

# REDUCED ORDER SOLUTION OF STRUCTURED LINEAR SYSTEMS ARISING IN CERTAIN PDE-CONSTRAINED OPTIMIZATION PROBLEMS \*

V. SIMONCINI<sup>†</sup>

**Abstract.** The solution of PDE-constrained optimal control problems is a computationally challenging task, and it involves the solution of structured algebraic linear systems whose blocks stem from the discretized optimality first-order conditions. In this paper we analyze the numerical solution of this large-scale system: we first perform a natural order reduction, and then we solve the reduced system iteratively by exploiting specifically designed preconditioning techniques. The analysis is accompanied by numerical experiments on two application problems.

**Key words.** Structured linear systems, iterative methods, PDE-constraints, optimization, preconditioning.

**AMS subject classifications.** 65F10.

**1. Introduction.** We consider solving the following structured linear system

$$\begin{bmatrix} A_1 & 0 & B_1^\top \\ 0 & A_2 & B_2^\top \\ B_1 & B_2 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} f_1 \\ f_2 \\ f_3 \end{bmatrix}, \quad (1.1)$$

where  $A_1$  is symmetric and positive definite matrix,  $A_2$  is symmetric positive semidefinite and highly singular, and  $B_2$  is square and nonsingular, so that the whole matrix is nonsingular;  $B_1^\top$  denotes the transpose of  $B_1$ . This type of structure is very common in a large number of application problems; here we are mainly interested in linear systems stemming from the discretization of PDE-constrained optimal control problems. In this context, several authors have addressed the problem of efficiently solving (1.1), and various preconditioners have very recently been explored that try to fully take into account the special coefficient matrix structure; we refer to, e.g., [38], [37], [31], [27], [24], [3] and references therein.

Within a discretize-then-optimize framework, systems similar to (1.1) need to be solved many times, usually either because they are used to tune a numerical model, or because they represent a step in an iterative (outer) nonlinear solution procedure. The computational cost associated with the solution of (1.1) is therefore a major indicator of the overall suitability of the numerical model for practical engineering purposes. We stress that the high singularity of  $A_2$ , with no other hypotheses, poses an additional challenge for the linear system solution.

A key fact in the solution of (1.1) is that the matrix blocks can be extremely large already for rather coarse mesh discretizations, because of a three or higher dimensional setting. This fact is particularly crucial in so-called “all-at-once” strategies, where the time frame is also discretized, and the matrix blocks in (1.1) include subblocks corresponding to different time instants [16],[24], [37]. In this framework, an initial reduction of the matrix size may lead to significant savings, as long as this reduction does not entail extra computational burden, such as additional inner solves or approximation of variables. In this paper we report on our numerical experience with

---

\*Version of January 18, 2011.

<sup>†</sup>Dipartimento di Matematica, Università di Bologna, Piazza di Porta S. Donato 5, I-40127 Bologna, Italy and CIRSA, Ravenna, Italy (valeria@dm.umibo.it).

a particularly simple and effective reduction. The reduction process fully relies on the assumption that  $A_1$  is easy to solve with (e.g.,  $A_1$  is diagonal), which is usually the case in the problems we are addressing. We discuss two realistic and popular specific problems, that require solving (1.1). The first one is a PDE-constrained optimal control problem with distributed constraints, with linear quadratic functional [38]. This simplified problem has become a popular benchmark for analyzing several numerical devices. We will see that in this case, the reduction process does not even need to solve with  $A_1$ . The second application we consider is a simplified Monge-Kantorovich mass transfer problem, used in the context of image registration [3]. Although we focus on two very specific applications, our discussion can be adapted to other settings, and we hope it may serve as a general platform for related and more general frameworks.

The solution of the reduced system exploits the obtained structure to build effective preconditioning techniques, which allow us to solve medium size 3D problems in a small amount of time on a commodity workstation. We also perform an ad-hoc spectral analysis that in some cases ensures optimality of the adopted preconditioning strategy. As it will be clear in the discussion of the specific problems, the use of the reduced formulation allows us to eliminate some redundant information of the original formulation, so that the preconditioning step need only be concerned with the relevant blocks in the matrix.

A synopsis of the paper is as follows. Section 2 shows the reduction step, based on the Schur complement formulation, and reviews and analyzes some general preconditioning strategies that are employed in later sections. Section 3 introduces the matrix setting of the PDE-constrained optimal control problem, and provides a detailed theoretical analysis of the proposed preconditioning approaches; section 3.1 reports on our numerical experience with these solution strategies. Section 4 is devoted to the matrix description of the simplified Monge-Kantorovich mass transfer problem, and of its numerical solution, while section 4.1 reports on our thorough numerical experience with these data. Finally, section 5 collects our concluding remarks.

Throughout the paper we shall use the Euclidean norm for vectors, and the induced norm for matrices. We shall denote by  $\text{spec}(A)$  the set of eigenvalues (spectrum) of a given matrix  $A$ .

**2. The reduced order problem and preconditioning strategies.** The size of each block of the coefficient matrix in (1.1) depends on the fineness of the given domain (2D or 3D in space possibly plus time), so that the problem dimension becomes very high already for rather coarse discretizations. The prospect of reducing the dimension of (1.1) is therefore of great interest, especially if this does not entail higher computational complexity. Indeed, we can decrease the number of blocks from three to two, by using the Schur complement formulation with respect to  $A_1$  (so-called “static condensation”): the first equation block gives  $x_1 = A_1^{-1}(f_1 - B_1^\top x_3)$ , which substituted into the original system yields the reduced system

$$\begin{bmatrix} A_2 & B_2^\top \\ B_2 & -B_1 A_1^{-1} B_1^\top \end{bmatrix} \begin{bmatrix} x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} f_2 \\ f_3 - B_1^\top A_1^{-1} f_1 \end{bmatrix} \quad \Leftrightarrow \quad Sx = b. \quad (2.1)$$

We recall here that we assume that  $A_2$  is positive semidefinite and singular, while  $B_1 A_1^{-1} B_1^\top$  is positive semidefinite, with the possibility of being nonsingular if  $B_1$  is full row-rank. It is interesting that this is a slightly different setting than the one often encountered in the literature:  $A_2$  is often required to be positive definite on the kernel of  $B_2$ , but since  $B_2$  is nonsingular here, such condition is empty. Thanks to this nonsingularity, both diagonal blocks can be highly singular. A sufficient condition

for the solvability of (2.1) is that  $B_2^\top + A_2 B_2^{-1} B_1 A_1^{-1} B_1^\top$  be nonsingular. Under the hypothesis of nonsingularity of the original system (1.1), such condition is satisfied.

