

## A COMPARISON OF DEFLATION AND THE BALANCING PRECONDITIONER\*

R. NABBEN<sup>†</sup> AND C. VUIK<sup>‡</sup>

**Abstract.** In this paper we compare various preconditioners for the numerical solution of partial differential equations. We compare the well-known balancing preconditioner used in domain decomposition methods with a so-called deflation preconditioner. We prove that the effective condition number of the deflated preconditioned system is always, i.e., for all deflation vectors and all restrictions and prolongations, below the condition number of the system preconditioned by the balancing preconditioner. Even more, we establish that both preconditioners lead to almost the same spectra. The zero eigenvalues of the deflation preconditioned system are replaced by eigenvalues which are one if the balancing preconditioner is used. Moreover, we prove that the A-norm of the errors of the iterates built by the deflation preconditioner is always below the A-norm of the errors of the iterates built by the balancing preconditioner. Depending on the implementation of the balancing preconditioner the amount of work of one iteration of the deflation preconditioned system is less than or equal to the amount of work of one iteration of the balancing preconditioned system. If the amount of work is equal, both preconditioners are sensitive with respect to inexact computations. Finally, we establish that the deflation preconditioner and the balancing preconditioner produce the same iterates if one uses certain starting vectors. Numerical results for porous media flows emphasize the theoretical results.

**Key words.** deflation, coarse grid correction, balancing, preconditioners, conjugate gradients, porous media flow, scalable parallel preconditioner

**AMS subject classifications.** 65F10, 65F50, 65N22

**DOI.** 10.1137/040608246

**1. Introduction.** The conjugate gradient method is the most-used method to solve large linear systems of equations

$$Ax = b$$

whose coefficient matrices  $A$  are sparse and symmetric positive definite. Such systems are encountered, for example, when a finite volume/difference/element method is used to discretize an elliptic partial differential equation.

The convergence rate of the conjugate gradient method (CG-method) is bounded as a function of the condition number of the system matrix to which it is applied. If the condition number of  $A$  is large, it is advisable to solve, instead, a preconditioned system  $M^{-1}Ax = M^{-1}b$ , where the symmetric positive definite preconditioner  $M$  is chosen such that  $M^{-1}A$  has a more clustered spectrum or a smaller condition number than that of  $A$ . Furthermore, system  $Mz = r$  must be cheap to solve relative to the improvement it provides in convergence rate. A final desirable property in a preconditioner is that it should parallelize well, especially on distributed memory computers.

---

\*Received by the editors May 14, 2004; accepted for publication (in revised form) September 6, 2005; published electronically February 3, 2006.

<http://www.siam.org/journals/sisc/27-5/60824.html>

<sup>†</sup>TU Berlin Institut für Mathematik, MA 3-3, Strasse des 17. Juni 136, D-10623 Berlin, Germany (nabben@math.tu-berlin.de).

<sup>‡</sup>Delft University of Technology, Faculty of Electrical Engineering, Mathematics and Computer Science, Delft Institute of Applied Mathematics, P.O. Box 5031, 2600 GA Delft, The Netherlands (c.vuik@math.tudelft.nl).

In [16] two different preconditioners are compared, namely a deflation preconditioner and an additive coarse grid correction preconditioner. It is shown that the deflation preconditioner leads to a smaller condition number than a coarse grid correction preconditioner like the BPS preconditioner by Bramble, Pasciak, and Schatz [1]. Here we compare the deflation preconditioner with the balancing preconditioner proposed by Mandel [11]. In the following we give a brief introduction into both preconditioning techniques.

To describe the deflation method we define the projection  $P_D$  by

$$(1.1) \quad P_D = I - AZ(Z^T AZ)^{-1}Z^T, \quad Z \in \mathbb{R}^{n \times r},$$

where the column space of  $Z$  is the deflation subspace, i.e., the space to be projected out of the residual, and  $I$  is the identity matrix of appropriate size. We assume that  $r \ll n$  and that  $Z$  has rank  $r$ . Under this assumption  $E \equiv Z^T AZ$  may be easily computed and factored and is symmetric positive definite. Since  $x = (I - P_D^T)x + P_D^T x$  and because

$$(1.2) \quad (I - P_D^T)x = Z(Z^T AZ)^{-1}Z^T Ax = ZE^{-1}Z^T b$$

can be immediately computed, we need only to compute  $P_D^T x$ . In light of the identity  $AP_D^T = P_D A$ , we can solve the deflated system

$$(1.3) \quad P_D A \tilde{x} = P_D b$$

for  $\tilde{x}$  using the CG-method, premultiply this by  $P_D^T$ , and add it to (1.2).

Obviously (1.3) is singular. But a positive semidefinite system can be solved by the CG-method as long as the right-hand side is consistent (i.e., as long as  $b = Ax$  for some  $x$ ) [9]. This is certainly true for (1.3), where the same projection is applied to both sides of the nonsingular system. Since the null space never enters the iteration, the corresponding zero eigenvalues do not influence the convergence [9, 22]. Motivated by this fact, we define the *effective condition number* of a positive semidefinite matrix  $C \in \mathbb{R}^{n \times n}$  with  $r$  zero eigenvalues to be the ratio of its largest to smallest *positive* eigenvalues:

$$\kappa_{eff}(C) = \frac{\lambda_n}{\lambda_{r+1}}.$$

It is possible to combine both a standard preconditioning and preconditioning by deflation (for details see [7]). The convergence is then described by the effective condition number of  $M^{-1}P_D A$ . For more details about the deflation preconditioner see [17, 15, 4, 13, 10, 25, 2, 24, 26, 7, 23, 16].

We compare the preconditioned deflation operator with the balancing preconditioner proposed by Mandel [11] and Mandel and Brezina [12] and analyzed by Dryja and Widlund [3], Pavarino and Widlund [18], and Toselli and Widlund [21]. As the FETI algorithm [5, 6] the balancing Neumann-Neumann preconditioner is one of the domain decomposition methods that have been most carefully implemented and severely tested on the very largest existing parallel computer systems.

Applied to some specific symmetric positive definite problems the balancing Neumann-Neumann preconditioner leads to moderately growing condition numbers if the size of the systems increases [20]. Moreover, if an appropriate scaling is used, the condition numbers are independent of jumps in the coefficients in the matrices [20].

In our notation the balancing preconditioner is given by

$$(1.4) \quad P_B = (I - ZE^{-1}Z^T A)M^{-1}(I - AZE^{-1}Z^T) + ZE^{-1}Z^T,$$

where  $Z \in \mathbb{R}^{n \times r}$ ,  $E = Z^T A Z$ , and  $M$  is a symmetric positive definite matrix. Note that  $P_B$  is symmetric and positive definite. For more details we refer the reader to [11] and the books [20, 19, 21].

As a first comparison of both preconditioners we observe that the balancing preconditioner needs per iteration 3 matrix vector products and the coarse grid operator is used 2 times. This makes the balancing preconditioner per iteration more expensive than the deflation approach. However, if an optimal implementation of the balancing preconditioner is used (see, e.g., [21]), the amount of work per iteration is the same.

In this article we give a detailed comparison of these two preconditioners. We prove that the effective condition number of the deflated preconditioned system  $M^{-1}P_D A$  is always below the condition number of the system preconditioned by the balancing preconditioner  $P_B A$ . Even more, we establish that the spectrum of  $P_B A$  is the same as  $M^{-1}P_D A$ , except the  $r$  zero eigenvalues are replaced by eigenvalues which are one.

This implies that for all matrices  $Z \in \mathbb{R}^{n \times r}$  and all positive definite preconditioners  $M^{-1}$  the effective condition number of the deflated preconditioned system is below or equal to the condition number of the system preconditioned by the balancing preconditioner! However, the condition number is not the only parameter which influences the convergence behavior of the CG-method. The convergence may be significantly faster if the eigenvalues of  $A$  are clustered [22]. But we obtain from the above-mentioned result that the clustering of the eigenvalues of the two different preconditioned systems is the same. However, we have a cluster at zero in one case and at one in the other case. These results are stated in section 2.

