\mathbf{x} = \mathbf{Q}_I \hat{\mathbf{x}} + \mathbf{M}_I \mathbf{b}$$

(cf. (3.11)), which indeed gives

$$\mathbf{A}\mathbf{x} = \mathbf{A}\mathbf{Q}_I \hat{\mathbf{x}} + \mathbf{A}\mathbf{M}_I \mathbf{b} = \mathbf{P}_I \mathbf{A} \hat{\mathbf{x}} + \mathbf{A}\mathbf{M}_I \mathbf{b} = (\mathbf{P}_I + \mathbf{A}\mathbf{M}_I) \mathbf{b} = \mathbf{b}.$$

This computation is mathematically equivalent to an explicit use of augmentation. Of course, in practice we stop the CG iteration for the deflated system once the solution is approximated sufficiently accurately. We then use the computed approximation  $\hat{\mathbf{x}}_n$  and (3.11) from Theorem 3.2 to obtain an approximation  $\mathbf{x}_n$  of the solution of the given system  $\mathbf{A}\mathbf{x} = \mathbf{b}$ . Note that, according to (3.11), the residual  $\hat{\mathbf{r}}_n = \hat{\mathbf{b}} - \hat{\mathbf{A}}\hat{\mathbf{x}}_n$  of the projected system (4.1) is equal to the residual  $\mathbf{r}_n = \mathbf{b} - \mathbf{A}\mathbf{x}_n$  of the original system (1.1).

**5. Non-Hermitian matrices and GMRES.** In this section we present mostly known results on applying versions of deflated GMRES to a general nonsingular matrix  $\mathbf{A}$ . We set  $\mathbf{B} = \mathbf{A}$  in the framework of section 2 and discuss some choices for  $\hat{\mathbf{A}}$ ,  $\hat{\mathbf{v}}$ , and  $\mathcal{U}$ .

Morgan [36] and also Chapman and Saad [6] presented variations of GMRES that can be mathematically described by (3.9) with  $\hat{\mathbf{A}} = \mathbf{A}$  and  $\hat{\mathbf{v}} = \mathbf{b} - \mathbf{A}\mathbf{x}_0$ . Hence they augmented the search space with an augmentation space  $\mathcal{U}$  but neither deflated the matrix nor projected the linear system onto a subspace of  $\mathbb{C}^N$ .

Erlangga and Nabben [19] used two matrices  $\mathbf{Y}, \mathbf{Z} \in \mathbb{C}^{N \times k}$  to define the abstract deflation operator  $\mathbf{P}_{\mathbf{YZ}} := \mathbf{I} - \mathbf{AZ}(\mathbf{Y}^H \mathbf{AZ})^{-1} \mathbf{Y}^H$  for non-Hermitian matrices  $\mathbf{A}$ . Of course, this choice needs the assumption of nonsingularity of  $\mathbf{Y}^H \mathbf{AZ}$ . Requiring  $\mathbf{Y}$  and  $\mathbf{Z}$  to have full rank obviously is not sufficient. They then applied GMRES to the deflated linear system  $\mathbf{P}_{\mathbf{YZ}} \mathbf{A} \hat{\mathbf{x}} = \mathbf{P}_{\mathbf{YZ}} \mathbf{b}$ .

De Sturler [10] introduced the GCRO method, which is a nested Krylov subspace method involving an outer and an inner iteration. The outer method is the GCR method [15, 16], while the inner iteration uses the projection

$$\mathbf{P}_A = \mathbf{I} - \mathbf{A}\mathbf{M}_A \mathbf{A}^H = \mathbf{I} - \mathbf{A}\mathbf{U}(\mathbf{U}^H \mathbf{A}^H \mathbf{A}\mathbf{U})^{-1} \mathbf{U}^H \mathbf{A}^H$$

to apply several steps of GMRES to the projected (or deflated) linear system

$$(5.1) \quad \hat{\mathbf{A}} \hat{\mathbf{x}} = \hat{\mathbf{b}}, \quad \text{where } \hat{\mathbf{A}} := \mathbf{P}_A \mathbf{A}, \quad \hat{\mathbf{b}} := \mathbf{P}_A \mathbf{b}.$$

In GCRO the matrix  $\mathbf{U}$  is determined from the corrections of the outer iteration. Clearly, the matrix  $\mathbf{E}_A = \mathbf{U}^H \mathbf{A}^H \mathbf{A}\mathbf{U}$  is nonsingular for any matrix  $\mathbf{U} \in \mathbb{C}^{N \times k}$  with rank  $\mathbf{U} = k > 0$ , so that all matrices in (3.6) are well defined. Note that the projection  $\mathbf{P}_A$  is equal to the abstract deflation operator  $\mathbf{P}_{\mathbf{YZ}}$  of Erlangga and Nabben with the choice  $\mathbf{Z} = \mathbf{U}$  and  $\mathbf{Y} = \mathbf{A}\mathbf{U}$ . For the application of  $\mathbf{P}_A$  only the matrix  $\mathbf{W} := \mathbf{A}\mathbf{U}$  is needed because  $\mathbf{P}_A = \mathbf{I} - \mathbf{W}(\mathbf{W}^H \mathbf{W})^{-1} \mathbf{W}^H$ . De Sturler further simplified this in [10, section 2] to  $\mathbf{P}_A = \mathbf{I} - \mathbf{C}\mathbf{C}^H$  by choosing a matrix  $\mathbf{C} \in \mathbb{C}^{N \times k}$  whose columns form an orthonormal basis of  $\text{im}(\mathbf{A}\mathbf{U})$ .

Here we concentrate on the GMRES method applied to the deflated system (5.1), and we first discuss some known results within the framework presented in section 2. Analogously to the approach for CG described in the previous section, the deflated system (5.1) results from the given system  $\mathbf{Ax} = \mathbf{b}$  by a left-multiplication with  $\mathbf{P}_A$  which projects onto  $(\mathbf{AU})^\perp$  along  $\mathbf{AU}$ . Note that  $\mathbf{P}_A$  is an orthogonal projection, since  $\mathbf{P}_A$  is Hermitian.

If we start GMRES with an initial guess  $\mathbf{x}_0$  and the corresponding initial residual  $\hat{\mathbf{r}}_0 = \hat{\mathbf{b}} - \hat{\mathbf{A}}\mathbf{x}_0 = \mathbf{P}_A(\mathbf{b} - \mathbf{Ax}_0)$ , then the iterate  $\hat{\mathbf{x}}_n$  and the residual  $\hat{\mathbf{r}}_n$  are characterized by the two conditions

$$\hat{\mathbf{x}}_n \in \mathbf{x}_0 + \mathcal{K}_n(\hat{\mathbf{A}}, \hat{\mathbf{r}}_0) \quad \text{and} \quad \hat{\mathbf{r}}_n = \hat{\mathbf{b}} - \hat{\mathbf{A}}\hat{\mathbf{x}}_n \perp \hat{\mathbf{A}}\mathcal{K}_n(\hat{\mathbf{A}}, \hat{\mathbf{r}}_0).$$

If the columns of  $\mathbf{V}_n$  form a basis of  $\mathcal{K}_n(\hat{\mathbf{A}}, \hat{\mathbf{r}}_0)$ , then the second condition means that

$$\begin{aligned} \mathbf{0} &= \mathbf{V}_n^H \hat{\mathbf{A}}^H \hat{\mathbf{r}}_n = \mathbf{V}_n^H \mathbf{A}^H \mathbf{P}_A^H \hat{\mathbf{r}}_n = \mathbf{V}_n^H \mathbf{A}^H \mathbf{P}_A \mathbf{P}_A (\mathbf{b} - \mathbf{Ax}_n) = \mathbf{V}_n^H \mathbf{A}^H \mathbf{P}_A (\mathbf{b} - \mathbf{Ax}_n) \\ &= \mathbf{V}_n^H \mathbf{A}^H \hat{\mathbf{r}}_n, \end{aligned}$$

or, equivalently,

$$\hat{\mathbf{r}}_n \perp \mathbf{A}\mathcal{K}_n(\hat{\mathbf{A}}, \hat{\mathbf{r}}_0).$$

Note that here the Krylov subspace is multiplied with  $\mathbf{A}$  instead of  $\hat{\mathbf{A}}$  and that this condition has precisely the form of the second condition in (3.10). Theorem 3.2 now implies that the mathematical characterization of GMRES applied to the deflated system  $\hat{\mathbf{A}}\hat{\mathbf{x}} = \hat{\mathbf{b}}$  is equivalent to the explicit use of augmentation, i.e., the conditions

$$(5.2) \quad \mathbf{x}_n \in \mathbf{x}_0 + \mathcal{K}_n(\hat{\mathbf{A}}, \hat{\mathbf{r}}_0) + \mathcal{U},$$

$$(5.3) \quad \mathbf{r}_n = \mathbf{b} - \mathbf{Ax}_n \perp \mathbf{A}\mathcal{K}_n(\hat{\mathbf{A}}, \hat{\mathbf{r}}_0) + \mathbf{AU},$$

in the sense that

$$(5.4) \quad \mathbf{x}_n = \mathbf{Q}_A \hat{\mathbf{x}}_n + \mathbf{M}_A \mathbf{A}^H \mathbf{b} \quad \text{and} \quad \mathbf{r}_n = \mathbf{b} - \mathbf{Ax}_n = \hat{\mathbf{b}} - \hat{\mathbf{A}}\hat{\mathbf{x}}_n = \hat{\mathbf{r}}_n.$$