This reduction procedure is classical, and it is often discarded because in general it requires explicitly dealing with the inverse in  $B_1 A_1^{-1} B_1^\top$ . However, it turns out that for various PDE-constrained optimization problems the matrix  $B_1 A_1^{-1} B_1^\top$  is very cheap to deal with, so that the reduction becomes extremely attractive.

The reduced problem (2.1) may also be derived by using the following decomposition

$$\left[ \begin{array}{c|cc} A_1 & 0 & B_1^\top \\ \hline 0 & A_2 & B_2^\top \\ B_1 & B_2 & 0 \end{array} \right] = \left[ \begin{array}{c|cc} I & 0 & 0 \\ \hline 0 & I & 0 \\ B_1 A_1^{-1} & 0 & I \end{array} \right] \left[ \begin{array}{c|cc} A_1 & 0 & B_1^\top \\ \hline 0 & A_2 & B_2^\top \\ 0 & B_2 & -B_1 A_1^{-1} B_1^\top \end{array} \right] \equiv \mathcal{L}\mathcal{U}.$$

The system in (1.1) can be rewritten as

$$\mathcal{L}\mathcal{U}x = \begin{bmatrix} f_1 \\ f_2 \\ f_3 \end{bmatrix} \Leftrightarrow \mathcal{U}x = \mathcal{L}^{-1} \begin{bmatrix} f_1 \\ f_2 \\ f_3 \end{bmatrix} \Leftrightarrow \mathcal{U}x = \begin{bmatrix} f_1 \\ f_2 \\ f_3 - B_2^\top A_1^{-1} f_1 \end{bmatrix},$$

from which the reduced system is obtained.

The system in (2.1) is once again in symmetric saddle point form, where the first diagonal block is only positive semidefinite, while the nonsingularity of the second diagonal block depends on the rank of  $B_1^\top$ . Note that the system may be written in a more familiar form after row and column permutation of the blocks, so that the two diagonal blocks are switched; a change of sign could also be applied, so as to have a positive definite leading block whenever  $B_1$  is full rank. We do not actually apply these changes, but we keep them in mind when choosing the preconditioners, as the choice of a suitable preconditioning strategy heavily depends on whether the second diagonal block is nonsingular. It is also important to realize that effective preconditioners can exploit the nonsingularity of the non-diagonal block.

A lot of effort is usually devoted to the determination of “ideal” preconditioners, with which the preconditioned problem can be solved in approximately the same number of iterations, irrespective of the size, leading to a mesh independent performance of the solver<sup>1</sup>; cf., e.g., [11], [20], [41]. On the other hand, we also stress that these ideal preconditioners may be difficult to either explicitly compute, or to cheaply approximate in a way so as to maintain their optimality properties. Therefore, great care should be put in choosing the ideal preconditioner, so that it can be then transformed into a feasible one.

*Nonsingular (2,2) block.* If the matrix  $B_1 A_1^{-1} B_1^\top$  is nonsingular, then a natural preconditioner is given by the following block diagonal matrix (see, e.g., [41, sec.3.2]):

$$P_d = \begin{bmatrix} \tilde{B}_2 C^{-1} \tilde{B}_2^\top & 0 \\ 0 & C \end{bmatrix}, \quad C \approx B_1 A_1^{-1} B_1^\top, \quad \tilde{B}_2 \approx B_2. \quad (2.2)$$

Note that the block  $\tilde{B}_2 C^{-1} \tilde{B}_2^\top$  is an approximation to the ideal matrix  $A_2 + \tilde{B}_2 C^{-1} \tilde{B}_2^\top$ . We refer to [34] for a detailed analysis of the algebraic spectral properties of the preconditioned matrix  $SP_d^{-1}$ . From a computational standpoint, we observe that

<sup>1</sup>The computational cost will still grow with the problem size, as basic operations are performed with longer vectors when finer discretizations are used.

systems of the form  $\tilde{B}_2 C^{-1} \tilde{B}_2^\top v = w$  arising when employing  $P_d$  can be solved in sequence as  $v = \tilde{B}_2^{-1} (C(\tilde{B}_2^{-\top} w))$ . Such computational advantage drove our choice of the first diagonal block in  $P_2$ , as opposed to the ideal one mentioned above. The actual matrices  $C$ ,  $\tilde{B}_2$  used in computation depend on the properties of the original blocks, and this will be discussed for the specific problems.

*Singular (2,2) block.* If the matrix  $B_1 A_1^{-1} B_1^\top$  is singular, and since the first diagonal block is also singular, a common choice consists of augmenting one of the diagonal blocks in the preconditioner. For instance, in the block diagonal case,

$$P_{ad} := \begin{bmatrix} D & 0 \\ 0 & \tilde{C}(D) \end{bmatrix}, \quad \tilde{C}(D) \approx C(D) := B_1 A_1^{-1} B_1^\top + B_2 D^{-1} B_2^\top, \quad (2.3)$$

where the symmetric and positive definite matrix  $D$  is chosen appropriately and  $\tilde{C}(D)$  is positive definite; see [26, 15, 12] and references therein. We also refer to [8] for the augmentation case when the (1,1) block of  $S$  is zero. For this augmented block diagonal preconditioner, we can apply the results in [12] to derive the following general spectral bounds for the preconditioned matrix  $P_{ad}^{-\frac{1}{2}} S P_{ad}^{-\frac{1}{2}}$ . Refined bounds will be obtained in section 4 for specific choices of the blocks.