There are other properties which influence the convergence behavior of the CG-method, e.g., the starting vector, the right-hand side, and the location of the clusters of eigenvalues. Therefore, a more detailed comparison is given in section 3. There we prove that the A-norm of the errors of the iterates built by the deflation preconditioner is always below the A-norm of the errors of the iterates built by the balancing preconditioner. Moreover, we establish that the deflation preconditioner and the balancing preconditioner produce the same iterates if one uses certain starting vectors. More precisely we show which terms in the preconditioned CG-method are the same for both methods and which terms are different. At the end of section 3 we prove that the condition of the balancing preconditioned system decreases if one takes a finer grid as a coarse grid.

In section 4 numerical results emphasize our theoretical results.

**2. Spectral properties.** In this section we compare the effective condition number for the deflation and balancing preconditioned matrices. In section 2.1 we give some definitions and preliminary results. Thereafter a comparison is made if the projection vectors are equal to eigenvectors (in section 2.2) and for general projection vectors (in section 2.3).

**2.1. Notations and preliminary results.** In the following we denote by  $\lambda_i(M)$  the eigenvalues of a matrix  $M$ . If the eigenvalues are real, the  $\lambda_i(M)$ 's are ordered increasingly.

For two Hermitian  $n \times n$  matrices  $A$  and  $B$  we write  $A \succeq B$ , if  $A - B$  is positive semidefinite.

Next we mention well-known properties of the eigenvalues of Hermitian matrices.

LEMMA 2.1. *Let  $A, B \in \mathbb{C}^{n \times n}$  be Hermitian. For each  $k = 1, 2, \dots, n$  we have*

$$\lambda_k(A) + \lambda_1(B) \leq \lambda_k(A + B) \leq \lambda_k(A) + \lambda_n(B).$$

From the above lemma we easily obtain the next lemma.

LEMMA 2.2. *If  $A, B \in \mathbb{C}^{n \times n}$  are positive semidefinite with  $A \succeq B$ , then  $\lambda_i(A) \geq \lambda_i(B)$ .*

Moreover, we will use the following lemma.

LEMMA 2.3. *Let  $A, B \in \mathbb{C}^{n \times n}$  be Hermitian and suppose that  $B$  has rank at most  $r$ . Then*

- $\lambda_k(A + B) \leq \lambda_{k+r}(A), \quad k = 1, 2, \dots, n - r,$
- $\lambda_k(A) \leq \lambda_{k+r}(A + B), \quad k = 1, 2, \dots, n - r.$

Lemma 2.1, Lemma 2.2, and Lemma 2.3 can be found, e.g., as Theorem 4.3.1, Corollary 7.7.4, and Theorem 4.3.6, respectively, in [8].

**2.2. Projection vectors chosen as eigenvectors.** In this section we compare the effective condition number of  $P_D A$  and  $P_B A$  if the projection vectors are equal to the eigenvectors of  $A$ .

DEFINITION 2.4. *Let  $\lambda_i$  be the eigenvalues of  $A$ . Choose the eigenvectors  $v_k$  of  $A$  such that  $v_k^T v_j = \delta_{kj}$ , and define  $Z = [v_1 \dots v_r]$ .*

THEOREM 2.5. *Using  $Z$  as given in Definition 2.4 and preconditioner  $M$  equal to the identity, the spectrum of  $P_B A$  is*

$$\text{spectrum}(P_B A) = \{1, \dots, 1, \lambda_{r+1}, \dots, \lambda_n\}.$$

*Proof.* For this choice of  $Z$  it appears that

$$(2.1) \quad E = Z^T A Z = \text{diag}(\lambda_1, \dots, \lambda_r).$$

We consider  $P_B A v_k$ . For  $k = 1, \dots, n$  we obtain

$$\begin{aligned} P_B A v_k &= \left( I - Z \text{diag} \left( \frac{1}{\lambda_1}, \dots, \frac{1}{\lambda_r} \right) Z^T A \right) \left( I - A Z \text{diag} \left( \frac{1}{\lambda_1}, \dots, \frac{1}{\lambda_r} \right) Z^T \right) \lambda_k v_k \\ &\quad + Z \text{diag} \left( \frac{1}{\lambda_1}, \dots, \frac{1}{\lambda_r} \right) Z^T \lambda_k v_k. \end{aligned}$$

Using the orthogonality properties of the eigenvectors one obtains

$$P_B A v_k = v_k, \text{ for } k = 1, \dots, r.$$

For  $k = r + 1, \dots, n$  the same orthogonality properties lead to

$$P_B A v_k = \lambda_k v_k, \text{ for } k = r + 1, \dots, n. \quad \square$$

In order to compare both approaches we note that

$$(2.2) \quad \kappa_{eff}(P_D A) = \frac{\lambda_n}{\lambda_{r+1}}$$

and

$$(2.3) \quad \kappa(P_B A) = \frac{\max\{1, \lambda_n\}}{\min\{1, \lambda_{r+1}\}}.$$

From (2.2) and (2.3) it follows that  $\kappa(P_B A) \geq \kappa_{eff}(P_D A)$ , so the convergence bound based on the effective condition number implies that deflated CG never converges slower than CG combined with the balancing preconditioner if both methods use the eigenvectors corresponding to the  $r$  smallest eigenvalues as projection vectors.

**2.3. Projection vectors chosen as general vectors.** In the previous section we showed that the deflation technique leads to a smaller effective condition number than the balancing preconditioner, if eigenvectors are used. However, computing the  $r$  smallest eigenvalues is mostly very expensive. Moreover, in multigrid methods and domain decomposition methods special interpolation and prolongation matrices are used to obtain grid independent convergence rates. So a comparison only for eigenvectors is not enough. But in this section we generalize the results of the last section. We prove that the effective condition number of the deflated preconditioned system is always, i.e., for all matrices  $Z \in \mathbb{R}^{n \times r}$  and all preconditioners  $M^{-1}$ , below the condition number of the system preconditioned by the balancing preconditioner. To do this we repeat some properties of the projection operator  $P_D$  used in the deflation method (see [7]). The operator  $P_D$  is defined as

$$(2.4) \quad P_D = I - AZE^{-1}Z^T, \text{ where } E = Z^T AZ.$$

Furthermore, the following identities hold:

$$P_D^2 = P_D, \quad P_D AZ = 0, \quad Z^T P_D = P_D^T Z = 0, \quad \text{and} \quad AP_D^T = P_D A.$$

We start with a result for the deflation preconditioner which helps to compare the deflation and the balancing preconditioner.

**PROPOSITION 2.6.** *Let  $A \in \mathbb{R}^{n \times n}$  be symmetric positive definite. Let  $Z \in \mathbb{R}^{n \times r}$  with  $\text{rank} Z = r$ . Then*

$$\sigma(P_D^T M^{-1} P_D A) = \sigma(M^{-1} P_D A).$$

*Proof.* Since  $AP_D^T = P_D A$  the following identities hold:

$$\sigma(P_D^T M^{-1} P_D A) = \sigma(M^{-1} P_D AP_D^T) = \sigma(M^{-1} P_D^2 A) = \sigma(M^{-1} P_D A). \quad \square$$

Using Proposition 2.6 we obtain the following theorem.