As mentioned in the beginning of section 2, conditions (5.2)–(5.3) are equivalent to the minimization problem

$$\text{find } \mathbf{x}_n \in \mathbf{x}_0 + \mathcal{S}_n \quad \text{s.t.} \quad \|\mathbf{b} - \mathbf{Ax}_n\|_2 = \min_{\mathbf{y} \in \mathbf{x}_0 + \mathcal{S}_n} \|\mathbf{b} - \mathbf{Ay}\|_2$$

with the search space  $\mathcal{S}_n = \mathcal{K}_n(\hat{\mathbf{A}}, \hat{\mathbf{r}}_0) + \mathcal{U}$ . In the setting of GCRO, where  $\mathbf{U}$  is determined from the GCR iteration, the equivalence between GMRES applied to  $\hat{\mathbf{A}}\hat{\mathbf{x}} = \hat{\mathbf{b}}$  and the minimization problem with an explicitly augmented search space has already been pointed out by de Sturler [10, Theorem 2.2]. The GCRO method was extended to an arbitrary rank- $k$  matrix  $\mathbf{U}$  in [32, section 2]. In the case where  $\mathcal{U}$  is an  $\mathbf{A}$ -invariant subspace the equivalence is straightforward and has been pointed out by Eiermann, Ernst, and Schneider [14, Lemma 4.3].

Again the deflated matrix  $\hat{\mathbf{A}}$  is singular, and we have to discuss whether the application of GMRES to the deflated system yields (in exact arithmetic) a well-defined sequence of iterates that terminates with a solution. This turns out to be significantly more difficult than in the case of the CG method. Properties of GMRES applied to singular systems have been analyzed by de Sturler [10] and by Brown and Walker [5]. The following result is an extension of [5, Theorem 2.6].

THEOREM 5.1. Consider an arbitrary matrix  $\widehat{\mathbf{A}} \in \mathbb{C}^{N \times N}$  and a vector  $\widehat{\mathbf{b}} \in \text{im}(\widehat{\mathbf{A}})$  (i.e., the linear algebraic system  $\widehat{\mathbf{A}}\widehat{\mathbf{x}} = \widehat{\mathbf{b}}$  is consistent). Then the following two conditions are equivalent:

1. For every initial guess  $\mathbf{x}_0 \in \mathbb{C}^N$  the GMRES method applied to the system  $\widehat{\mathbf{A}}\widehat{\mathbf{x}} = \widehat{\mathbf{b}}$  is well defined at each iteration step  $n$  and it terminates with a solution of the system.
2.  $\ker(\widehat{\mathbf{A}}) \cap \text{im}(\widehat{\mathbf{A}}) = \{\mathbf{0}\}$ .

*Proof.* It has been shown in [5, Theorem 2.6] that condition 2 implies condition 1. We prove the reverse by contradiction. We assume that  $\ker(\widehat{\mathbf{A}}) \cap \text{im}(\widehat{\mathbf{A}}) \neq \{\mathbf{0}\}$ , and we will construct an initial guess for which GMRES does not terminate with the solution. For a nonzero vector  $\mathbf{y} \in \ker(\widehat{\mathbf{A}}) \cap \text{im}(\widehat{\mathbf{A}})$  there exists a nonzero vector  $\widehat{\mathbf{y}} \in \mathbb{C}^N$ , such that  $\mathbf{y} = \widehat{\mathbf{A}}\widehat{\mathbf{y}}$ , and since  $\widehat{\mathbf{A}}\widehat{\mathbf{x}} = \widehat{\mathbf{b}}$  is consistent, there exists a vector  $\widehat{\mathbf{x}} \in \mathbb{C}^N$  with  $\widehat{\mathbf{b}} = \widehat{\mathbf{A}}\widehat{\mathbf{x}}$ . Then the initial guess  $\mathbf{x}_0 := \widehat{\mathbf{x}} - \widehat{\mathbf{y}}$  gives  $\mathbf{r}_0 = \widehat{\mathbf{b}} - \widehat{\mathbf{A}}\mathbf{x}_0 = \widehat{\mathbf{b}} - \widehat{\mathbf{A}}\widehat{\mathbf{x}} + \widehat{\mathbf{A}}\widehat{\mathbf{y}} = \mathbf{y}$ . But since  $\mathbf{y} \in \ker(\widehat{\mathbf{A}})$ , we obtain  $\widehat{\mathbf{A}}\mathbf{r}_0 = \mathbf{0}$ , so that the GMRES method terminates at the first iteration with the approximation  $\mathbf{x}_0$ , for which  $\mathbf{r}_0 = \mathbf{y} \neq \mathbf{0}$ . Thus, for this particular initial guess  $\mathbf{x}_0$  the GMRES method cannot determine the solution of  $\widehat{\mathbf{A}}\widehat{\mathbf{x}} = \widehat{\mathbf{b}}$ .  $\square$

The situation that the GMRES method terminates without finding the exact solution is often called a *breakdown* of GMRES. The above proof leads to the following characterization of all initial guesses that lead to a breakdown of GMRES at the first iteration.

COROLLARY 5.2. Let  $\widehat{\mathbf{A}} \in \mathbb{C}^{N \times N}$  and  $\widehat{\mathbf{x}}, \widehat{\mathbf{b}} \in \mathbb{C}^{N \times N}$  such that  $\widehat{\mathbf{A}}\widehat{\mathbf{x}} = \widehat{\mathbf{b}}$ . Then the GMRES method breaks down at the first iteration for all initial guesses

$$\mathbf{x}_0 \in \mathcal{X}_0 := \{\widehat{\mathbf{x}} - \widehat{\mathbf{y}} \mid \widehat{\mathbf{A}}\widehat{\mathbf{y}} \in \ker(\widehat{\mathbf{A}}) \setminus \{\mathbf{0}\}\}.$$

We next have a closer look at condition 2 in Theorem 5.1. If we had  $\ker(\widehat{\mathbf{A}}) = \ker(\widehat{\mathbf{A}}^H)$ , then  $\text{im}(\widehat{\mathbf{A}})^\perp = \ker(\widehat{\mathbf{A}}^H)$  would imply

$$\{\mathbf{0}\} = \text{im}(\widehat{\mathbf{A}})^\perp \cap \text{im}(\widehat{\mathbf{A}}) = \ker(\widehat{\mathbf{A}}^H) \cap \text{im}(\widehat{\mathbf{A}}) = \ker(\widehat{\mathbf{A}}) \cap \text{im}(\widehat{\mathbf{A}}),$$

so that condition 2 would hold. Thus condition 2 in Theorem 5.1 is fulfilled for any Hermitian matrix  $\widehat{\mathbf{A}}$ . For a general non-Hermitian matrix, however, it seems difficult to determine a deflated matrix with  $\ker(\widehat{\mathbf{A}}) = \ker(\widehat{\mathbf{A}}^H)$ . However, for the deflated system (5.1) we can derive another condition that is equivalent with condition 2 (and hence condition 1) in Theorem 5.1.

COROLLARY 5.3. For the deflated system (5.1), condition 2 in Theorem 5.1 is satisfied if and only if  $\mathcal{U} \cap (\mathbf{A}\mathcal{U})^\perp = \{\mathbf{0}\}$ . In particular, the latter condition is satisfied when  $\mathcal{U}$  is an exactly  $\mathbf{A}$ -invariant subspace, i.e., when  $\mathbf{A}\mathcal{U} = \mathcal{U}$ .

*Proof.* Using the properties of the projection  $\mathbf{P}_\mathbf{A}$  from Lemma 3.1 and the fact that  $\mathbf{A}$  is nonsingular, we obtain

$$\begin{aligned} \ker(\widehat{\mathbf{A}}) &= \ker(\mathbf{P}_\mathbf{A}\mathbf{A}) = \mathbf{A}^{-1}\ker(\mathbf{P}_\mathbf{A}) = \mathcal{U}, \\ \text{im}(\widehat{\mathbf{A}}) &= \text{im}(\mathbf{P}_\mathbf{A}\mathbf{A}) = \text{im}(\mathbf{P}_\mathbf{A}) = (\mathbf{A}\mathcal{U})^\perp. \end{aligned}$$

If  $\mathbf{A}\mathcal{U} = \mathcal{U}$ , then  $\mathcal{U} \cap (\mathbf{A}\mathcal{U})^\perp = \{\mathbf{0}\}$  holds trivially.  $\square$

For a nonsingular matrix, condition 2 in Theorem 5.1 always holds trivially, and hence a breakdown of GMRES can only occur if the method is applied to a linear algebraic system with a singular matrix. (This fact has been known since the method's introduction in 1986 [47].) Breakdowns have also been analyzed by de Sturler [10]

in the context of the GCRO method. (See the end of section 6.1 below for further comments.) We want to point out that it is unlikely that a random initial guess lies in the subspace  $\mathcal{X}_0$  specified in Corollary 5.2 for which the GMRES method breaks down in the first step. However, a general  $\mathcal{U}$  may lead to a breakdown. To illustrate the problem of breakdowns in our context, we give an example that is adapted from [5, Example 1.1].