PROPOSITION 2.1. [12, Prop. 3.3]. *Let*

$$\hat{S} := \begin{bmatrix} \hat{A}_2 & \hat{B}_2^\top \\ \hat{B}_2 & -\hat{T} \end{bmatrix},$$

and let  $\lambda_{\max}, \gamma_{\max}$  be the largest eigenvalues of  $\hat{A}_2$  and  $\hat{T}$ , respectively. Let  $\sigma_{\min}, \sigma_{\max}$  be the largest and smallest (nonzero) singular values of the square matrix  $\hat{B}_2$ . Then  $\text{spec}(\hat{S}) \subseteq \mathcal{I}^- \cup \mathcal{I}^+$ , where

$$\mathcal{I}^- = \left[ \frac{1}{2}(-\gamma_{\max} - \sqrt{\gamma_{\max}^2 + 4\sigma_{\max}^2}), \frac{1}{2}(\lambda_{\max} - \sqrt{\lambda_{\max}^2 + 4\sigma_{\min}^2}) \right]$$

and

$$\mathcal{I}^+ = \left[ \frac{1}{2}(-\gamma_{\max} + \sqrt{\gamma_{\max}^2 + 4\sigma_{\min}^2}), \frac{1}{2}(\lambda_{\max} + \sqrt{\lambda_{\max}^2 + 4\sigma_{\max}^2}) \right].$$

This quite general result already indicates that the length of both positive and negative intervals mainly depends on the magnitude of the singular values of  $\hat{B}_2$ , which in turn strongly depends on the choice of  $D$ , as the following result shows.

PROPOSITION 2.2. *Let  $\tilde{C}(D) = C(D)$  and  $\hat{S} = P_{ad}^{-\frac{1}{2}} S P_{ad}^{-\frac{1}{2}}$ . With the notation of Proposition 2.1, let  $\sigma$  be a singular value of  $\hat{B}_2$ . Then either  $\sigma = 1$  or  $\sigma^2 = \mu/(\mu+1) < 1$ , where  $\mu$  are the finite eigenvalues of the problem  $B_2 D^{-1} B_2^\top x = \mu B_1 A_1^{-1} B_1^\top x$ , which thus satisfy*

$$\frac{\sigma_{\min}(B_2)^2}{\|D\| \|B_1 A_1^{-1} B_1^\top\|} \leq \mu \leq \frac{\|D^{-1}\| \|B_2\|^2}{\lambda_{\min}(B_1 A_1^{-1} B_1^\top)},$$

where  $\lambda_{\min}(B_1 A_1^{-1} B_1^\top)$  is the smallest nonzero eigenvalue of  $B_1 A_1^{-1} B_1^\top$ .

*Proof.* The singular values of  $\hat{B}_2$  are the square roots of the eigenvalues in the eigenvalue problem  $B_2 D^{-1} B_2^\top x = \lambda (B_1 A_1^{-1} B_1^\top + B_2 D^{-1} B_2^\top) x$ . For  $x$  in the null

space of  $B_1^\top$ , we obtain  $\lambda = 1$ . Otherwise, assuming that  $\lambda \neq 1$ , we can rewrite this eigenvalue problem as  $B_2 D^{-1} B_2^\top x = \mu B_1 A_1^{-1} B_1^\top x$  with  $\lambda = \mu/(1 + \mu)$ , and the result follows.  $\square$

The proposition above shows that for  $\tilde{C}(D) = C(D)$ ,  $\sigma_{\max} = 1$  in Proposition 2.1, so that the external extremes of the two intervals  $\mathcal{I}^-, \mathcal{I}^+$  only depend on the norms of the diagonal blocks. The lower bound for  $\mu$ , and thus for  $\sigma_{\min}$  in Proposition 2.2, provides a good indicator towards the selection of  $D$  so as to improve the clustering of the two spectral intervals  $\mathcal{I}^-, \mathcal{I}^+$ . We will explore this further in section 4.

Regardless of the singularity of  $B_1 A_1^{-1} B_1$ , the following indefinite (constraint-type) block preconditioner may be considered:

$$\mathcal{P}_{indef} = \begin{bmatrix} D & \tilde{B}_2^\top \\ \tilde{B}_2 & -C \end{bmatrix}, \quad \tilde{B}_2 \approx B_2, \quad (2.4)$$

where the symmetric matrix  $C$  is appropriately chosen, as in the case of  $P_{ad}$ . We restrict the discussion to the choice  $D = 0$ . Since  $B_2$  is nonsingular, systems with  $\mathcal{P}_{indef}$  can be solved by solving with the non-diagonal blocks, in sequence. In particular,  $C$  is not required to be nonsingular. The indefiniteness of the symmetric preconditioner enforces the use either of a nonsymmetric solver such as GMRES or Simplified QMR (see [36] and references therein), or of a variant of the projected PCG [13]; see also [18] for additional alternatives. Experiments in literature show that the preconditioner  $\mathcal{P}_{indef}$  is particularly well suited for various classes of optimization problems, such as constrained quadratic optimization, cf., e.g., [5], [9], [17], [10], [19], [23], and this is fully confirmed by our numerical experiments in section 3.1. On the other hand, it should be mentioned that the positive definite block diagonal preconditioner  $P_d$  allows one to use a short-term symmetric solver, such as MINRES, whose convergence behavior is better understood than that of nonsymmetric solvers.

**3. A PDE-constrained optimal control problem.** In this section we discuss the algebraic setting of a benchmark PDE-constrained optimal control problem described in [38]. We briefly introduce the problem here, while we refer to [38] for a more complete presentation.

The problem can be stated as follows<sup>2</sup>. Given the bounded domain  $\Omega \subset \mathbb{R}^d$ ,  $d = 2, 3$ , and  $\hat{u}$  (the desired state) defined in  $\hat{\Omega} \subseteq \Omega$ , find  $u$  such that

$$\min_{u, f} \frac{1}{2} \|u - \hat{u}\|_{L_2(\hat{\Omega})}^2 + \beta \|f\|_{L_2(\Omega)}^2 \quad (3.1)$$

$$\text{s.t.} \quad -\nabla^2 u = f \quad \text{in } \Omega \quad (3.2)$$

$$u = \hat{u} \quad \text{on } \partial\Omega, \quad (3.3)$$

where  $\beta$  is the regularization parameter. This linear quadratic optimal control problem is a simplified model, as more complex formulations may include, e.g., box constraints for  $f$ , more complex PDE constraints and/or other PDE boundary conditions; see, e.g., [17], [25], [31] for some sample treatments of the resulting algebraic problems, and [39] for a general discussion of the PDE and optimal control framework. After appropriate discretization, the discrete cost functional associated with (3.1) is

<sup>2</sup>With a little abuse of notation, the same letters will be used here to denote scalar functions and vectors.