**THEOREM 2.7.** *Let  $A \in \mathbb{R}^{n \times n}$  be symmetric positive definite. Let  $Z \in \mathbb{R}^{n \times r}$  with  $\text{rank} Z = r$ . Then the preconditioner defined in (1.1) and (1.4) satisfies*

$$(2.5) \quad \lambda_n(M^{-1} P_D A) \leq \lambda_n(P_B A),$$

$$(2.6) \quad \lambda_{r+1}(M^{-1} P_D A) \geq \lambda_1(P_B A).$$

*Proof.* We can write  $P_B$  as

$$P_B = P_D^T M^{-1} P_D + ZE^{-1}Z^T.$$

Thus

$$A^{\frac{1}{2}} P_B A^{\frac{1}{2}} = A^{\frac{1}{2}} P_D^T M^{-1} P_D A^{\frac{1}{2}} + A^{\frac{1}{2}} ZE^{-1}Z^T A^{\frac{1}{2}}.$$

Since  $A^{\frac{1}{2}} ZE^{-1}Z^T A^{\frac{1}{2}}$  is a symmetric positive semidefinite matrix of rank  $r$ , we obtain with Lemma 2.2

$$\lambda_i(P_B A) = \lambda_i(A^{\frac{1}{2}} P_B A^{\frac{1}{2}}) \geq \lambda_i(A^{\frac{1}{2}} P_D^T M^{-1} P_D A^{\frac{1}{2}}) = \lambda_i(P_D^T M^{-1} P_D A).$$

Using Lemma 2.3 we get

$$\lambda_{r+1}(P_D^T M^{-1} P_D A) \geq \lambda_1(P_B A).$$

Using Proposition 2.6 we get the desired result.  $\square$

It follows from Theorem 2.7 that

$$\kappa(P_B A) \geq \kappa_{eff}(M^{-1}P_D A),$$

so the convergence bound based on the effective condition number implies that preconditioned deflated CG never converges slower than CG preconditioned by the balancing preconditioner.

It appears that the results given in Theorem 2.5 can be generalized to general projection vectors.

THEOREM 2.8. *Suppose that the spectrum of  $M^{-1}P_D A$  is given by*

$$spectrum(M^{-1}P_D A) = \{0, \dots, 0, \mu_{r+1}, \dots, \mu_n\}.$$

Then

$$spectrum(P_B A) = \{1, \dots, 1, \mu_{r+1}, \dots, \mu_n\}.$$

*Proof.* We know that  $M^{-1}P_D A Z = 0$ , so the eigenvectors corresponding to the zero eigenvalues of  $M^{-1}P_D A$  are equal to  $\{z_1, \dots, z_r\}$ . On the other hand, it is easy to check that

$$P_B A Z = (P_D^T M^{-1} P_D + Z E^{-1} Z^T) A Z = P_D^T M^{-1} P_D A Z + Z E^{-1} Z^T A Z = Z.$$

This implies that  $\{z_1, \dots, z_r\}$  are the eigenvectors corresponding to the eigenvalues of  $P_B A$ , which are equal to 1.

Now we consider the eigenvalue  $\mu_i$ , with  $r + 1 \leq i \leq n$ . Suppose  $v_i$  is the corresponding eigenvector of  $M^{-1}P_D A$ , and thus  $M^{-1}P_D A v_i = \mu_i v_i$ . Since

$$M^{-1}P_D A v_i = M^{-1}P_D^2 A v_i = M^{-1}P_D A P_D^T v_i = \mu_i v_i \neq 0,$$

the vector  $P_D^T v_i$  is nonzero. Using this vector, it follows that

$$(2.7) \quad P_B A (P_D^T v_i) = (P_D^T M^{-1} P_D + Z E^{-1} Z^T) A P_D^T v_i$$

$$(2.8) \quad = P_D^T M^{-1} P_D A P_D^T v_i + Z E^{-1} Z^T A P_D^T v_i$$

$$(2.9) \quad = P_D^T M^{-1} P_D^2 A v_i = P_D^T M^{-1} P_D A v_i = \mu_i P_D^T v_i.$$

So the vectors  $P_D^T v_i$  are eigenvectors of  $P_B A$  corresponding to the eigenvalues  $\mu_i$ .  $\square$

Thus both preconditioners lead to almost the same spectra with the same clustering. The zero eigenvalues of the deflation preconditioned system are replaced by eigenvalues which are one if the balancing preconditioner is used.

**3. Comparing the norm of the residuals.** In order to make a more detailed comparison of the deflation operator and the balancing preconditioner for general projection vectors we start to compare the vector spaces which contain the approximations of both methods. Using CG with  $P_B$  as preconditioner and start vector  $x_{0,B} = 0$ , it is well known that

$$x_{k,B} \in K^k\{P_B A, P_B b\},$$

where the Krylov subspace  $K^k\{P_B A, P_B b\} = span\{P_B b, P_B A P_B b, \dots, (P_B A)^{k-1} P_B b\}$ .

THEOREM 3.1. *The Krylov space used in the CG method with  $P_B$  as preconditioner and start vector  $x_{0,B} = 0$  has the following property:*

$$(3.1) \quad K^k\{P_B A, P_B b\} \subset span\{Z E^{-1} Z^T b, P_D^T M^{-1} P_D b, \dots, P_D^T (M^{-1} P_D A)^{k-1} M^{-1} P_D b\}.$$

*Proof.* To start the proof we first note that

$$P_B b = P_D^T M^{-1} P_D b + Z E^{-1} Z^T b.$$

Therefore, the property holds for  $k = 1$ . For  $k = 2$  we note that

$$P_B A P_B b = (P_D^T M^{-1} P_D + Z E^{-1} Z^T) A (P_D^T M^{-1} P_D + Z E^{-1} Z^T) b.$$

Writing out the various terms on the right-hand side one obtains

$$Z E^{-1} Z^T A Z E^{-1} Z^T b = Z E^{-1} Z^T b,$$

$$P_D^T M^{-1} P_D A P_D^T M^{-1} P_D b = P_D^T M^{-1} P_D P_D A M^{-1} P_D b = P_D^T M^{-1} P_D A M^{-1} P_D b,$$

where we have used that  $A P_D^T = P_D A$  and  $P_D^2 = P_D$ . Finally, the terms

$$P_D^T M^{-1} P_D A Z E^{-1} Z^T \text{ and } Z E^{-1} Z^T A P_D^T M^{-1} P_D$$

are both zero because they contain the combination  $P_D A Z = 0$  or  $Z^T A P_D^T = (P_D A Z)^T = 0$ . Repeating this argument for  $(P_B A)^i P_B b$  for  $i = 3, \dots, k - 1$  proves the theorem.  $\square$

With respect to the approximation using preconditioned CG combined with deflation, we note that  $x = (I - P_D^T)x + P_D^T x = Z E^{-1} Z^T b + P_D^T x$ . So after  $k$  iterations of preconditioned CG applied to  $A P_D^T x = P_D A x = P_D b$  we get the approximation  $\tilde{x}_{k,D}$ . The approximation  $x_{k,D}$  of the solution vector  $x$  is then given by  $x_{k,D} = Z E^{-1} Z^T b + P_D^T \tilde{x}_{k,D}$ . The vector  $x_{k,D}$  is contained in the following space:

$$x_{k,D} \in Z E^{-1} Z^T b + \text{span}\{P_D^T M^{-1} P_D b, \dots, P_D^T (M^{-1} P_D A)^{k-1} M^{-1} P_D b\}.$$

This implies that both approximations are elements of the same space. So the difference in quality of the approximation only depends on which norm is minimized.