*Example 5.4.* Consider a linear algebraic system with

$$\mathbf{A} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 1 \\ 0 \end{bmatrix},$$

so that the unique solution is given by the vector  $[0, 1]^T$ . Let the augmentation space be defined by  $\mathbf{U}_1 = [1, 0]^T$ , and then

$$\mathbf{P}_A = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}, \quad \widehat{\mathbf{A}} = \mathbf{P}_A \mathbf{A} = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}, \quad \widehat{\mathbf{b}} = \mathbf{P}_A \mathbf{b} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}.$$

If  $\mathbf{x}_0$  is the zero vector, then  $\widehat{\mathbf{r}}_0 = \widehat{\mathbf{b}}$  and  $\widehat{\mathbf{A}}\widehat{\mathbf{r}}_0 = \mathbf{0}$ , and thus GMRES applied to the deflated system terminates at the very first iteration with the approximation  $\mathbf{x}_0$ . Since  $\widehat{\mathbf{A}}\mathbf{x}_0 \neq \widehat{\mathbf{b}}$ , this is a breakdown of GMRES. Furthermore, applying the correction (3.7) to  $\widehat{\mathbf{x}}_0 = \mathbf{x}_0$  does not yield the solution of the original system  $\mathbf{A}\mathbf{x} = \mathbf{b}$  because

$$\mathbf{Q}_A \mathbf{x}_0 + \mathbf{M}_A \mathbf{A}^H \mathbf{b} = \mathbf{M}_A \mathbf{A}^H \mathbf{b} = \mathbf{U}_1 \mathbf{U}_1^H \mathbf{A}^H \mathbf{b} = \mathbf{0} \neq \begin{bmatrix} 0 \\ 1 \end{bmatrix}.$$

Corollary 5.3 states that the GMRES method applied to the deflated system (5.1) cannot break down if  $\mathcal{U}$  is an  $\mathbf{A}$ -invariant subspace. The following example shows that care has also to be taken with approximate  $\mathbf{A}$ -invariant subspaces.

*Example 5.5.* Let  $\alpha > 0$  be a small positive number. Then  $\mathbf{v} := [0, 1, \alpha]^T$  is an eigenvector of the matrix

$$\mathbf{A} := \begin{bmatrix} 0 & 1 & -\alpha^{-1} \\ 1 & 0 & \alpha^{-1} \\ 0 & 0 & 1 \end{bmatrix}$$

corresponding to the eigenvalue 1. Instead of  $\mathbf{v}$  we use the perturbed vector  $\mathbf{U}_2 := [0, 1, 0]^T$  as a basis for the deflation space  $\mathcal{U} = \text{im}(\mathbf{U}_2)$  and obtain

$$\mathbf{A}\mathbf{U}_2 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \quad \mathbf{P}_A = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad \mathbf{P}_A \mathbf{A} = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & \alpha^{-1} \\ 0 & 0 & 1 \end{bmatrix}.$$

For  $\mathbf{x}, \mathbf{b} \in \mathbb{C}^3$  with  $\mathbf{A}\mathbf{x} = \mathbf{b}$  the GMRES method then breaks down in the first step for all  $\mathbf{x}_0 \in \{\mathbf{x} + \beta[1, 0, 0]^T \mid \beta \neq 0\}$ . Note that  $\|\mathbf{U}_2 - \mathbf{v}\|_2 = \alpha$  can be chosen arbitrarily small. A better measure for the quality of an approximate invariant subspace would be the largest principal angle between  $\mathcal{U}$  and  $\mathbf{A}\mathcal{U}$ .

**6. Hermitian matrices and variants of MINRES.** We will now apply the results presented in sections 2 and 5 to the case where  $\mathbf{A}$  is Hermitian, nonsingular, and possibly indefinite. For a Hermitian matrix the GMRES method considered in section 5 is mathematically equivalent to the MINRES method, which is based on the Hermitian Lanczos algorithm and thus uses efficient three-term recurrences.

**6.1. The RMINRES method.** This subsection discusses the RMINRES method developed by Wang, de Sturler, and Paulino [57]. This method fits into the framework of section 2, and the results presented in section 5 apply. Wang, de Sturler, and Paulino were interested in solving sequences of linear algebraic systems that exhibit only small changes from one matrix in the sequence to the next one, and they suggested to reuse information from previous solves. The RMINRES method consists of two main parts that can basically be analyzed separately: an augmented and deflated MINRES solver which is based on GCRO and an extraction procedure for the augmentation and deflation data. In the second part Wang, de Sturler, and Paulino determined harmonic Ritz vectors that correspond to harmonic Ritz values close to zero and used these approximate eigenspaces for augmenting the Krylov subspace. Here, we omit the extraction of the augmentation and deflation space and concentrate on the method for solving the systems. We refer to this as the solver part of the RMINRES method. We point out that the extracted spaces can be arbitrary if there are no restrictions on the changes of the matrices in the sequence of linear algebraic systems. However, the RMINRES method has been presented in [57] with an application in topology optimization where the extracted approximated eigenvectors of one matrix are still good approximations to eigenvectors of the next matrix. Furthermore, we will not address the preconditioning technique outlined in [57] and assume that the given linear algebraic system is already in the preconditioned form.

As in section 5, we set  $\mathbf{B} = \mathbf{A}$  and consider first the resulting deflated system of the form (5.1),

$$\widehat{\mathbf{A}}\widehat{\mathbf{x}} = \widehat{\mathbf{b}}, \quad \text{where } \widehat{\mathbf{A}} := \mathbf{P}_A\mathbf{A}, \quad \widehat{\mathbf{b}} := \mathbf{P}_A\mathbf{b}.$$

If we apply MINRES to this linear algebraic system with an initial guess  $\mathbf{x}_0$  and the corresponding initial residual  $\widehat{\mathbf{r}}_0 = \widehat{\mathbf{b}} - \widehat{\mathbf{A}}\mathbf{x}_0 = \mathbf{P}_A(\mathbf{b} - \mathbf{A}\mathbf{x}_0)$ , then the iterate  $\widehat{\mathbf{x}}_n$  and the residual  $\widehat{\mathbf{r}}_n$  are characterized by the two conditions

$$(6.1) \quad \widehat{\mathbf{x}}_n \in \mathbf{x}_0 + \mathcal{K}_n(\widehat{\mathbf{A}}, \widehat{\mathbf{r}}_0) \quad \text{and} \quad \widehat{\mathbf{r}}_n = \widehat{\mathbf{b}} - \widehat{\mathbf{A}}\widehat{\mathbf{x}}_n \perp \widehat{\mathbf{A}}\mathcal{K}_n(\widehat{\mathbf{A}}, \widehat{\mathbf{r}}_0).$$

This is essentially the approach of Kilmer and de Sturler [32, section 2]. Olshanskii and Simoncini [42] recently used a different approach where the MINRES method is applied to the deflated system  $\mathbf{P}_I\mathbf{A}\widehat{\mathbf{x}} = \mathbf{P}_I\mathbf{b}$  with the special initial guess  $\mathbf{x}_0 = \mathbf{U}(\mathbf{U}^H\mathbf{A}\mathbf{U})^{-1}\mathbf{U}^H\mathbf{b}$ . We note that the presentation in [42] is slightly different but the above can be seen with minor algebraic modifications to the relations in and preceding Proposition 3.1 in [42].

An attentive reader has certainly noticed that the deflated matrix  $\widehat{\mathbf{A}} = \mathbf{P}_A\mathbf{A} = \mathbf{A} - \mathbf{A}\mathbf{M}_A\mathbf{A}^2$  is in general *not* Hermitian, even when  $\mathbf{A}$  is Hermitian. However, as pointed out in [32, footnote on p. 2153] and [57, footnote on p. 2446], a straightforward computation shows that

$$(6.2) \quad \mathcal{K}_n(\mathbf{P}_A\mathbf{A}, \mathbf{P}_A\mathbf{v}) = \mathcal{K}_n(\mathbf{P}_A\mathbf{A}\mathbf{P}_A, \mathbf{P}_A\mathbf{v})$$

holds for every vector  $\mathbf{v} \in \mathbb{C}^N$  because  $\mathbf{P}_A$  is a projection. The matrix  $\mathbf{P}_A\mathbf{A}\mathbf{P}_A$  is obviously Hermitian (since  $\mathbf{A}$  and  $\mathbf{P}_A$  are Hermitian), and hence the Krylov subspaces we work with are also generated by a Hermitian matrix. It is therefore possible to implement a MINRES-like method for the deflated system, which is based on three-term recurrences and which is characterized by the conditions (6.1). As presented in section 5, these conditions combined with the correction step (5.4) are equivalent to the explicit use of augmentation, i.e., conditions (5.2)–(5.3).