given by

$$\min_{u_h, f_h} \frac{1}{2} \|u_h - \hat{u}\|_2^2 + \beta \|f_h\|_2^2 = \min_{u, f} \frac{1}{2} u^\top \bar{M} u - u^\top u + \alpha_u + \beta f^\top M f$$

with  $\alpha_u = \|\hat{u}\|_2^2$ ,  $M$  the mass matrix, and  $\bar{M}$  a portion of the mass matrix corresponding to the part of the domain where  $\hat{u}$  is defined. The constraint  $-\nabla^2 u = f$  in  $\Omega$  gives the algebraic equation  $Ku = Mf + d_b$ , where  $K$  is the stiffness matrix, in this case the discretization of the Laplacian, and  $d_b$  accounts for the boundary conditions [38]. The solution with the Lagrange multipliers method deals with the following Lagrange function

$$\mathcal{L}(f, u, \ell) = \frac{1}{2} u^\top \bar{M} u - u^\top b + \alpha_u + \beta f^\top M f + \ell^\top (Ku - Mf - d_b),$$

and the associated optimality first order conditions give the following algebraic saddle point linear system

$$\begin{bmatrix} 2\beta M & 0 & -M \\ 0 & \bar{M} & K^\top \\ -M & K & 0 \end{bmatrix} \begin{bmatrix} f \\ u \\ \ell \end{bmatrix} = \begin{bmatrix} 0 \\ b \\ d_b \end{bmatrix}. \quad (3.4)$$

A lot of attention has been recently devoted to structured matrices of this form. Known block preconditioners, such as block diagonal, constraint-style and Bramble-Pasciak preconditioners have been successfully adapted to this setting ([31], [38], [28]) and to somewhat more general ones, see. e.g., [17], [20]. However, it has been somehow overlooked that the structure in (3.4) perfectly matches that in (1.1), so that the order reduction in (2.1) can be performed; but see [35], and also [2], [41, sec. 4.1] for similar considerations. Note that the presence of the nonsingular matrix  $M$  in two of the three places in the first block row of (3.4) allows one to readily eliminate the first unknown, as  $2\beta f = \ell$ . Therefore, the reduction process described in (2.1) is in this case particularly convenient. The reduction gives the system

$$\begin{bmatrix} \bar{M} & K^\top \\ K & -\frac{1}{2\beta} M \end{bmatrix} \begin{bmatrix} u \\ \ell \end{bmatrix} = \begin{bmatrix} b \\ d_b \end{bmatrix},$$

where  $\bar{M}$  may be singular, in case  $\hat{u}$  is not defined in the whole  $\Omega$  [38]; note also the particularly simple structure of the (2,2) block compared with the general case in (2.1). In addition, the matrix  $M$  is nonsingular and its inverse can either be cheaply computed (if  $M$  is diagonal) or easily approximated; see, e.g., [40]. Although the matrix structure for more complex PDE-constrained optimization problems may vary significantly, the reduction performed above is made possible simply by the presence of the regularizing term  $\beta\|u\|$ , which is often employed; see, e.g., [27]. We thus stress that the reduction step is extremely cheap to perform, and does not lead to a more difficult problem. Indeed, the main ingredients of the preconditioning step are the same as those used to solve the original problem (1.1). The simplified strategy however, eliminates some of the ‘‘redundant’’ information in the system. We will show that the resulting reduced system can be solved efficiently at low cost.

An ideal block diagonal preconditioner in this setting is given by  $P_d$  in (2.2) with the following blocks:

$$P_d = \begin{bmatrix} KC^{-1}K^\top & 0 \\ 0 & C \end{bmatrix}, \quad C = \frac{1}{2\beta} M. \quad (3.5)$$

In practice, a more cost efficient preconditioner is obtained for  $\tilde{K} \approx K$  a, say, (algebraic) multigrid operator, and  $\tilde{C} \approx C$ , with for instance  $\tilde{C} = \frac{1}{2\beta} \text{diag}(M)$ ; Note that with the discretization used in [38], the diagonal of  $M$  is constant. Other more sophisticated approaches have been devised, if more accurate approximations are desired, e.g., [28]. Under natural hypotheses on the given blocks (cf. [38]), the following corollary of Proposition 2.1 can be derived, ensuring mesh independence of the spectrum of the preconditioned matrix.

**COROLLARY 3.1.** *Under the hypotheses of Proposition 2.1, let  $h$  be the mesh parameter and  $d \in \{2, 3\}$ . Assume that  $C = \frac{\alpha h^d}{2\beta} I$  for some  $\alpha > 0$  independent of  $h$ , and that there exist  $\gamma_1, \gamma_2, \delta_1, \delta_2, \bar{\delta}_2$  all positive and independent of  $h$  such that*

$$\gamma_1 h^d \leq \lambda(K) \leq \gamma_2 h^{d-2}, \quad \delta_1 h^d \leq \lambda(M) \leq \delta_2 h^d, \quad 0 \leq \lambda(\bar{M}) \leq \bar{\delta}_2 h^d.$$

Finally, assume that  $\tilde{K}$  is chosen so that there exist  $\tau_1, \tau_2 > 0$  such that  $\tau_1 x^* \tilde{K} x \leq x^* K x \leq \tau_2 x^* \tilde{K} x$  for all nonzero vectors  $x$ . Then the spectrum of  $SP_d^{-1}$  with  $\tilde{K}$  in place of  $K$  in  $P_d$  is bounded independently of  $h$ .

*Proof.* Following Proposition 2.1, the matrix  $\hat{S}$  takes the form

$$\hat{S} = \begin{bmatrix} \hat{K}^{-1} C^{\frac{1}{2}} \bar{M} C^{\frac{1}{2}} \hat{K}^{-1} & \hat{K}^{-1} C^{\frac{1}{2}} K C^{-\frac{1}{2}} \\ C^{-\frac{1}{2}} K C^{\frac{1}{2}} \hat{K}^{-1} & -\frac{1}{2\beta} C^{-\frac{1}{2}} M C^{-\frac{1}{2}} \end{bmatrix}.$$

We will show that the spectral intervals of Proposition 2.1 applied to  $\hat{S}$  do not depend on  $h$ . To this end, we only need to show that the relevant eigenvalues of the (1,1) and (2,2) blocks are bounded independently of  $h$ , similarly for the singular values of the nondiagonal block. Since  $C$  is a multiple of the identity, these singular values satisfy  $\sigma_i(\hat{K}^{-1} C^{\frac{1}{2}} K C^{-\frac{1}{2}}) = \lambda_i(\hat{K}^{-1} K)$ , and these are bounded from below and from above by  $\tau_1$  and  $\tau_2$ , respectively, showing mesh independence.