LEMMA 3.2. *For the deflation iterates  $x_{k,D}$  and  $\tilde{x}_{k,D}$  with start vector  $\tilde{x}_{0,D} = 0$  the following optimality property holds:*

$$(3.2) \quad \|x - x_{k,D}\|_A = \|P_D^T(x - \tilde{x}_{k,D})\|_A = \min_{\xi \in K^k\{M^{-1}P_D A, M^{-1}P_D b\}} \|P_D^T(x - \xi)\|_A.$$

*Proof.* The first equality follows from the fact that  $x = (I - P_D^T)x + P_D^T x$ . If CG is applied to the preconditioned system

$$(3.3) \quad L^{-1} P_D A L^{-T} y = L^{-1} P_D b,$$

the following expression holds:

$$(3.4) \quad \|\tilde{y} - y_k\|_{L^{-1}P_D A L^{-T}} = \min_{\eta \in K^k\{L^{-1}P_D A L^{-T}, L^{-1}P_D b\}} \|\tilde{y} - \eta\|_{L^{-1}P_D A L^{-T}},$$

where  $\tilde{y}$  is a solution of (3.3). Note that  $\tilde{x} = L^{-T} \tilde{y}$  is a solution of  $P_D A x = P_D b$ . Rewriting (3.4) with  $\xi = L^{-T} \eta$  leads to

$$\|L^T(\tilde{x} - \tilde{x}_{k,D})\|_{L^{-1}P_D A L^{-T}} = \min_{\xi \in K^k\{M^{-1}P_D A, M^{-1}P_D b\}} \|L^T(\tilde{x} - \xi)\|_{L^{-1}P_D A L^{-T}}.$$

Using the equalities

$$\begin{aligned} \|L^T(\tilde{x} - \tilde{x}_{k,D})\|_{L^{-1}P_D A L^{-T}}^2 &= (\tilde{x} - \tilde{x}_{k,D})^T P_D A (\tilde{x} - \tilde{x}_{k,D}) \\ &= (\tilde{x} - \tilde{x}_{k,D})^T P_D^2 A (\tilde{x} - \tilde{x}_{k,D}) = \|P_D^T(\tilde{x} - \tilde{x}_{k,D})\|_A^2 = \|P_D^T(x - \tilde{x}_{k,D})\|_A^2 \end{aligned}$$

leads to the proof of the lemma.  $\square$

**THEOREM 3.3.** *Let  $x_{k,D}$  and  $\tilde{x}_{k,D}$  be the deflation iterates with start vector  $\tilde{x}_{0,D} = 0$ . For every  $x_k \in \text{span}\{ZE^{-1}Z^Tb, P_D^T M^{-1}P_D b, \dots, P_D^T (M^{-1}P_D A)^{k-1} M^{-1}P_D b\}$  the following inequality holds:*

$$\|x - x_{k,D}\|_A \leq \|x - x_k\|_A.$$

*Proof.* We decompose  $x_k$  as follows:

$$x_k = \alpha ZE^{-1}Z^Tb + P_D^T \xi, \text{ where } \xi \in K^k\{M^{-1}P_D A, M^{-1}P_D b\}.$$

Substituting this into  $\|x - x_k\|_A^2$  shows that

$$\|x - x_k\|_A^2 = \|x - \alpha ZE^{-1}Z^Tb - P_D^T \xi\|_A^2.$$

Using the equation  $x = (I - P_D^T)x + P_D^T x = ZE^{-1}Z^Tb + P_D^T x$  we obtain

$$\begin{aligned} \|x - x_k\|_A^2 &= \|(1 - \alpha)ZE^{-1}Z^Tb - P_D^T(x - \xi)\|_A^2 \\ &= (1 - \alpha)^2\|ZE^{-1}Z^Tb\|_A^2 + \|P_D^T(x - \xi)\|_A^2 \\ &\quad + (1 - \alpha)b^T ZE^{-1}Z^T A P_D^T(x - \xi) \\ &\quad + (1 - \alpha)(x - \xi)^T P_D A ZE^{-1}Z^T b. \end{aligned}$$

The last two terms are equal to zero, because  $Z^T A P_D^T = (P_D A Z)^T = 0$ . For  $x_{k,D}$  we know that  $\alpha = 1$ . This together with Lemma 3.2 implies

$$\|x - x_{k,D}\|_A^2 \leq (1 - \alpha)^2\|ZE^{-1}Z^Tb\|_A^2 + \|P_D^T(x - \xi)\|_A^2 = \|x - x_k\|_A^2,$$

where  $\xi \in K^k\{M^{-1}P_D A, M^{-1}P_D b\}$ .  $\square$

Theorems 3.1 and 3.3 imply the following theorem.

**THEOREM 3.4.** *The iterates  $x_{k,D}$  and  $x_{k,B}$  of the CG-method with start vector zero and preconditioned by the deflation preconditioner and the balancing preconditioner, respectively, satisfy*

$$\|x - x_{k,D}\|_A \leq \|x - x_{k,B}\|_A.$$

Next we are able to prove that using a certain start vector the iterates  $x_{k,D}$  are equal to the  $x_{k,B}$ .

**THEOREM 3.5.** *Using  $x_{0,B} = ZE^{-1}Z^Tb$  and  $\tilde{x}_{0,D} = 0$  it follows that  $x_{k,D} = x_{k,B}$ .*

*Proof.* Using the start vector  $x_{0,B} = ZE^{-1}Z^Tb$  it appears that

$$r_{0,B} = b - Ax_{0,B} = (I - AZE^{-1}Z^T)b = P_D b.$$

This implies that the Krylov subspace is given by  $K^k\{P_B A, P_B P_D b\}$ . For  $k = 1$  it follows from  $P_D^2 = P_D$  and  $Z^T P_D = Z^T(I - AZE^{-1}Z^T) = 0$  that

$$(P_D^T M^{-1}P_D + ZE^{-1}Z^T)P_D b = P_D^T M^{-1}P_D b.$$

For  $k = 2$  we know from the proof of Theorem 3.1 that

$$P_B A P_B P_D b = ZE^{-1}Z^T P_D b + P_D^T M^{-1}P_D A M^{-1}P_D^2 b.$$

Note that  $Z^T P_D = Z^T(I - AZE^{-1}Z^T) = 0$ , and thus

$$P_B A P_B P_D b = P_D^T M^{-1}P_D A M^{-1}P_D b.$$



Repeating this argument shows that

$$\begin{aligned} K^k\{P_B A, P_B P_D b\} &= \text{span}\{P_D^T M^{-1} P_D b, \dots, P_D^T (M^{-1} P_D A)^{k-1} M^{-1} P_D b\} \\ &= P_D^T K^k\{M^{-1} P_D A, M^{-1} P_D b\}. \end{aligned}$$

We again use the fact that CG combined with the balancing preconditioner minimizes

$$(x - x_{k,B})^T A(x - x_{k,B}),$$

where

$$x_{k,B} = ZE^{-1}Z^T b + P_D^T \xi, \text{ and } \xi \in K^k\{M^{-1}P_D A, M^{-1}P_D b\}$$

due to the choice of the start vector  $x_{0,B} = ZE^{-1}Z^T b$ . We have that

$$x - x_{k,B} = x - ZE^{-1}Z^T b - P_D^T \xi = P_D^T(x - \xi).$$

But by Lemma 3.2 the optimal  $\xi$  is nothing other than  $\tilde{x}_{k,D}$ . Thus we obtain

$$x - x_{k,B} = P_D^T(x - \tilde{x}_{k,D}).$$

Since  $x = ZE^{-1}Z^T b + P_D^T x$  we get

$$x_{k,D} = ZE^{-1}Z^T b + P_D^T \tilde{x}_{k,D} = x_{k,B}. \quad \square$$

Using the identity  $x_{k,D} = x_{k,B}$  it is easy to see that Theorem 2.11 of [16] implies that the balancing preconditioner with  $x_{0,B} = ZE^{-1}Z^T b$  never converges slower than the additive coarse grid preconditioner.

In the following we give a more detailed analysis of the preconditioned CG-method for both preconditioners if the above start vectors are used. We prove which quantities in the preconditioned CG-algorithm (PCG) are the same for both preconditioners and which are different. To make this paper self-contained we repeat the PCG-algorithm.