The latter conditions are the basis of the solver part of the RMINRES method by Wang, de Sturler, and Paulino in [57, section 3]. We summarize the above and give two mathematically equivalent characterizations of the RMINRES solver applied to the system  $\mathbf{Ax} = \mathbf{b}$  with an initial guess  $\mathbf{x}_0$ :

1. The original approach used in [57] incorporates explicit augmentation, which means to construct iterates  $\mathbf{x}_n$  satisfying the two conditions

$$(6.3) \quad \begin{aligned} \mathbf{x}_n &\in \mathbf{x}_0 + \mathcal{K}_n(\mathbf{P}_A \mathbf{A}, \mathbf{P}_A \mathbf{r}_0) + \mathcal{U}, \\ \mathbf{r}_n &= \mathbf{b} - \mathbf{Ax}_n \perp \mathbf{A} \mathcal{K}_n(\mathbf{P}_A \mathbf{A}, \mathbf{P}_A \mathbf{r}_0) + \mathbf{AU}. \end{aligned}$$

2. A mathematically equivalent approach is to apply MINRES to the deflated system

$$(6.4) \quad \mathbf{P}_A \mathbf{A} \hat{\mathbf{x}} = \mathbf{P}_A \mathbf{b}$$

and correct the resulting iterates  $\hat{\mathbf{x}}_n$  according to  $\mathbf{x}_n = \mathbf{Q}_A \hat{\mathbf{x}}_n + \mathbf{M}_A \mathbf{Ab}$ .

Note that on an algorithmic level the second approach exhibits lower computational cost since the correction in the space  $\mathcal{U}$  is only carried out once at the end, while the RMINRES solver requires one update per iteration.

Since the solver part of RMINRES is mathematically equivalent to MINRES (and GMRES) applied to the deflated system, Corollary 5.3 also applies to RMINRES. In particular, the method can break down for specific initial guesses if (and only if)  $\mathcal{U} \cap (\mathbf{AU})^\perp \neq \{\mathbf{0}\}$ . Breakdowns cannot occur if  $\mathcal{U}$  is an exact  $\mathbf{A}$ -invariant subspace, but this is an unrealistic assumption in practical applications. Note that the matrix  $\mathbf{A}$  in Example 5.4 is Hermitian, thus it also serves as an example for a breakdown of the RMINRES solver. That the RMINRES method can break down may already be guessed from the fact that this method is based on the GCRO method and thus potentially suffers from the breakdown conditions for GCRO derived in [10]. However, the possibility of breakdowns has not been mentioned in [57], and in the example of a GCRO breakdown given in [10] the matrix  $\mathbf{A}$  is not Hermitian. Hence this example cannot be used in the context of the RMINRES method, which is intended for Hermitian matrices.

In the next subsection we show how to suitably modify the RMINRES approach to avoid breakdowns.

**6.2. Avoiding breakdowns in deflated MINRES.** We have seen in section 5 that if  $\ker(\hat{\mathbf{A}}) = \ker(\hat{\mathbf{A}}^H)$ , then condition 1 in Theorem 5.1 is satisfied. Consequently, if we can determine a Hermitian deflated matrix  $\hat{\mathbf{A}}$  and a corresponding consistent deflated system, MINRES applied to this system cannot break down for any initial guess.

Using the projections  $\mathbf{P}_A$  and  $\mathbf{Q}_A$  from (3.6) we decompose the solution  $\mathbf{x}$  of  $\mathbf{Ax} = \mathbf{b}$  as

$$(6.5) \quad \mathbf{x} = \mathbf{P}_A \mathbf{x} + (\mathbf{I} - \mathbf{P}_A) \mathbf{x} = \mathbf{P}_A \mathbf{x} + \mathbf{AM}_A \mathbf{Ax} = \mathbf{P}_A \mathbf{x} + \mathbf{AM}_A \mathbf{b},$$

$$(6.6) \quad \mathbf{x} = \mathbf{Q}_A \mathbf{x} + (\mathbf{I} - \mathbf{Q}_A) \mathbf{x} = \mathbf{Q}_A \mathbf{x} + \mathbf{M}_A \mathbf{A}^2 \mathbf{x} = \mathbf{Q}_A \mathbf{x} + \mathbf{M}_A \mathbf{Ab}.$$

Using (6.6), the system  $\mathbf{Ax} = \mathbf{b}$  becomes  $\mathbf{A}(\mathbf{Q}_A \mathbf{x} + \mathbf{M}_A \mathbf{Ab}) = \mathbf{b}$ . With the definition of  $\mathbf{P}_A$  and  $\mathbf{AQ}_A = \mathbf{P}_A \mathbf{A}$  (cf. Lemma 3.1) we see that this is equivalent to

$$\mathbf{P}_A \mathbf{Ax} = \mathbf{P}_A \mathbf{b}.$$



We now substitute for  $\mathbf{x}$  from (6.5) and obtain  $\mathbf{P}_A \mathbf{A}(\mathbf{P}_A \mathbf{x} + \mathbf{A} \mathbf{M}_A \mathbf{b}) = \mathbf{P}_A \mathbf{b}$  which is equivalent to

$$(6.7) \quad \mathbf{P}_A \mathbf{A} \mathbf{P}_A \mathbf{x} = \mathbf{P}_A \mathbf{Q}_A^H \mathbf{b}.$$

We can show the following result for the MINRES method applied to this symmetric system.

**THEOREM 6.1.** *For each initial guess  $\mathbf{x}_0 \in \mathbb{C}^N$  the MINRES method applied to the system (6.7) yields (in exact arithmetic) a well-defined iterate  $\bar{\mathbf{x}}_n$  at every step  $n \geq 1$  until it terminates with a solution. Moreover, the sequence of iterates*

$$(6.8) \quad \mathbf{x}_n := \mathbf{Q}_A (\mathbf{P}_A \bar{\mathbf{x}}_n + \mathbf{A} \mathbf{M}_A \mathbf{b}) + \mathbf{M}_A \mathbf{A} \mathbf{b}$$

*is well defined. It terminates (in exact arithmetic) with the exact solution  $\mathbf{x}$  of the original linear system  $\mathbf{A} \mathbf{x} = \mathbf{b}$ , and its residuals are given by  $\mathbf{r}_n = \mathbf{b} - \mathbf{A} \mathbf{x}_n = \mathbf{P}_A \mathbf{Q}_A^H \mathbf{b} - \mathbf{P}_A \mathbf{A} \mathbf{P}_A \bar{\mathbf{x}}_n$ .*

*Proof.* The first part follows from the fact that the system (6.7) is a consistent system with a Hermitian matrix  $\mathbf{P}_A \mathbf{A} \mathbf{P}_A$ , so that we can apply Theorem 5.1. It remains to show the second part. The  $n$ th residual of the original system  $\mathbf{A} \mathbf{x} = \mathbf{b}$  is given by

$$\begin{aligned} \mathbf{r}_n &= \mathbf{b} - \mathbf{A} \mathbf{x}_n = \mathbf{b} - \mathbf{A} (\mathbf{Q}_A (\mathbf{P}_A \bar{\mathbf{x}}_n + \mathbf{A} \mathbf{M}_A \mathbf{b}) + \mathbf{M}_A \mathbf{A} \mathbf{b}) \\ &= \mathbf{b} - \mathbf{A} \mathbf{Q}_A (\mathbf{P}_A \bar{\mathbf{x}}_n + \mathbf{A} \mathbf{M}_A \mathbf{b}) - \mathbf{A} \mathbf{M}_A \mathbf{A} \mathbf{b} \\ &= (\mathbf{I} - \mathbf{A} \mathbf{M}_A \mathbf{A}) \mathbf{b} - \mathbf{P}_A \mathbf{A} (\mathbf{P}_A \bar{\mathbf{x}}_n + \mathbf{A} \mathbf{M}_A \mathbf{b}) \\ &= \mathbf{P}_A \mathbf{b} - \mathbf{P}_A \mathbf{A} \mathbf{P}_A \bar{\mathbf{x}}_n - \mathbf{P}_A \mathbf{A}^2 \mathbf{M}_A \mathbf{b} \\ &= \mathbf{P}_A (\mathbf{I} - \mathbf{A}^2 \mathbf{M}_A) \mathbf{b} - \mathbf{P}_A \mathbf{A} \mathbf{P}_A \bar{\mathbf{x}}_n \\ &= \mathbf{P}_A \mathbf{Q}_A^H \mathbf{b} - \mathbf{P}_A \mathbf{A} \mathbf{P}_A \bar{\mathbf{x}}_n. \end{aligned}$$

We see that  $\mathbf{r}_n$  is equal to the  $n$ th MINRES residual for the system (6.7). In particular, this implies that the exact solution of (1.1) is given by (6.8) once an exact solution  $\bar{\mathbf{x}}_n$  of (6.7) has been determined by MINRES.  $\square$

When MINRES is applied to the deflated system (6.7), the Hermitian iteration matrix  $\mathbf{P}_A \mathbf{A} \mathbf{P}_A$  can again be replaced by the non-Hermitian matrix  $\mathbf{P}_A \mathbf{A}$  (cf. section 6.1).