For the (2,2) block we have  $\gamma_{\max} = \lambda_{\max}(\frac{1}{2\beta} C^{-\frac{1}{2}} M C^{-\frac{1}{2}})$  so that

$$\gamma_{\max} = \frac{2\beta}{\alpha h^d} \frac{1}{2\beta} \lambda_{\max}(M) \leq \frac{\delta_2 h^d}{\alpha h^d} = \frac{\delta_2}{\alpha}.$$

Finally, for the first diagonal block we have  $\lambda_{\max} = \lambda_{\max}(\hat{K}^{-1} C^{\frac{1}{2}} \bar{M} C^{\frac{1}{2}} \hat{K}^{-1})$  so that

$$\lambda_{\max} = \frac{\alpha h^2}{2\beta} \lambda_{\max}(\hat{K}^{-1} \bar{M} \hat{K}^{-1}) \leq \frac{\alpha h^d}{2\beta} \bar{\delta}_2 h^d \frac{\tau_2^2}{\gamma_1^2 h^{2d}} = \frac{\alpha \bar{\delta}_2 \tau_2^2}{2\beta \gamma_1^2}.$$

The result thus follows.  $\square$

The existence of  $\tau_1, \tau_2$  is called *spectral equivalence condition* between  $K$  and  $\tilde{K}$ , and it is often met, at least numerically, when  $\tilde{K}$  is chosen as a multigrid type operator [6], [7], [29].

Finally, we notice that our choice of  $P_d$  in (3.5) is closely related to the *robust* theoretical preconditioner in [41, sec.3.2] for the same problem, where the matrix  $\bar{M} + K C^{-1} K^\top$  is suggested for the (1,1) block instead of  $K C^{-1} K^\top$ . Indeed, under the hypotheses of Corollary 3.1, the two matrices are spectrally equivalent.

As mentioned in section .2, an alternative to  $P_d$  is an indefinite preconditioner. The one in (2.4) takes a particularly convenient form in this case:

$$P_{indef} = \begin{bmatrix} 0 & \tilde{K} \\ \tilde{K} & -C \end{bmatrix}, \quad P_{indef}^{-1} = \begin{bmatrix} \tilde{K}^{-1} C \tilde{K}^{-1} & \tilde{K}^{-1} \\ \tilde{K}^{-1} & 0 \end{bmatrix},$$

showing that the application of  $P_{indef}^{-1}$  to a vector requires two solves with  $\tilde{K}$ , one multiplication with  $C$  and one sum of vectors. Since there are no solves with  $C$ , we can select  $C = \frac{1}{2\beta}M$ . The following proposition gives an explicit form for the preconditioned matrix, from which spectral information can be deduced. The proof is a simple matrix-matrix multiplication and is therefore omitted.

**PROPOSITION 3.2.** *Consider the preconditioner  $P_{indef}$  with  $C = \frac{1}{2\beta}M$ . If  $\tilde{K} = K$ , the preconditioned matrix has the form*

$$SP_{indef}^{-1} = \begin{bmatrix} G + I & \bar{M}K^{-1} \\ 0 & I \end{bmatrix}, \quad G = \bar{M}K^{-1}CK^{-1},$$

and  $G$  singular. Therefore,  $\text{spec}(SP_{indef}^{-1}) = \{1\} \cup \text{spec}(G + I)$ .

Since  $G$  is singular, unit eigenvalues also arise in the (1,1) block of  $SP_{indef}^{-1}$ , so that some of the unit eigenvalues may have geometric multiplicity larger than one; we refer to [4, sec. 10.2] for a detailed description of the class of indefinite preconditioners and for other relevant references in the context of spectral analysis.

We observe that again for  $\tilde{K} = K$  and since  $C$  is symmetric and positive definite,  $G$  is similar to a symmetric and positive semidefinite matrix, from which it follows that the eigenvalues of  $SP_{indef}^{-1}$  are real and not smaller than one. The following results ensures that the spectrum of  $SP_{indef}^{-1}$  (for  $\tilde{K} = K$ ) is bounded independently of the mesh parameter, under natural conditions on the involved matrices; cf., e.g., [38, Th.2.1].

**COROLLARY 3.3.** *Under the hypotheses of Proposition 3.2, let  $h$  be the mesh parameter and  $d \in \{2, 3\}$ . Assume there exist  $\gamma_1, \gamma_2, \delta_1, \delta_2, \bar{\delta}_2$  all positive and independent of  $h$  such that*

$$\gamma_1 h^d \leq \lambda(K) \leq \gamma_2 h^{d-2}, \quad \delta_1 h^d \leq \lambda(M) \leq \delta_2 h^d, \quad 0 \leq \lambda(\bar{M}) \leq \bar{\delta}_2 h^d.$$

then  $\lambda(SP_{indef}^{-1}) = 1 + \eta$  with  $0 \leq \eta \leq \delta_2 \bar{\delta}_2 / (2\beta \gamma_1^2)$ .

*Proof.* From Proposition 3.2 we know that either  $\eta = 0$  or  $\eta$  is a positive eigenvalue of  $\bar{M}K^{-1}MK^{-1}$ . Since  $\eta \leq \|\bar{M}K^{-1}MK^{-1}\| \leq \frac{1}{2\beta} \|\bar{M}\| \|K^{-1}\| \|M\| \|K^{-1}\|$ , the result follows.  $\square$

For  $\beta \ll 1$ , the eigenvalues of  $SP_{indef}^{-1}$  are bounded by a possibly quite large quantity; cf. also [2] for similar remarks, but on a problem with nonsingular first diagonal block. Our experiments showed that indeed the eigenvalues do increase for small  $\beta$ , in agreement with this bound. This justifies the need for a few extra iterations in the experiments reported in Tables 3.1-3.2. On the other hand, since for these data  $\bar{M}$  has very low rank, then only few non-unit eigenvalues arise in  $SP_{indef}^{-1}$ , at least for  $\tilde{K} = K$ , explaining the overall low number of iterations.