PCG-ALGORITHM for  $Ax = b$  with preconditioner  $M^{-1}$ .

$$r_0 := b - Ax_0, \quad z_0 = M^{-1}r_0, \quad p_0 := z_0$$

For  $j = 0, 1, \dots$  until convergence, do

$$\alpha_j := (r_j, z_j) / (Ap_j, p_j)$$

$$x_{j+1} := x_j + \alpha_j p_j$$

$$r_{j+1} := r_j - \alpha_j Ap_j$$

$$z_{j+1} := M^{-1}r_{j+1}$$

$$\beta_j := (r_{j+1}, z_{j+1}) / (r_j, z_j)$$

$$p_{j+1} := z_{j+1} + \beta_j p_j$$

end

Moreover, we need the next proposition.

PROPOSITION 3.6. *Let  $P_D, P_B,$  and  $M^{-1}$  be defined as above. Then*

$$(3.5) \quad P_D^T P_B P_D = P_D^T M^{-1} P_D = P_D^T P_B = P_B P_D.$$

*Proof.* Since  $P_D = I - AZE^{-1}Z^T$  we have  $P_D^T Z = Z - ZE^{-1}Z^T AZ = 0$ . Hence,

$$P_D^T P_B P_D = P_D^{T^2} M^{-1} P_D^2 + P_D^T Z E^{-1} Z P_D = P_D^T M^{-1} P_D.$$

Similarly,

$$P_D^T P_B = P_D^{T^2} M^{-1} P_D + P_D Z E^{-1} Z P_D = P_D^T M^{-1} P_D.$$

Since  $P_D^T M^{-1} P_D$  is symmetric we also have  $P_D^T M^{-1} P_D = P_B P_D$ .  $\square$

Now we can prove the following theorem.

**THEOREM 3.7.** *Using the PCG-algorithm with the balancing preconditioner  $P_B$  and  $x_{0,B} = ZE^{-1}Z^Tb$  on one side and with the deflation preconditioner  $M^{-1}P_D$  and  $\tilde{x}_{0,D} = 0$  on the other side we have, for all  $j$ ,*

$$\begin{aligned} (r_{j,D}, z_{j,D}) &= (r_{j,B}, z_{j,B}), \\ (P_D A p_{j,D}, p_{j,D}) &= (A p_{j,B}, p_{j,B}), \\ r_{j+1,D} &= r_{j+1,B}, \\ z_{j+1,B} &= P_D^T z_{j+1,D}, \\ p_{j+1,B} &= P_D^T p_{j+1,D}, \\ \beta_{j,B} &= \beta_{j,D}, \\ x_{j+1,B} &= x_{j+1,D} = ZE^{-1}Z^T + P_D^T \tilde{x}_{j+1,D}. \end{aligned}$$

*Proof.* If we use PCG for

$$P_D A x = P_D b$$

with preconditioner  $M^{-1}$  and start vector  $x_0 = 0$ , we obtain

$$\begin{aligned} x_{0,D} &= 0, \quad r_{0,D} = P_D b, \quad z_{0,D} = M^{-1} P_D b, \\ p_{0,D} &= z_{0,D} = M^{-1} P_D b, \quad \alpha_{0,D} = \frac{(r_{0,D}, z_{0,D})}{(P_D A p_{0,D}, p_{0,D})}, \\ \tilde{x}_{1,D} &= 0 + \alpha_{0,D} M^{-1} P_D b, \\ x_{1,D} &= ZE^{-1}Z^T b + \alpha_{0,D} P_D^T M^{-1} P_D b. \end{aligned}$$

If we use PCG for

$$A x = b$$

with preconditioner  $P_B$  and start vector  $x_0 = ZE^{-1}Z^Tb$ , we obtain

$$\begin{aligned} x_{0,B} &= ZE^{-1}Z^Tb, \quad r_{0,B} = P_D b, \quad z_{0,B} = P_B P_D b, \\ p_{0,B} &= z_{0,B} = P_B P_D b, \quad \alpha_{0,B} = \frac{(r_{0,B}, z_{0,B})}{(A p_{0,B}, p_{0,B})}, \\ x_{1,B} &= ZE^{-1}Z^Tb + \alpha_{0,B} P_B P_D b. \end{aligned}$$

Obviously, we have for all iterates

$$(3.6) \quad P_D r_{j+1,D} = P_D(P_D b - P_D A x_{j+1,D}) = r_{j+1,D}.$$

The identity

$$P_B P_D = P_D^T M^{-1} P_D^2 + ZE^{-1}Z^T P_D = P_D^T M^{-1} P_D$$

is frequently used in the following analysis.

Next, we prove the following identities by induction:

$$\begin{aligned} (r_{j,D}, z_{j,D}) &= (r_{j,B}, z_{j,B}), \quad (r_{j+1,D}, z_{j+1,D}) = (r_{j+1,B}, z_{j+1,B}), \\ (P_D A p_{j,D}, p_{j,D}) &= (A p_{j,B}, p_{j,B}), \end{aligned}$$

$$\begin{aligned}
r_{j+1,D} &= r_{j+1,B}, \\
z_{j+1,B} &= P_D^T z_{j+1,D}, \\
p_{j+1,B} &= P_D^T p_{j+1,D}, \\
\beta_{j,B} &= \beta_{j,D}, \\
x_{j+1,B} &= x_{j+1,D} = ZE^{-1}Z^T + P_D^T \tilde{x}_{j+1,D}.
\end{aligned}$$

In the following we use Proposition 3.6 and (3.6). For  $j = 0$  we have

$$(r_{0,D}, z_{0,D}) = b^T P_D^T M^{-1} P_D b = b^T P_D^T P_B P_D b = (r_{0,B}, z_{0,B}).$$

$$\begin{aligned}
(P_D A p_{0,D}, p_{0,D}) &= b^T P_D^T M^{-1} P_D A M^{-1} P_D b \\
&= b^T P_D^T M^{-1} P_D A P_D^T M^{-1} P_D b \\
&= b^T P_D^T P_B A P_B P_D b \\
&= (A p_{0,B}, p_{0,B}).
\end{aligned}$$

Hence,  $\alpha_{0,D} = \alpha_{0,B}$ .

$$\begin{aligned}
r_{1,D} &= P_D b - \alpha_{0,D} P_D A M^{-1} P_D b = P_D b - \alpha_{0,D} A P_D^T M^{-1} P_D b \\
&= P_D b - \alpha_{0,B} A P_B P_D b \\
&= r_{1,B}.
\end{aligned}$$

$$\begin{aligned}
x_{1,D} &= ZE^{-1}Z^T b + \alpha_{0,D} P_D^T M^{-1} P_D b \\
&= ZE^{-1}Z^T b + \alpha_{0,B} P_B P_D b = x_{1,B}.
\end{aligned}$$

$$\begin{aligned}
P_D^T z_{1,D} &= P_D^T M^{-1} r_{1,D} = P_D^T M^{-1} P_D r_{1,D} \\
&= P_B P_D r_{1,D} = P_B r_{1,B} = z_{1,B}.
\end{aligned}$$

Thus

$$\begin{aligned}
(r_{1,B}, z_{1,B}) &= (r_{1,D}, P_D^T z_{1,D}) \\
&= (P_D r_{1,D}, z_{1,D}) = (r_{1,D}, z_{1,D}).
\end{aligned}$$

Hence  $\beta_{0,D} = \beta_{0,B}$ . Next,

$$\begin{aligned}
p_{1,B} &= z_{1,B} + \beta_{0,B} p_{0,B} \\
&= P_D^T z_{1,D} + \beta_{0,B} P_B P_D b \\
&= P_D^T z_{1,D} + \beta_{0,B} P_D^T M^{-1} P_D b \\
&= P_D^T (z_{1,D} + \beta_{0,D} p_{0,D}) \\
&= P_D^T p_{1,D}.
\end{aligned}$$

Now assume that the above identities hold for  $j - 1$  and that  $(r_{j,B}, z_{j,B}) = (r_{j,D}, z_{j,D})$  holds. We then have

$$\begin{aligned}
(A p_{j,B}, p_{j,B}) &= (A P_D^T p_{j,D}, P_D^T p_{j,D}) \\
&= p_{j,D}^T P_D A P_D^T p_{j,D} \\
&= (P_D A p_{j,D}, p_{j,D}).
\end{aligned}$$