The following theorem shows that a modification of the initial guess suffices to make the solver part of the RMINRES method mathematically equivalent to MINRES applied to the system (6.7).

**THEOREM 6.2.** *We consider the following two approaches:*

1. *The solver part of the RMINRES method applied to  $\mathbf{A} \mathbf{x} = \mathbf{b}$  with the initial guess  $\hat{\mathbf{x}}_0 := \mathbf{P}_A \mathbf{x}_0 + \mathbf{A} \mathbf{M}_A \mathbf{b}$  and resulting iterates  $\mathbf{x}_n$  and residuals  $\mathbf{r}_n = \mathbf{b} - \mathbf{A} \mathbf{x}_n$ .*
2. *The MINRES method applied to (6.7) with the initial guess  $\mathbf{x}_0$  and resulting iterates  $\bar{\mathbf{x}}_n$  and residuals  $\bar{\mathbf{r}}_n := \mathbf{P}_A \mathbf{Q}_A^H \mathbf{b} - \mathbf{P}_A \mathbf{A} \mathbf{P}_A \bar{\mathbf{x}}_n$ .*

*Both approaches are equivalent in the sense that  $\mathbf{x}_n = \mathbf{Q}_A (\mathbf{P}_A \bar{\mathbf{x}}_n + \mathbf{A} \mathbf{M}_A \mathbf{b}) + \mathbf{M}_A \mathbf{A} \mathbf{b}$  and  $\mathbf{r}_n = \bar{\mathbf{r}}_n$ .*

*Proof.* Let us start with the MINRES method applied to (6.7), which constructs iterates  $\bar{\mathbf{x}}_n = \mathbf{x}_0 + \mathbf{V}_n \mathbf{y}_n$ , where  $\mathbf{V}_n \in \mathbb{C}^{N \times n}$  is of full rank  $n$  such that  $\text{im}(\mathbf{V}_n) = \mathcal{K}_n(\mathbf{P}_A \mathbf{A} \mathbf{P}_A, \mathbf{P}_A \mathbf{Q}_A^H \mathbf{r}_0)$ . Then  $\mathbf{P}_A \mathbf{V}_n = \mathbf{V}_n$  and the corrected iterates are

$$(6.9) \quad \begin{aligned} \mathbf{x}_n &= \mathbf{Q}_A (\mathbf{P}_A (\mathbf{x}_0 + \mathbf{V}_n \mathbf{y}_n) + \mathbf{A} \mathbf{M}_A \mathbf{b}) + \mathbf{M}_A \mathbf{A} \mathbf{b} = \mathbf{Q}_A (\hat{\mathbf{x}}_0 + \mathbf{V}_n \mathbf{y}_n) + \mathbf{M}_A \mathbf{A} \mathbf{b} \\ &= \mathbf{Q}_A \hat{\mathbf{x}}_n + \mathbf{M}_A \mathbf{A} \mathbf{b} \end{aligned}$$

with  $\widehat{\mathbf{x}}_n := \widehat{\mathbf{x}}_0 + \mathbf{V}_n \mathbf{y}_n$ . For  $n > 0$  the  $n$ th residual of  $\widehat{\mathbf{x}}_n$  with respect to the system (6.7) is

$$\begin{aligned}\bar{\mathbf{r}}_n &= \mathbf{P}_A \mathbf{Q}_A^H \mathbf{b} - \mathbf{P}_A \mathbf{A} \mathbf{P}_A \widehat{\mathbf{x}}_n = \mathbf{P}_A (\mathbf{Q}_A^H \mathbf{b} - \mathbf{A} \mathbf{P}_A \widehat{\mathbf{x}}_n) \\ &= \mathbf{P}_A (\mathbf{b} - \mathbf{A} (\mathbf{P}_A \widehat{\mathbf{x}}_0 + \mathbf{A} \mathbf{M}_A \mathbf{b} + \mathbf{V}_n \mathbf{y}_n)) = \mathbf{P}_A \mathbf{b} - \mathbf{P}_A \mathbf{A} \widehat{\mathbf{x}}_n =: \widehat{\mathbf{r}}_n.\end{aligned}$$

This is the residual of  $\widehat{\mathbf{x}}_n$  with respect to the system (6.4). We also have

$$\widehat{\mathbf{r}}_0 = \mathbf{P}_A \mathbf{b} - \mathbf{P}_A \mathbf{A} \widehat{\mathbf{x}}_0 = \mathbf{P}_A \mathbf{Q}_A^H \mathbf{b} - \mathbf{P}_A \mathbf{A} \mathbf{P}_A \mathbf{x}_0 = \bar{\mathbf{r}}_0,$$

and thus the starting vectors of the Krylov subspace for both methods are equal. Because of (6.2) the Krylov subspaces are also equal. From the definition of the Krylov subspaces we immediately obtain

$$\bar{\mathbf{r}}_n \perp \mathbf{P}_A \mathbf{A} \mathbf{P}_A \mathcal{K}_n(\mathbf{P}_A \mathbf{A} \mathbf{P}_A, \bar{\mathbf{r}}_0) \iff \widehat{\mathbf{r}}_n \perp \mathbf{P}_A \mathcal{K}_n(\mathbf{P}_A \mathbf{A}, \widehat{\mathbf{r}}_0).$$

We can now see that the iterates  $\widehat{\mathbf{x}}_n$  are the iterates of MINRES applied to (6.4) with the initial guess  $\widehat{\mathbf{x}}_0$ . Along with the correction (6.9) this was shown to be equivalent to the RMINRES solver applied to  $\mathbf{A} \mathbf{x} = \mathbf{b}$  with the initial guess  $\widehat{\mathbf{x}}_0$  (cf. section 6.1).  $\square$

This means that (in exact arithmetic) breakdowns in the solver part of the RMINRES method can be prevented by either adapting the right-hand side to  $\mathbf{Q}^H \mathbf{b}$  and correcting the approximate solution at the end according to Theorem 6.1 or by choosing the adapted initial guess  $\widehat{\mathbf{x}}_0$  defined in Theorem 6.2. Both choices do not increase the computational cost significantly since these computations only need to be carried out once. A similar special initial guess has also been used in [54] to obtain a robust deflation-based preconditioner for the CG method; compare the A-DEF2 method in [54, Table 2].

**6.3. Numerical experiments.** In this subsection, we will show the numerical behavior of selected Krylov subspace methods discussed above. Detailed numerical experiments with the deflated CG method (cf. section 4) and equivalent approaches have been presented in [54]. Here, we will focus on the solver part of the RMINRES method and the deflated MINRES method in order to numerically illustrate the phenomenon of breakdowns that have only been described theoretically so far (cf. sections 6.1 and 6.2). Both methods are implemented in MATLAB with three-term Lanczos recurrences and Givens rotations for solving the least squares problem. All residuals have been computed explicitly in each iteration.

*Example 6.3.* In this example we use a matrix  $\mathbf{A} = \mathbf{W}^H \mathbf{D} \mathbf{W} \in \mathbb{R}^{2m \times 2m}$ ,  $m = 50$ , where  $\mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_{2m})$  with  $\lambda_j = \sqrt{j}$ ,  $\lambda_{m+j} = -\sqrt{j}$  for  $j = 1, \dots, m$  and  $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_{2m}]$  is a randomly generated orthogonal matrix. We consider a matrix  $\mathbf{U} = [u_1, \dots, u_k]$  whose columns are pairwise orthogonal eigenvectors of  $\mathbf{A}$ , i.e.,  $\mathbf{A} \mathbf{U} = \mathbf{U} \mathbf{D}_U$  and  $\mathbf{U}^H \mathbf{U} = \mathbf{I}_k$  with a diagonal matrix  $\mathbf{D}_U = \text{diag}(\lambda_{j_1}, \dots, \lambda_{j_k})$  for  $0 < j_1 < \dots < j_k < 2m$ . This means that  $\mathcal{U} = \text{im}(\mathbf{U})$  is an exact  $\mathbf{A}$ -invariant subspace. Then a straightforward computation reveals that  $\mathbf{P}_A = \mathbf{Q}_A = \mathbf{I} - \mathbf{U} \mathbf{U}^H$ , which is obviously Hermitian, and

$$\mathbf{P}_A \mathbf{A} \mathbf{P}_A = \mathbf{P}_A \mathbf{A} \mathbf{Q}_A = \mathbf{P}_A^2 \mathbf{A} = \mathbf{P}_A \mathbf{A}, \quad \mathbf{P}_A \mathbf{Q}_A^H = \mathbf{P}_A^2 = \mathbf{P}_A.$$

By comparing the correction steps of RMINRES and deflated MINRES (cf. sections 6.1 and 6.2) and using  $\mathbf{P}_A \mathbf{A} \mathbf{M}_A = \mathbf{0}$ , we can see that both methods are mathematically equivalent if  $\mathcal{U}$  is an exact invariant subspace.