As already mentioned, indefinite preconditioners try to mimic as much as possible the original matrix,  $S$  here. Practical computationally feasible preconditioners however, replace the diagonal blocks of the original matrix with convenient approximations. In addition to taking zero as (1,1) block, our choice of  $P_{indef}$  has the particular feature that also the nondiagonal blocks are approximated by means of  $\tilde{K}$ ; this is similar in spirit to the approach taken in [5] and [33]. For general symmetric and positive definite  $\tilde{K}$ , we can write the preconditioned problem as

$$SP_{indef}^{-1} = \begin{bmatrix} I + \frac{1}{2\beta} \bar{M} \tilde{K}^{-1} M \tilde{K}^{-1} & \bar{M} \tilde{K}^{-1} \\ 0 & I \end{bmatrix} + \begin{bmatrix} E & 0 \\ E & E \end{bmatrix}, \quad E = K \tilde{K}^{-1} - I.$$

This explicit form shows how the spectrum of  $SP_{indef}^{-1}$  changes from  $\text{spec}(I + \frac{1}{2\beta}\bar{M}\tilde{K}^{-1}M\tilde{K}^{-1}) \cup \{1\}$  as a (possibly nonlinear) function of  $\|E\|$ . Note that if for instance  $\tilde{K}^{-1}$  is spectrally equivalent to  $K$ , then  $\text{spec}(I + \frac{1}{2\beta}\bar{M}\tilde{K}^{-1}M\tilde{K}^{-1})$  satisfies a result similar to that of Corollary 3.3, since

$$\|\bar{M}\tilde{K}^{-1}M\tilde{K}^{-1}\| \leq \|\bar{M}\| \|\tilde{K}^{-1}K\|^2 \|K^{-1}\|^2 \|M\|.$$

Moreover,  $\|E\|$  does not depend on the mesh parameter either, ensuring mesh independence of the spectrum of  $SP_{indef}^{-1}$ . As is well known, however, eigenvalues alone may not completely explain the behavior of a nonsymmetric solver, and additional information on the eigenvectors should also be included. Nonetheless, our numerical experiments did not seem to indicate that the eigenbasis played a role. We refer to [32] for a more complete invariant subspace analysis of similarly perturbed structured matrices.

**3.1. Numerical experiments.** We consider the algebraic problem stemming from the finite element discretization with  $Q_1$  elements on a uniform mesh of a distributed control problem analyzed in [38]. In particular, we focus on the data described in [38, sec. 4], the 2D and 3D problems labelled Target 2. The domain is  $\Omega = [0, 1]^d$ ,  $d = 2, 3$ . In 2D the subset  $\hat{\Omega}$  is given by a circle centered at  $(\frac{5}{8}, \frac{3}{4})$  and radius  $\frac{1}{5}$ , union the boundary of  $\Omega$ . In 3D,  $\hat{\Omega}$  is a sphere with center at  $(\frac{5}{8}, \frac{3}{4}, \frac{7}{10})$  and radius  $\frac{1}{4}$ , union the boundary of  $\Omega$ . The target  $\hat{u}$ , defined in  $\hat{\Omega}$  is set equal to a constant value within the circle or sphere, and zero on  $\partial\Omega$ .

We analyze the performance of MINRES preconditioned by  $P_d$ , and of GMRES right preconditioned by  $P_{indef}$ ; in both cases the action of  $K^{-1}$  is substituted with that of  $\tilde{K}^{-1}$  computed as the Algebraic Multigrid method HSL\_MI20 ([6]) with all default parameters except for `control.one_pass_coarsen=1`, `control.v_iterations=5`, and `control.st_parameter=.85`. The stopping criterion is based on the relative residual norm. However, it should be noticed that while during the GMRES iterations it is possible to monitor the Euclidean norm of the residual, in MINRES the  $P_d^{-1}$ -norm of the residual can be monitored. This discrepancy implies that the final solutions may have different accuracies. For the smaller dimension problems we observed that the GMRES solution was at least one order of magnitude more accurate than the MINRES solution.

The number of iterations and CPU time to meet the stopping criteria with a tolerance of  $10^{-8}$  are reported in Table 3.1 and Table 3.2, for the 2D and the 3D examples, respectively. For various mesh discretizations we considered two typical values of  $\beta$ , namely  $\beta = 10^{-2}, 10^{-5}$ , according to [38]. Note that the reported  $n$  implies that the actual problem (1.1) has size  $3n$ , whereas the reduced one in (2.1) has size  $2n$ .

The digits show the mesh independence of the solver preconditioned with  $P_{indef}$ , in agreement with the discussion towards the end of section 3. Also MINRES preconditioned by  $P_d$  appears to have a mesh independent performance; thanks to Corollary 3.1 this was to be expected since, at least numerically, all hypotheses of Corollary 3.1 seem to be satisfied.

The cost of applying  $P_{indef}$  and  $P_d$  is very similar, therefore the much lower number of iterations when using  $P_{indef}$  makes this approach far more appealing, with CPU times that are one order of magnitude lower for  $\beta = 10^{-2}$ , and less than 20% than those of the problem solved with MINRES and  $P_d$  for  $\beta = 10^{-5}$ . It is also interesting to note that the number of iterations with  $P_{indef}$  seems to be rather insensitive to

$\beta$	$n$	GMRES w/ $P_{indef}$		MINRES w/ $P_d$	
		its	CPU Time	its	CPU Time
$10^{-2}$	961	3	0.15	31	0.18
	3969	3	0.17	28	0.37
	16129	3	0.29	25	1.22
	65025	3	0.88	25	4.28
	261121	4	4.54	27	21.31
$10^{-5}$	961	10	0.19	35	0.20
	3969	10	0.25	32	0.41
	16129	10	0.65	29	1.19
	65025	10	2.21	27	4.52
	261121	10	9.80	27	19.74

TABLE 3.1

2D problem. Number of iterations and CPU Time for solvers with  $P_{indef}$  and  $P_d$  and AMG, for various dimensions and values of the parameter  $\beta$ .

$\beta$	$n$	GMRES w/ $P_{indef}$		MINRES w/ $P_d$	
		its	CPU Time	its	CPU Time
$10^{-2}$	343	3	0.16	59	0.20
	3375	3	0.20	63	1.09
	29791	3	0.94	63	11.20
	250047	3	9.96	63	129.48
$10^{-5}$	343	8	0.13	57	0.20
	3375	9	0.30	61	1.03
	29791	9	2.04	64	11.28
	250047	9	23.02	65	132.36

TABLE 3.2

3D problem. Number of iterations and CPU Time for solvers with  $P_{indef}$  and  $P_d$ , for various dimensions and values of the parameter  $\beta$ .

the physical dimension of the problem (2D or 3D), whereas the performance of  $P_d$  degrades significantly when passing from the 2D to the 3D case.