Hence,  $\alpha_{j,D} = \alpha_{j,B}$ . Since

$$\begin{aligned} x_{j+1,B} &= x_{j,B} + \alpha_{j,B} p_{j,B}, \\ \tilde{x}_{j+1,D} &= \tilde{x}_{j,D} + \alpha_{j,D} p_{j,D}, \end{aligned}$$

we obtain

$$\begin{aligned} x_{j+1,D} &= ZE^{-1}Z^T + P_D \tilde{x}_{j,D} + \alpha_{j,D} P_D^T p_{j,D} \\ &= x_{j,B} + \alpha_{j,B} p_{j,B} \\ &= x_{j+1,B}. \end{aligned}$$

$$\begin{aligned} r_{j+1,B} &= r_{j,B} - \alpha_{j,B} A p_{j,B} \\ &= P_D r_{j,D} - \alpha_{j,D} A P_D^T p_{j,D} \\ &= P_D r_{j,D} - \alpha_{j,D} P_D A p_{j,D} \\ &= P_D r_{j+1,D} \\ &= r_{j+1,D}. \end{aligned}$$

Moreover,

$$\begin{aligned} z_{j+1,B} &= P_B r_{j+1,B} = P_B P_D r_{j+1,D} \\ &= P_D^T M^{-1} P_D r_{j+1,D} = P_D^T M^{-1} r_{j+1,D} \\ &= P_D^T z_{j+1,D}. \end{aligned}$$

$$\begin{aligned} (r_{j+1,B}, z_{j+1,B}) &= (P_D r_{j+1,D}, P_D^T z_{j+1,D}) \\ &= r_{j+1,D}^T P_D^T z_{j+1,D} \\ &= (r_{j+1,D}, z_{j+1,D}). \end{aligned}$$

Hence,  $\beta_{j,B} = \beta_{j,D}$ . Next we have

$$\begin{aligned} p_{j+1,B} &= z_{j+1,B} + \beta_{j,B} p_{j,B} \\ &= P_D^T z_{j+1,D} + \beta_{j,D} P_D^T p_{j,D} \\ &= P_D^T p_{j+1,D}, \end{aligned}$$

which completes the proof.  $\square$

In the following we show how the eigenvalues and the condition number of the system preconditioned by balancing behave if we choose a finer coarse grid. Therefore, let  $Z_1 \in \mathbb{R}^{n \times r}$  and  $Z_2 \in \mathbb{R}^{n \times s}$  with  $\text{rank} Z_1 = r$  and  $\text{rank} Z_2 = s$ . Define

$$\begin{aligned} E_1 &:= Z_1^T A Z_1 \quad \text{and} \quad E_2 := Z_2^T A Z_2, \\ P_{D_1} &= I - A Z_1 E_1^{-1} Z_1^T \quad \text{and} \quad P_{D_2} = I - A Z_2 E_2^{-1} Z_2^T. \end{aligned}$$

Moreover, let

$$(3.7) \quad P_{B_1} = P_{D_1}^T M^{-1} P_{D_1} + Z_1 E_1^{-1} Z_1^T \quad \text{and} \quad P_{B_2} = P_{D_2}^T M^{-1} P_{D_2} + Z_2 E_2^{-1} Z_2^T.$$

We then have the following theorem.

**THEOREM 3.8.** *Let  $A$  and  $M \in \mathbb{R}^{n \times n}$  be symmetric positive definite. Let  $P_{B_1}$  and  $P_{B_2}$  be defined as in (3.7). If  $\text{Im} Z_1 \subseteq \text{Im} Z_2$ , then*

$$(3.8) \quad \lambda_n(P_{B_1} A) \geq \lambda_n(P_{B_2} A),$$

$$(3.9) \quad \lambda_{r+1}(P_{B_1} A) \leq \lambda_{s+1}(P_{B_2} A).$$

Moreover,

$$\text{cond}(P_{B_1}A) \geq \text{cond}(P_{B_2}A).$$

*Proof.* Theorem 2.12 in [16] states that

$$\begin{aligned} \lambda_n(M^{-1}P_{D_1}A) &\geq \lambda_n(M^{-1}P_{D_2}A), \\ \lambda_{r+1}(M^{-1}P_{D_1}A) &\leq \lambda_{s+1}(M^{-1}P_{D_2}A). \end{aligned}$$

Thus, with Theorem 2.8 we get

$$\begin{aligned} \text{cond}(P_{B_1}A) &= \frac{\max(1, \lambda_n(M^{-1}P_{D_1}A))}{\min(1, \lambda_{r+1}(M^{-1}P_{D_1}A))} \\ &\geq \frac{\max(1, \lambda_n(M^{-1}P_{D_2}A))}{\min(1, \lambda_{s+1}(M^{-1}P_{D_2}A))} \\ &= \text{cond}(P_{B_2}A). \quad \square \end{aligned}$$

If a finer grid is used as a coarse grid in the balancing preconditioner, the amount of work to solve the coarse grid system is increased. But then Theorem 3.8 states that the condition number of the system preconditioned by the balancing preconditioner decreases. In general this leads to fewer iterations, although more work is needed on the coarse grid problem.

**4. Numerical experiments.** In all our numerical experiments, the multiplication  $y = E^{-1}b$  is done by solving  $y$  from  $Ey = b$ , where  $E$  is decomposed in its Cholesky factor. The choice of the boundary conditions is such that all problems have as exact solution the vector with components equal to 1. In order to make the convergence behavior representative for general problems, we chose a random vector as the starting solution instead of the zero start vector.

**4.1. Artificial test problems.** We apply both methods (deflation and balancing) to the Poisson equation. It appears that in the numerical experiments  $\|x - x_{k,D}\|_A \leq \|x - x_{k,B}\|_A$ , but for well-scaled problems the differences are very small. From Theorem 2.8 it follows that the spectrum of the balancing preconditioner consists of two parts: in one part the eigenvalues are equal to 1, and in the other part the eigenvalues are equal to the nonzero eigenvalues of the deflated matrix. This suggests that if the eigenvalues equal to 1 are interior eigenvalues, the convergence is close to the convergence of the preconditioned deflation method; otherwise these eigenvalues may influence the convergence.

**Scaling properties.** Note that  $P_D A$  is scaling invariant, whereas  $P_B A$  is not scaling invariant. This means that if deflation is applied to a system  $\gamma A x = \gamma b$  the effective condition number of  $P_{D\gamma A} \gamma A = (I - \gamma A Z (Z^T \gamma A Z)^{-1} Z^T) \gamma A$  is independent of the scalar  $\gamma$ ; i.e.,

$$\kappa_{eff}(P_{D\gamma A} \gamma A) = \frac{\gamma \lambda_n(P_{D A A})}{\gamma \lambda_{r+1}(P_{D A A})} = \kappa_{eff}(P_{D A A}),$$

whereas the condition number of  $P_B \gamma A$  depends on the choice of  $\gamma$ ,

$$\kappa(P_{B\gamma A} \gamma A) \neq \kappa(P_{B A A}).$$

To check this in practice, we do experiments with balancing using various values of  $\gamma$ . From Figure 4.1 it appears that the convergence of the balancing preconditioner

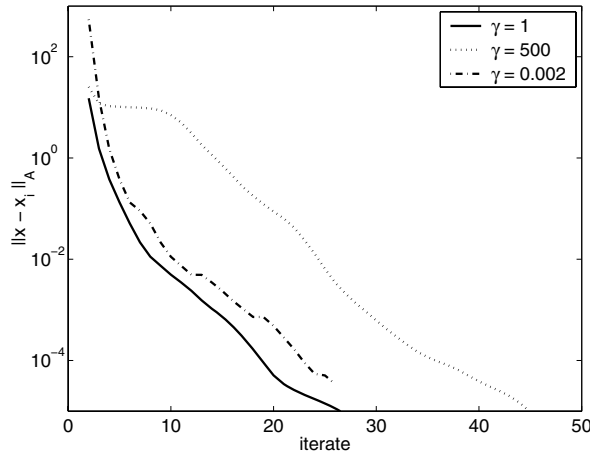


FIG. 4.1. Comparison of the balancing preconditioner for various values of  $\gamma$ . The convergence of deflation is identical to the convergence of the balancing preconditioner with  $\gamma = 1$ .

is worse if  $\gamma \neq 1$ . We note that the deflation method (for all values of  $\gamma$ ) has the same convergence as the balancing method with  $\gamma = 1$ . In most preconditioning techniques an automatic scaling is used, which makes the combination scaling invariant. However, the resulting spectrum is not always clustered around 1. We used this example to illustrate what happens if the spectrum is far away from 1.