We solve the system  $\mathbf{A} \mathbf{x} = \mathbf{b}$  with a random right-hand side  $\mathbf{b}$  and the initial guess  $\mathbf{x}_0 = \mathbf{0}$ . In Figure 6.1 we show the relative residual norms of the solvers

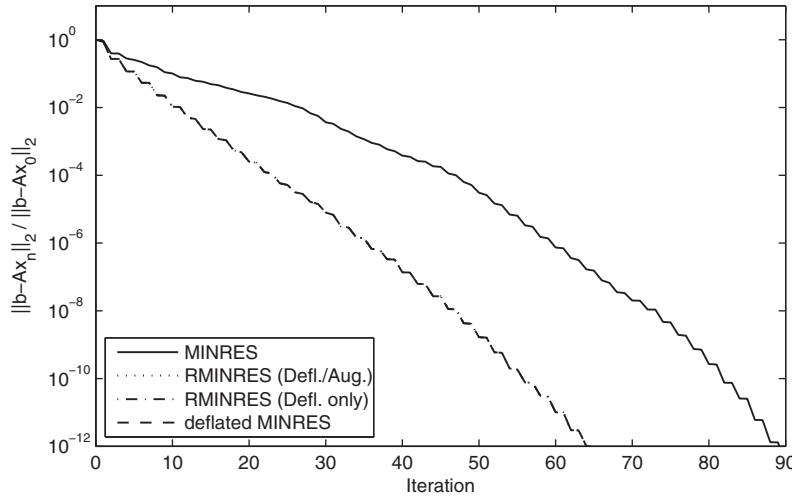


FIG. 6.1. Convergence history for Example 6.3. The convergence curves of both RMINRES solver implementations and the deflated MINRES method coincide.

- MINRES (solid line),
- RMINRES with explicit augmentation and deflation (dotted line) according to Wang, de Sturler, and Paulino [57]; cf. (6.3),
- RMINRES with deflation only (dash-dotted line), i.e., the residual norms of MINRES applied to the system  $\mathbf{P}_A \mathbf{A} \mathbf{x} = \mathbf{P}_A \mathbf{b}$ ; cf. (6.4),
- deflated MINRES (dashed line), i.e., the residual norms of MINRES applied to the system  $\mathbf{P}_A \mathbf{A} \mathbf{P}_A \mathbf{x} = \mathbf{P}_A \mathbf{Q}_A^H \mathbf{b}$ ; cf. section 6.2.

For the last three methods we used the matrix  $\mathbf{U} = [\mathbf{w}_1, \dots, \mathbf{w}_5, \mathbf{w}_{51}, \dots, \mathbf{w}_{55}]$  which contains the eigenvectors associated with the 10 eigenvalues of  $\mathbf{A}$  of smallest absolute value. Thus the deflation space  $\mathcal{U}$  has dimension 10. We have shown above that the two implementations of RMINRES and the deflated MINRES method are mathematically equivalent, and in this example the three convergence curves corresponding to these methods indeed coincide; see Figure 6.1.

*Example 6.4.* We now investigate breakdowns and near-breakdowns of the RMINRES method using a set of artificially constructed examples. Of course, the occurrence of an exact breakdown as in the following examples will be rare in practical applications.

For our construction we use the same matrix  $\mathbf{A}$  as in Example 6.3, and we construct a subspace  $\mathcal{U}$  for which  $\mathcal{U} \cap (\mathbf{A}\mathcal{U})^\perp \neq \{\mathbf{0}\}$ . Thus, the condition that guarantees a breakdown-free RMINRES computation is violated; cf. section 6.1. To construct the subspace  $\mathcal{U}$  we choose an integer  $k$ ,  $0 < k < m$ , and we define  $\mathbf{W}_1 = [\mathbf{w}_{i_1}, \dots, \mathbf{w}_{i_k}]$  and  $\mathbf{W}_2 = [\mathbf{w}_{m+i_1}, \dots, \mathbf{w}_{m+i_k}]$  for indices  $0 < i_1 < \dots < i_k < m$ . With  $\mathbf{D}_U = \text{diag}(\lambda_{i_1}, \dots, \lambda_{i_k})$  we obtain  $\mathbf{A}\mathbf{W}_1 = \mathbf{W}_1\mathbf{D}_U$  and  $\mathbf{A}\mathbf{W}_2 = -\mathbf{W}_2\mathbf{D}_U$  because of the symmetry of the spectrum of  $\mathbf{A}$ . We now choose the matrix  $\mathbf{U} = \mathbf{W}_1 + \mathbf{W}_2$ . Applying  $\mathbf{A}$  yields  $\mathbf{A}\mathbf{U} = (\mathbf{W}_1 - \mathbf{W}_2)\mathbf{D}_U$  and using the fact that  $\mathbf{W}$  is unitary shows that  $\mathbf{U}^H \mathbf{A}\mathbf{U} = \mathbf{0}$ , or, equivalently,  $\mathcal{U} \subset (\mathbf{A}\mathcal{U})^\perp$ . The proof of Theorem 5.1 gives us a way to construct an initial guess which leads to an immediate breakdown of RMINRES. For an arbitrary  $\mathbf{0} \neq \mathbf{u} \in \mathcal{U}$  we choose  $\mathbf{x}_0 = \mathbf{A}^{-1}(\mathbf{b} - \mathbf{u})$ . Because of  $\mathcal{U} \perp \mathbf{A}\mathcal{U}$  we have  $\mathbf{P}_A \mathbf{u} = \mathbf{u}$  and the initial residual of RMINRES is  $\mathbf{r}_0 = \mathbf{P}_A \mathbf{b} - \mathbf{P}_A \mathbf{A} \mathbf{x}_0 = \mathbf{u}$ . The breakdown then occurs in the first iteration because  $\mathbf{P}_A \mathbf{A} \mathbf{r}_0 = \mathbf{P}_A \mathbf{A} \mathbf{u} = \mathbf{0}$  since

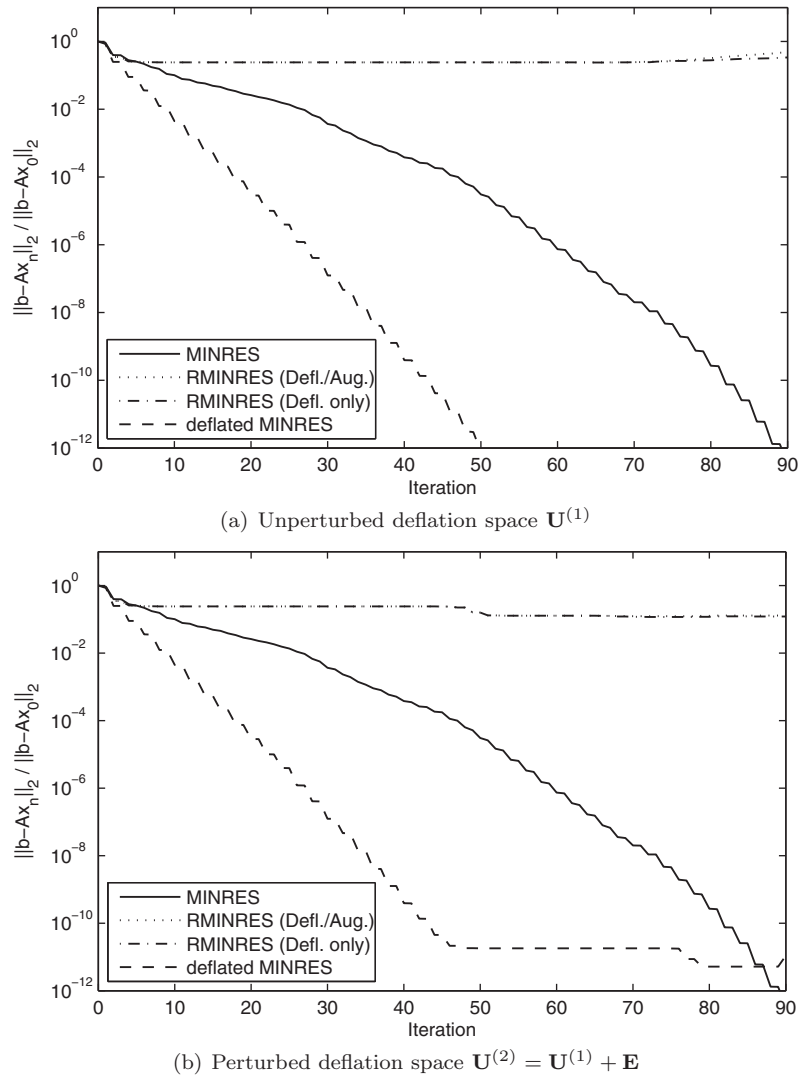


FIG. 6.2. Convergence history for Example 6.4. The convergence curves of both RMINRES solver implementations coincide.