We remark that in [38], the indefinite preconditioning strategy implemented within the Projected Preconditioned Conjugate Gradient (PPCG) method was used on the  $3n$  size problem [13]. If one is not willing to use a memory consuming nonsymmetric solver with  $P_{indef}$ , then PPCG is a viable effective strategy. Since the number of iterations of GMRES preconditioned by  $P_{indef}$  is extremely low, we did not find it necessary to employ PPCG; see also the discussion in [9].

**4. A simplified Monge-Kantorovich mass transfer problem.** In [3] the authors have considered the numerical solution of a simplified Monge-Kantorovich mass transfer problem. The proposed approach applies a Newton-type iteration to approximate the solution to a nonlinear system of equations, which considers an “all-at-once” discretization strategy of the original problem, stated in two space dimensions and in time. Here we provide a very brief review of the setting that leads to a  $3 \times 3$  block structured system similar to that in (1.1), while we refer to [3] for a full description of the application and of the employed discretization strategy. The linear system stems from a parameter identification problem, which in the most general form may be stated as a constrained optimization problem as  $\min_{u,m} \frac{1}{2} \|Qu - b\|^2 + \alpha_R R(m - m_r)$ , with the

constraint  $\mathcal{A}(m)u - q = 0$ . The latter driving equation is selected as the hyperbolic equation  $u_t + \nabla \cdot (um) = 0$ ,  $t \in [0, T]$ , and  $m$  represents the set of parameters to be identified, which are space dependent. As of the other quantities involved,  $b$  contains the observed data affected by noise;  $Q$  is the projector onto the portion of the space(-time) domain associated with  $b$ ;  $R$  is the regularization functional, together with its weight  $\alpha_R$ . According to [3],  $\alpha_R = 10$  in our experiments. After discretization and for appropriate choices of the operators, the discretized optimization problem discussed in [3] is posed as

$$\min_{u, m} \frac{1}{2} \|Qu - b\|^2 + \frac{1}{2} \xi u^\top L \text{diag}(m) m \quad (4.1)$$

$$\text{s.t. } u_t + \nabla \cdot (um) = 0, \quad (4.2)$$

with  $\xi = \alpha_R T h_t h_x^2$ , where  $h_t, h_x$  are the mesh parameters in time and space, respectively (in our experiments we shall use  $h_t = T/n_t$  and  $h_x = 1/(n_x - 1)$  for a finite difference discretization with  $n_t$  and  $n_x$  points, and the same number  $n_x$  of points in the  $x$  and  $y$  space directions). Here  $A(m)$  represents the matrix discretization in space and time of the constraint, using an implicit Lax-Friedrichs scheme. The associated Lagrangian function is given by  $\mathcal{L}(u, m, p) = \frac{1}{2} \|Qu - b\|^2 + \frac{1}{2} \xi u^\top L \text{diag}(m) m + p^\top V(A(m)u - q)$ , where  $V$  is a diagonal matrix such that  $Vq$  is a grid discretization of the continuous (in space) Lagrange multiplier. By using a Newton-type approximation, one obtains a sequence of linear systems in the form:

$$\begin{bmatrix} Q^\top Q & 0 & A^\top V \\ 0 & \xi \text{diag}(L^\top u) & G^\top V \\ VA & VG & 0 \end{bmatrix} \begin{bmatrix} \delta u \\ \delta m \\ \delta p \end{bmatrix} = - \begin{bmatrix} \mathcal{L}_u \\ \mathcal{L}_m \\ \mathcal{L}_p \end{bmatrix},$$

where  $G$  is the discretization of the Jacobian of  $A(m)u$  with respect to  $m$ , and  $A = A(m)$ . Clearly, both these matrices depend on the current approximation of  $m$  and  $u$ , and thus the whole system changes at each Newton iteration. By reordering rows and columns we obtain a linear system in the form (1.1):

$$\begin{bmatrix} \xi \text{diag}(L^\top u) & 0 & G^\top V \\ 0 & Q^\top Q & A^\top V \\ VG & VA & 0 \end{bmatrix} \begin{bmatrix} \delta m \\ \delta u \\ \delta p \end{bmatrix} = - \begin{bmatrix} \mathcal{L}_m \\ \mathcal{L}_u \\ \mathcal{L}_p \end{bmatrix},$$

in which the first diagonal block is diagonal and nonsingular, while the second block  $Q^\top Q$  is diagonal and highly singular, since we assume that  $b$  is only available in a small portion of the domain. Here  $A^\top V$  is square and nonsingular, whereas  $G^\top V$  is tall rectangular and possibly rank deficient, depending, e.g., on the boundary conditions imposed to define the discretized constraint, here periodic conditions. The reduction described in section 2 yields

$$\begin{bmatrix} Q^\top Q & A^\top V \\ VA & -VG(\xi \text{diag}(L^\top u))^{-1} G^\top V \end{bmatrix} \begin{bmatrix} \delta u \\ \delta p \end{bmatrix} = - \begin{bmatrix} \mathcal{L}_u \\ \mathcal{L}_p + VG(\xi \text{diag}(L^\top u))^{-1} \mathcal{L}_m \end{bmatrix}.$$

Therefore, the linear system above can be restated as that in (2.1) with  $A_2 = Q^\top Q$ ,  $B_2 = VA$ ,  $A_1 = \xi \text{diag}(L^\top u)$  and  $B_1 = VG$ ; note that here  $B_2$  is nonsymmetric. The reduced system is once again obtained at a negligible cost, while the dimension reduction is very significant, since in this application problem  $\xi \text{diag}(L^\top u)$  is twice as large as  $Q^\top Q$ .