**Inaccurate solution.** If the dimensions of matrix  $E$  become large (i.e., many projection vectors are used), it seems to be a good idea to compute  $E^{-1}$  approximately (by an iterative method or a recursive procedure). It appears that the balancing preconditioner is insensitive to the accuracy of the approximation of  $E^{-1}$ , while deflation is sensitive to it.

To illustrate this we consider the same Poisson problem. In both examples seven projection vectors are used. We replace  $E^{-1}$  by  $\tilde{E}^{-1} = (I + \epsilon R)E^{-1}(I + \epsilon R)$ , where  $R$  is a symmetric  $r \times r$  matrix with random elements chosen from the interval  $[-\frac{1}{2}, \frac{1}{2}]$ . From Figure 4.2 it follows that the convergence of the deflation preconditioner is good as long as  $|\epsilon| < 10^{-6}$ .

**Starting solution for the balancing preconditioner.** In Theorem 3.5 we have proven that  $x_{k,B} = x_{k,D}$  if  $x_{0,B} = ZE^{-1}Z^Tb$  and  $\tilde{x}_{0,D} = 0$ . In this paragraph we illustrate this by numerical examples. In Figure 4.3 we plot the convergence of the balancing preconditioner with start vector  $x_{0,B} = ZE^{-1}Z^Tb$ . It appears that the choice  $\gamma = 500$  leads to the same results as  $\gamma = 1$  (and deflation). Furthermore, the convergence for the choice  $\gamma = 0.002$  is initially also the same, but later on the convergence becomes worse. This can be explained by rounding errors. Using the choice  $\gamma = 0.002$  the eigenvalues equal to 1 are large with respect to the other eigenvalues. Initially, due to the start vector the components of the corresponding eigenvectors are zero or small. During the iterations, the perturbations in large eigencomponents increase, which leads to the same convergence as if the method is started with  $x_{0,B} = 0$ . To enlarge the rounding error effect we have also done experiments where the matrix  $E^{-1}$  is replaced by  $\tilde{E}^{-1}$  with  $\epsilon = 10^{-2}$ . The results are given in Figure 4.4. Note that the same effect now appears for both values of  $\gamma \neq 1$ .

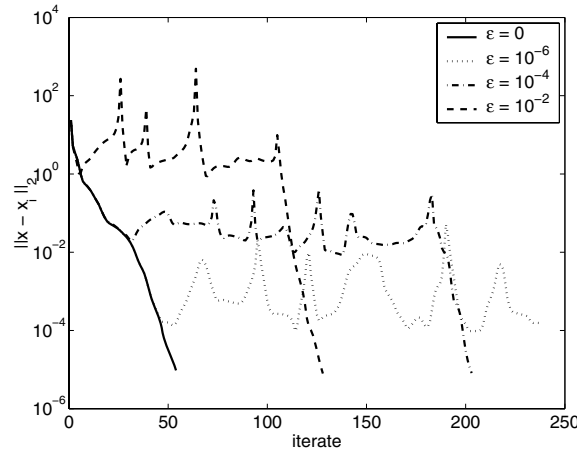


FIG. 4.2. Convergence behavior of the deflated ICCG including perturbations.

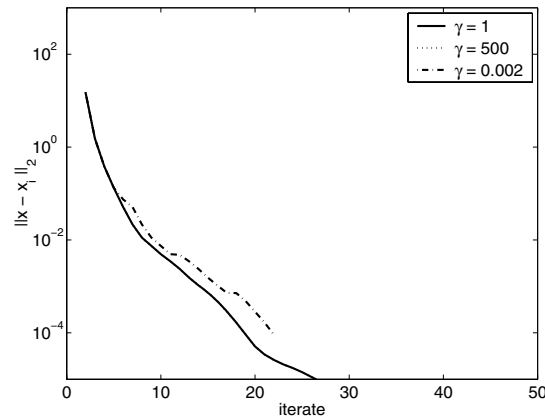


FIG. 4.3. Convergence of the balancing preconditioner with  $x_{0,B} = ZE^{-1}Z^Tb$  and  $\epsilon = 0$ . The convergence of deflation is identical to the convergence of the balancing preconditioner with  $\gamma = 1$ .

**4.2. Porous media flows.** In this section we simulate a porous media oil flow in a 3-dimensional layered geometry, where the layers vary in thickness and orientation (see Figures 4.5 and 4.6 for a 4-layer problem). The fluid pressure and permeability are denoted by  $p$  and  $\sigma$ , respectively. The pressure  $p$  satisfies the equation

$$(4.1) \quad -\operatorname{div}(\sigma \nabla p) = 0 \text{ on } \Omega,$$

with boundary conditions

$$p = 1 \text{ on } \partial\Omega^D \text{ (Dirichlet) and } \frac{\partial p}{\partial n} = 0 \text{ on } \partial\Omega^N \text{ (Neumann),}$$

where  $\partial\Omega = \partial\Omega^D \cup \partial\Omega^N$ . In this problem  $\partial\Omega^D$  is the top boundary of the domain. Figure 4.5 shows a part of the earth's crust. The depth of this part varies between 3 and 6 kilometers, whereas horizontally its dimensions are  $40 \times 60$  kilometers. The upper layer is a mixture of sandstone and shale and has a permeability of  $10^{-4}$ . Below this layer, shale and sandstone layers are present with permeabilities of  $10^{-7}$  and 10,

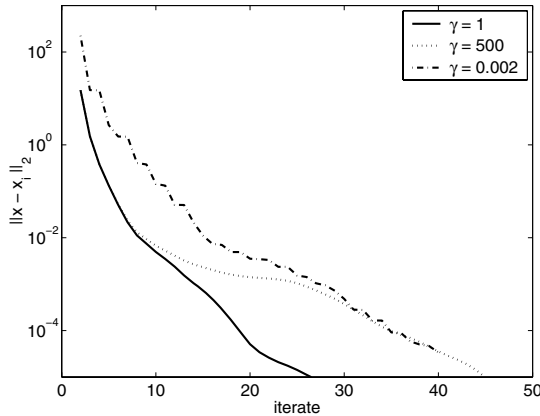


FIG. 4.4. Convergence of the balancing preconditioner with  $x_{0,B} = ZE^{-1}Z^Tb$  and  $\epsilon = 10^{-2}$ .

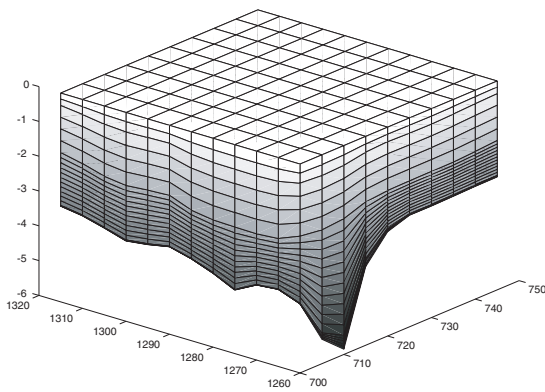


FIG. 4.5. The geometry of an oil flow problem with 4 layers.