$\mathbf{A}\mathbf{u} \in \mathbf{AU} = \ker(\mathbf{P}\mathbf{A})$ . For these constructed initial guesses the RMINRES method indeed breaks down immediately in numerical experiments, whereas the deflated MINRES method finds the solution after one step. There is no need to plot these results.

Of greater interest are situations with perturbed data. Interestingly, randomly perturbed initial guesses lead to a breakdown of RMINRES with the previously constructed deflation space as well. In Figure 6.2(a) we show the relative residual norms of the solvers listed above applied to the same  $\mathbf{A}$  and  $\mathbf{b}$  as in the previous example and with the matrix  $\mathbf{U}^{(1)} = [\mathbf{w}_1 + \mathbf{w}_{51}, \dots, \mathbf{w}_{10} + \mathbf{w}_{60}]$ .

Furthermore, breakdowns also occur when we perturb the deflation space. Figure 6.2(b) shows the results for a perturbed matrix  $\mathbf{U}^{(2)} = \mathbf{U}^{(1)} + \mathbf{E}$  with a random  $\mathbf{E} \in \mathbb{C}^{100 \times 10}$  and  $\|\mathbf{E}\|_2 = 10^{-10}$ . The used initial guess is the same perturbed initial guess as in the experiment conducted for Figure 6.2(a).

Note that both RMINRES implementations suffer from a breakdown after a few steps with both matrices  $\mathbf{U}^{(1)}$  and  $\mathbf{U}^{(2)}$ . With the unperturbed matrix  $\mathbf{U}^{(1)}$  the deflated MINRES method converges to the solution with a relative residual smaller than  $10^{-12}$ , while in the case of the perturbed matrix  $\mathbf{U}^{(2)}$  the method stagnates with a relative residual of order  $10^{-11}$ . This stagnation of deflated MINRES seems to be related to an unfavorable spectrum of  $\mathbf{P}_A \mathbf{A} \mathbf{P}_A$  for these specifically constructed and perturbed matrices like  $\mathbf{U}^{(2)}$ . It is unlikely that the stagnation is caused by roundoff errors because the stagnation also occurs (up to iteration 100) when full recurrences (GMRES) are used instead of short recurrences (MINRES). Perturbing the matrix  $\mathbf{U}$  from Example 6.3, whose columns are exact eigenvectors of  $\mathbf{A}$ , does not cause stagnation. This behavior is still subject to further research.

Note that the construction of  $\mathcal{U} = \text{im}(\mathbf{U})$  in Example 6.4 is such that  $\mathbf{A}\mathcal{U} \perp \mathcal{U}$  which cannot be achieved with a (nearly)  $\mathbf{A}$ -invariant subspace if  $\mathbf{A}$  is Hermitian. In [57] an approximation to an invariant subspace of a previous matrix in a sequence of linear algebraic systems is used. In this situation care has to be taken that the extracted space is still a good approximation to an invariant subspace of the current matrix. However, in the experiments of [57, section 7] this seems to be fulfilled since stagnation has not been observed.

**7. Conclusions.** In this paper we first analyzed theoretically the link between basic theoretical properties of deflated and augmented Krylov subspace methods whose residuals satisfy a Galerkin condition, including the minimum residual methods whose inclusion into the class of Galerkin methods requires a replacement of the standard inner product. We proved that augmentation can be achieved without explicitly augmenting the Krylov subspace but instead projecting the residuals appropriately and using a correction formula for the approximate solutions. We discussed this result in detail for the CG method and GMRES/MINRES methods, the main representatives of our class. It turned out that for these methods some of our results had been mentioned before in the literature.

The projections which arise from the augmentation can also be used to obtain a deflated system. We have seen that a left-multiplication of the original system with the corresponding projection yields a deflated system for which the CG method and GMRES/MINRES methods implicitly achieve augmentation. We proved that for non-singular Hermitian matrices the MINRES method for the deflated system is equivalent to the solver part of the RMINRES method introduced in [57]. While CG never breaks down, GMRES, MINRES, and thus RMINRES may suffer from breakdowns when used with the deflated systems. We stated necessary and sufficient conditions to characterize breakdowns of these minimal residual methods. For Hermitian matrices, we introduced the deflated MINRES method which also uses a Hermitian deflated matrix and proved that it cannot break down. These results were illustrated numerically.

Our framework covers methods based on a specific type of Galerkin condition; see (2.1)–(2.2). It does not include methods based on other conditions, in particular those that in practical methods are realized using the non-Hermitian Lanczos algorithms. Examples for such methods are BICG [21] and its variants including CGS [50], BI-CGSTAB [55], and IDR(s) [51]. Extending our framework to such methods remains a subject of further work.

Moreover, in this paper we did not discuss or recommend practical choices of deflation or augmentation spaces. Finding spaces that lead to an improved convergence behavior of the deflated or augmented method is a highly challenging task that should be attacked with a specific application in mind. Similar to preconditioning,

there exists no single best strategy for choosing deflation or augmentation spaces in practice. Often one deflates (approximations of) eigenvectors corresponding to the smallest eigenvalues of the given matrix. For symmetric or Hermitian positive definite matrices, this strategy can be shown to reduce the effective condition number, which in turn leads to improved convergence bounds and actually faster convergence of the iterative solver; see, e.g., [56]. For nonsymmetric or non-Hermitian matrices, however, the question of effective choices of deflation or augmentation spaces is largely open.

**Acknowledgments.** The authors wish to thank the anonymous referees for their comments which helped to improve the presentation.

## REFERENCES

- [1] J. BAGLAMA, D. CALVETTI, G. H. GOLUB, AND L. REICHEL, *Adaptively preconditioned GMRES algorithms*, SIAM J. Sci. Comput., 20 (1998), pp. 243–269.
- [2] A. H. BAKER, E. R. JESSUP, AND T. MANTEUFFEL, *A technique for accelerating the convergence of restarted GMRES*, SIAM J. Matrix Anal. Appl., 26 (2005), pp. 962–984.
- [3] M. BENZI, *Preconditioning techniques for large linear systems: A survey*, J. Comput. Phys., 182 (2002), pp. 418–477.
- [4] P. N. BROWN, *A theoretical comparison of the Arnoldi and GMRES algorithms*, SIAM J. Sci. Statist. Comput., 12 (1991), pp. 58–78.
- [5] P. N. BROWN AND H. F. WALKER, *GMRES on (nearly) singular systems*, SIAM J. Matrix Anal. Appl., 18 (1997), pp. 37–51.
- [6] A. CHAPMAN AND Y. SAAD, *Deflated and augmented Krylov subspace techniques*, Numer. Linear Algebra Appl., 4 (1997), pp. 43–66.
- [7] E. J. CRAIG, *The  $N$ -step iteration procedures*, J. Math. Phys., 34 (1955), pp. 64–73.
- [8] J. CULLUM, *Peaks, plateaus, numerical instabilities in a Galerkin/minimal residual pair of methods for solving  $Ax = b$* , Appl. Numer. Math., 19 (1995), pp. 255–278.
- [9] 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.
- [10] E. DE STURLER, *Nested Krylov methods based on GCR*, J. Comput. Appl. Math., 67 (1996), pp. 15–41.
- [11] E. DE STURLER, *Truncation strategies for optimal Krylov subspace methods*, SIAM J. Numer. Anal., 36 (1999), pp. 864–889.
- [12] Z. DOSTÁL, *Conjugate gradient method with preconditioning by projector*, Int. J. Comput. Math., 23 (1988), pp. 315–323.
- [13] M. EIERMANN AND O. ERNST, *Geometric aspects in the theory of Krylov space methods*, Acta Numer., 10 (2001), pp. 251–312.
- [14] M. EIERMANN, O. G. ERNST, AND O. SCHNEIDER, *Analysis of acceleration strategies for restarted minimal residual methods*, J. Comput. Appl. Math., 123 (2000), pp. 261–292.
- [15] S. C. EISENSTAT, H. C. ELMAN, AND M. H. SCHULTZ, *Variational iterative methods for non-symmetric systems of linear equations*, SIAM J. Numer. Anal., 20 (1983), pp. 345–357.
- [16] H. C. ELMAN, *Iterative Methods for Large, Sparse, Nonsymmetric Systems of Linear Equations*, Ph.D. thesis, Yale University, New Haven, CT, 1982.
- [17] J. ERHEL, K. BURRAGE, AND B. POHL, *Restarted GMRES preconditioned by deflation*, J. Comput. Appl. Math., 69 (1996), pp. 303–318.
- [18] J. ERHEL AND F. GUYOMARC'H, *An augmented conjugate gradient method for solving consecutive symmetric positive definite linear systems*, SIAM J. Matrix Anal. Appl., 21 (2000), pp. 1279–1299.
- [19] Y. A. ERLANGGA AND R. NABBEN, *Deflation and balancing preconditioners for Krylov subspace methods applied to nonsymmetric matrices*, SIAM J. Matrix Anal. Appl., 30 (2008), pp. 684–699.
- [20] Y. A. ERLANGGA AND R. NABBEN, *Multilevel projection-based nested Krylov iteration for boundary value problems*, SIAM J. Sci. Comput., 30 (2008), pp. 1572–1595.
- [21] R. FLETCHER, *Conjugate gradient methods for indefinite systems*, in Proceedings of the 6th Biennial Dundee Conference on Numerical Analysis, University of Dundee, Dundee, 1975, Lecture Notes in Math. 506, Springer, Berlin, 1976, pp. 73–89.
- [22] J. FRANK AND C. VUIK, *On the construction of deflation-based preconditioners*, SIAM J. Sci. Comput., 23 (2001), pp. 442–462.

- [23] R. W. FREUND, *On conjugate gradient type methods and polynomial preconditioners for a class of complex non-Hermitian matrices*, Numer. Math., 57 (1990), pp. 285–312.
- [24] V. M. FRIDMAN, *The method of minimum iterations with minimum errors for a system of linear algebraic equations with a symmetrical matrix*, USSR Comput. Math. Math. Phys., 2 (1963), pp. 362–363.
- [25] L. GIRAUD, S. GRATTON, X. PINEL, AND X. VASSEUR, *Flexible GMRES with deflated restarting*, SIAM J. Sci. Comput., 32 (2010), pp. 1858–1878.
- [26] A. GREENBAUM, *Iterative Methods for Solving Linear Systems*, Frontiers in Applied Mathematics 17, SIAM, Philadelphia, 1997.
- [27] M. H. GUTKNECHT AND M. ROZLOŽNÍK, *By how much can residual minimization accelerate the convergence of orthogonal residual methods?*, Numer. Algorithms, 27 (2001), pp. 189–213.
- [28] M. H. GUTKNECHT AND M. ROZLOŽNÍK, *A framework for generalized conjugate gradient methods—with special emphasis on contributions by Rüdiger Weiss*, Appl. Numer. Math., 41 (2002), pp. 7–22.
- [29] M. R. HESTENES AND E. STIEFEL, *Methods of conjugate gradients for solving linear systems*, J. Research Nat. Bur. Standards, 49 (1952), pp. 409–436.
- [30] E. F. KAASSCHIETER, *Preconditioned conjugate gradients for solving singular systems*, J. Comput. Appl. Math., 24 (1988), pp. 265–275.
- [31] S. A. KHARCHENKO AND A. Y. YEREMIN, *Eigenvalue translation based preconditioners for the GMRES( $k$ ) method*, Numer. Linear Algebra Appl., 2 (1995), pp. 51–77.
- [32] M. E. KILMER AND E. DE STURLER, *Recycling subspace information for diffuse optical tomography*, SIAM J. Sci. Comput., 27 (2006), pp. 2140–2166.
- [33] L. Y. KOLOTLINA, *Twofold deflation preconditioning of linear algebraic systems. I. Theory*, J. Math. Sci., 89 (1998), pp. 1652–1689.
- [34] L. MANSFIELD, *On the conjugate gradient solution of the Schur complement system obtained from domain decomposition*, SIAM J. Numer. Anal., 27 (1990), pp. 1612–1620.
- [35] L. MANSFIELD, *Damped Jacobi preconditioning and coarse grid deflation for conjugate gradient iteration on parallel computers*, SIAM J. Sci. Statist. Comput., 12 (1991), pp. 1314–1323.
- [36] R. B. MORGAN, *A restarted GMRES method augmented with eigenvectors*, SIAM J. Matrix Anal. Appl., 16 (1995), pp. 1154–1171.
- [37] R. B. MORGAN, *Restarted block-GMRES with deflation of eigenvalues*, Appl. Numer. Math., 54 (2005), pp. 222–236.
- [38] R. NABBEN AND C. VUIK, *A comparison of deflation and coarse grid correction applied to porous media flow*, SIAM J. Numer. Anal., 42 (2004), pp. 1631–1647.
- [39] R. NABBEN AND C. VUIK, *A comparison of deflation and the balancing preconditioner*, SIAM J. Sci. Comput., 27 (2006), pp. 1742–1759.
- [40] R. NABBEN AND C. VUIK, *A comparison of abstract versions of deflation, balancing and additive coarse grid correction preconditioners*, Numer. Linear Algebra Appl., 15 (2008), pp. 355–372.
- [41] R. A. NICOLAIDES, *Deflation of conjugate gradients with applications to boundary value problems*, SIAM J. Numer. Anal., 24 (1987), pp. 355–365.
- [42] M. A. OLSHANSKII AND V. SIMONCINI, *Acquired clustering properties and solution of certain saddle point systems*, SIAM J. Matrix Anal. Appl., 31 (2010), pp. 2754–2768.
- [43] C. C. PAIGE AND M. A. SAUNDERS, *Solutions of sparse indefinite systems of linear equations*, SIAM J. Numer. Anal., 12 (1975), pp. 617–629.
- [44] Y. SAAD, *Krylov subspace methods for solving large unsymmetric systems*, Math. Comp., 37 (1981), pp. 105–126.
- [45] Y. SAAD, *Analysis of augmented Krylov subspace methods*, SIAM J. Matrix Anal. Appl., 18 (1997), pp. 435–449.
- [46] Y. SAAD, *Iterative Methods for Sparse Linear Systems*, 2nd ed., SIAM, Philadelphia, 2003.
- [47] Y. SAAD AND M. H. SCHULTZ, *GMRES: A generalized minimal residual algorithm for solving nonsymmetric linear systems*, SIAM J. Sci. Statist. Comput., 7 (1986), pp. 856–869.
- [48] Y. SAAD, M. YEUNG, J. ERHEL, AND F. GUYOMARCH, *A deflated version of the conjugate gradient algorithm*, SIAM J. Sci. Comput., 21 (2000), pp. 1909–1926.
- [49] V. SIMONCINI AND D. B. SZYLD, *Recent computational developments in Krylov subspace methods for linear systems*, Numer. Linear Algebra Appl., 14 (2007), pp. 1–59.
- [50] P. SONNEVELD, *CGS, a fast Lanczos-type solver for nonsymmetric linear systems*, SIAM J. Sci. Statist. Comput., 10 (1989), pp. 36–52.
- [51] P. SONNEVELD AND M. B. VAN GIJZEN, *IDR( $s$ ): A family of simple and fast algorithms for solving large nonsymmetric systems of linear equations*, SIAM J. Sci. Comput., 31 (2008), pp. 1035–1062.
- [52] E. STIEFEL, *Relaxationsmethoden bester Strategie zur Lösung linearer Gleichungssysteme*, Comm. Math. Helv., 29 (1955), pp. 157–179.

- [53] J. M. TANG, S. P. MACLACHLAN, R. NABBEN, AND C. VUIK, *A comparison of two-level preconditioners based on multigrid and deflation*, SIAM J. Matrix Anal. Appl., 31 (2010), pp. 1715–1739.
- [54] J. M. TANG, R. NABBEN, C. VUIK, AND Y. A. ERLANGGA, *Comparison of two-level preconditioners derived from deflation, domain decomposition and multigrid methods*, J. Sci. Comput., 39 (2009), pp. 340–370.
- [55] 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.
- [56] C. VUIK, R. NABBEN, AND J. TANG, *Deflation acceleration for domain decomposition preconditioners*, in Proceedings of the 8th European Multigrid Conference, TU Delft, Delft, The Netherlands, P. Wesseling, C. Oosterlee, and P. Hemker, eds., 2006.
- [57] S. WANG, E. DE 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.
- [58] R. WEISS, *Convergence Behavior of Generalized Conjugate Gradient Methods*, Ph.D. thesis, University of Karlsruhe, Karlsruhe, Germany, 1990.
- [59] R. WEISS, *Error-minimizing Krylov subspace methods*, SIAM J. Sci. Comput., 15 (1994), pp. 511–527.
- [60] R. WEISS, *Properties of generalized conjugate gradient methods*, Numer. Linear Algebra Appl., 1 (1994), pp. 45–63.
- [61] M. YEUNG, J. TANG, AND C. VUIK, *On the convergence of GMRES with invariant-subspace deflation*, Report 10-14, Delft University of Technology, Delft, The Netherlands, 2010.