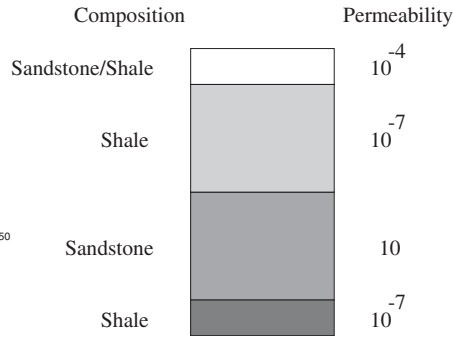


FIG. 4.6. Permeabilities for each layer.

TABLE 4.1  
The results for the oil flow problem.

Method	Deflation	Balancing
Iterations	36	36
CPU time in seconds	6.3	9.8

respectively. An incomplete Cholesky factorization with no fill in is used as preconditioner [14]. We consider a problem with nine layers. Five sandstone layers are separated by four shale layers. Due to the Dirichlet boundary condition at the top, the preconditioned matrix has four small eigenvalues. We use four *physical* projection vectors [26] and stop if  $\|r_k\|_2 \leq 10^{-5}$ . Trilinear hexahedral elements are used and the total number of gridpoints is equal to 148185. The results are given in Table 4.1. It appears that the norms of the residuals for both preconditioners are the same. Due to extra work per iteration the balancing preconditioner costs more CPU time. We note that it is possible to implement the balancing preconditioner such that the costs are the same as those of deflation [11]. However, this implementation also leads to the same difficulties as deflation if the matrix  $E^{-1}$  is perturbed. The computations



are performed on an AMD Athlon, 1.4 GHz processor with 256 Mb of RAM. The code is compiled with FORTRAN g77 on LINUX.

**5. Conclusions.** In this paper we compared various preconditioners for the numerical solution of partial differential equations.

We have given a detailed comparison of the well-known balancing preconditioner used in domain decomposition methods and the deflation preconditioner.

We proved that both preconditioners lead to almost the same spectra. The zero eigenvalues of the deflation preconditioned system are replaced by eigenvalues which are one if the balancing preconditioner is used. Thus the effective condition number of the deflated preconditioned system is always, i.e., for all deflation vectors and all restrictions and prolongations, below or equal to the condition number of the system preconditioned by the balancing preconditioner. Moreover, we proved that the A-norm of the errors of the iterates built by the deflation preconditioner is always below the A-norm of the errors of the iterates built by the balancing preconditioner. Hence, the CG-method applied to the deflated preconditioned system never converges slower than the CG-method applied to the system preconditioned by the balancing preconditioner. Additionally, the amount of work of one iteration of the deflation preconditioned system is less than the amount of work of one iteration of the balancing preconditioned system. Hence the deflation preconditioner leads to fewer iterations and each iteration having less work.

Moreover, we established that the deflation preconditioner and the balancing preconditioner produce the same iterates if one uses certain starting vectors. More precisely, we showed which terms in the PCG-method are the same for both methods and which terms are different. Numerical results for porous media flows emphasized the theoretical results.

#### REFERENCES

- [1] J. H. BRAMBLE, J. E. PASCIAK, AND A. H. SCHATZ, *The construction of preconditioners for elliptic problems by substructuring*. I, Math. Comp., 47 (1986), pp. 103–134.
- [2] H. DE GERSEM AND K. HAMEYER, *A deflated iterative solver for magnetostatic finite element models with large differences in permeability*, Eur. Phys. J. Appl. Phys., 13 (2000), pp. 45–49.
- [3] M. DRYJA AND O. B. WIDLUND, *Schwarz methods of Neumann-Neumann type for three-dimensional elliptic finite element problems*, Comm. Pure Appl. Math., 48 (1995), pp. 121–155.
- [4] 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.
- [5] C. FARHAT AND F.-X. ROUX, *A method of finite element tearing and interconnecting and its parallel solution algorithm*, Internat. J. Numer. Methods Engrg., 32 (1991), pp. 1205–1227.
- [6] C. FARHAT AND F.-X. ROUX, *Implicit parallel processing in structural mechanics*, Comput. Mech. Adv., 2 (1994), pp. 1–124.
- [7] J. FRANK AND C. VUIK, *On the construction of deflation-based preconditioners*, SIAM J. Sci. Comput., 23 (2001), pp. 442–462.
- [8] R. A. HORN AND C. R. JOHNSON, *Matrix Analysis*, Cambridge University Press, Cambridge, UK, 1985.
- [9] E. F. KAASSCHIETER, *Preconditioned conjugate gradients for solving singular systems*, J. Comput. Appl. Math., 24 (1988), pp. 265–275.
- [10] L. YU. KOLOTILINA, *Preconditioning of systems of linear algebraic equations by means of twofold deflation*. I. *Theory*, J. Math. Sci. (New York), 89 (1998), pp. 1652–1689.
- [11] J. MANDEL, *Balancing domain decomposition*, Comm. Numer. Methods Engrg., 9 (1993), pp. 233–241.
- [12] J. MANDEL AND M. BREZINA, *Balancing domain decomposition for problems with large jumps in coefficients*, Math. Comp., 216 (1996), pp. 1387–1401.

- [13] 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.
- [14] J. A. MEIJERINK AND H. A. VAN DER VORST, *An iterative solution method for linear systems of which the coefficient matrix is a symmetric  $M$ -matrix*, Math. Comp., 31 (1977), pp. 148–162.
- [15] R. B. MORGAN, *A restarted GMRES method augmented with eigenvectors*, SIAM J. Matrix Anal. Appl., 16 (1995), pp. 1154–1171.
- [16] 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.
- [17] R. A. NICOLAIDES, *Deflation of conjugate gradients with applications to boundary value problems*, SIAM J. Numer. Anal., 24 (1987), pp. 355–365.
- [18] L. F. PAVARINO AND O. B. WIDLUND, *Balancing Neumann-Neumann methods for incompressible Stokes equations*, Comm. Pure Appl. Math., 55 (2002), pp. 302–335.
- [19] A. QUARTERONI AND A. VALLI, *Domain Decomposition Methods for Partial Differential Equations*, Oxford Science Publications, Oxford, UK, 1999.
- [20] B. SMITH, P. BJØRSTAD, AND W. GROPP, *Domain Decomposition*, Cambridge University Press, Cambridge, UK, 1996.
- [21] A. TOSELLI AND W. WIDLUND, *Domain Decomposition Methods*, Springer Ser. Comput. Math. 34, Springer-Verlag, Berlin, 2005.
- [22] A. VAN DER SLUIS AND H. A. VAN DER VORST, *The rate of convergence of conjugate gradients*, Numer. Math., 48 (1986), pp. 543–560.
- [23] F. VERMOLEN, C. VUIK, AND A. SEGAL, *Deflation in preconditioned conjugate gradient methods for finite element problems*, in Conjugate Gradient and Finite Element Methods, M. Křížek, P. Neittaanmäki, R. Glowinski, and S. Korotov, eds., Springer-Verlag, Berlin, 2004, pp. 103–129.
- [24] C. VUIK, A. SEGAL, J. A. MEIJERINK, AND G. T. WIJMA, *The construction of projection vectors for a deflated ICCG method applied to problems with extreme contrasts in the coefficients*, J. Comput. Phys., 172 (2001), pp. 426–450.
- [25] C. VUIK, A. SEGAL, AND J. A. MEIJERINK, *An efficient preconditioned CG method for the solution of a class of layered problems with extreme contrasts in the coefficients*, J. Comput. Phys., 152 (1999), pp. 385–403.
- [26] C. VUIK, A. SEGAL, L. EL YAAKOUBI, AND E. DUFOUR, *A comparison of various deflation vectors applied to elliptic problems with discontinuous coefficients*, Appl. Numer. Math., 41 (2002), pp. 219–233.