Since both diagonal blocks are only semi-definite, we consider using the augmented block diagonal preconditioner  $\mathcal{P}_{ad}$  in (2.3). We are thus left to appropriately select the block  $D$  in (2.3). The matrix  $Q^\top Q$  is diagonal, with all zero diagonal entries in the top part, namely  $Q^\top Q = \text{blkdiag}(0, \Omega)$  with  $\Omega$  diagonal and nonsingular; if not already in this form, such ordering of the zero diagonal elements may be achieved by row and column permutation of the data in the reduced system. Let  $\mathcal{I}_Q$  be the projector onto the null space of  $Q^\top Q$  (this is nothing but a diagonal matrix with one's where the diagonal elements of  $Q^\top Q$  are zero). Following the approach in [3], we define

$$D = Q^\top Q + \gamma \mathcal{I}_Q = \begin{bmatrix} \gamma I & 0 \\ 0 & \Omega \end{bmatrix}, \quad (4.3)$$

with  $\gamma > 0$  to be chosen. Proposition 2.2 suggests that if  $\gamma^{-1} \approx \frac{\|B_1 A_1^{-1} B_1^\top\|}{\sigma_{\min}(B_2)^2}$ , then all singular values of the (1,2) block of the preconditioned matrix will be close to one, and thus the spectrum of  $SP_{ad}^{-1}$  will be nicely clustered. Since estimating  $\sigma_{\min}(B_2)$  is expensive, we content ourselves with the consideration that the finite eigenvalues  $\mu$  and corresponding eigenvectors  $x$  of the pencil  $(B_2 D^{-1} B_2^\top, B_1 A_1^{-1} B_1^\top)$  satisfy

$$\mu = \frac{x^\top B_2 D^{-1} B_2^\top x}{x^\top B_1 A_1^{-1} B_1^\top x},$$

(cf. Proposition 2.2) where  $B_2$  depends on  $O(h_t/h_x)$  and  $B_1$  is independent of the mesh [3]. Therefore, for  $h_x \approx h_t$ , if  $D$  mimics  $A_1$  then all  $\mu$  are not too far from one, and the matrix  $\widehat{B}_2$  is not ill-conditioned. This reasoning leads us to set

$$\gamma := \|A_1\|, \quad (4.4)$$

and our numerical experiments seem to support this choice. This choice is in agreement with the argument in [3, sec. 6] where it is stated that  $\gamma$  should scale like  $h_x^2 h_t$ , the way  $A_1$  does indeed. General qualitative arguments for other choices of this augmentation parameter can also be found in [14].

To continue our analysis of  $P_{ad}$ , we rewrite the eigenvalue problem  $P_{ad}^{-\frac{1}{2}} S P_{ad}^{-\frac{1}{2}} u = \lambda u$  as

$$\left[ \begin{array}{c|c} 0 & \widehat{B}_{12}^\top \\ \hline I & \widehat{B}_{22}^\top \\ \widehat{B}_{12} & \widehat{B}_{22} \end{array} \middle| \begin{array}{c} \widehat{T} \\ -\widehat{T} \end{array} \right] \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \lambda \begin{bmatrix} x \\ y \\ z \end{bmatrix}, \quad \text{with} \quad \begin{aligned} \widehat{B}_{12} &= \frac{1}{\sqrt{\gamma}} C^{-\frac{1}{2}} B_{12}, \\ \widehat{B}_{22} &= C^{-\frac{1}{2}} B_{22} \Omega^{-\frac{1}{2}}, \end{aligned} \quad (4.5)$$

and  $\widehat{T} = C^{-\frac{1}{2}} B_1 A_1^{-1} B_1^\top C^{-\frac{1}{2}}$ ; here we have split the original  $B_2$  as  $[B_{12}, B_{22}]$  conformingly with  $A_2$  and also  $D$  in (4.3). This explicit form shows that a necessary condition for the nonsingularity of the reduced system is that  $\widehat{B}_{12}$ , and thus also  $B_{12}$ , be full column rank. For this type of structure it is possible to derive sharp estimates for its positive and negative spectral intervals.

PROPOSITION 4.1. *Consider the preconditioner  $P_{ad}$  in (2.3) with*

$$\widetilde{C}(D) = C(D) = B_1 A_1^{-1} B_1^\top + B_2 D^{-1} B_2^\top \quad \text{and} \quad D = Q^\top Q + \gamma \mathcal{I}_Q, \quad \gamma > 0.$$

*Then the eigenvalues of  $SP_{ad}^{-1}$  are contained in  $\mathcal{I}^- \cup \mathcal{I}^+$ , with*

$$\mathcal{I}^- = \left[ -1, \frac{1}{2}(1 - \sqrt{5}) \right], \quad \mathcal{I}^+ = \left[ \frac{-1 + \sqrt{1 + 4\sigma_{\min}^2}}{2}, \frac{1}{2}(1 + \sqrt{5}) \right],$$

where  $\sigma_{\min} \leq 1$  is the smallest nonzero singular value of  $\widehat{B}_{12}$  in (4.5).

*Proof.* We readily notice that if  $z = 0$  in (4.5), then from the first block equation it also follows that  $x = 0$  ( $\lambda \neq 0$ , because of the nonsingularity of the matrix). Moreover, from the second and third equations we obtain  $y = \lambda y$  and  $\widehat{B}_{22}y = 0$ . Since  $\widehat{B}_{22}$  has full column rank, then the second equation is only satisfied for  $y = 0$ , which cannot hold as it would imply a zero eigenvector. From now on we can thus assume that  $z \neq 0$ .

With the notation introduced in Proposition 2.1, we first notice that the (1,1) block in (4.5), namely  $\widehat{A}_2 = \text{blkdiag}(0, I)$ , satisfies  $\lambda_{\max} = 1$ . Moreover, since  $\widehat{T} + \widehat{B}_2\widehat{B}_2^\top = I$  and both matrices on the left-hand side are positive semidefinite, we have  $\gamma_{\max} \leq 1$ . A similar reasoning shows that  $\sigma_{\max}(\widehat{B}_2) \leq 1$ , where  $\widehat{B}_2 = [\widehat{B}_{12}, \widehat{B}_{22}]$ . Therefore, Proposition 2.1 provides the right extreme of  $\mathcal{I}^+$ . Finding a sharp left positive extreme is more cumbersome, as a mere application of Proposition 2.1 would give a loose bound.

Let  $\lambda > 0$  be an eigenvalue of (4.5). Substituting the first and second matrix equations in the third one for  $\lambda \neq 0$ ,  $\lambda \neq 1$ , we obtain

$$\frac{1}{\lambda}\widehat{B}_{12}\widehat{B}_{12}^\top z + \frac{1}{\lambda-1}\widehat{B}_{22}\widehat{B}_{22}^\top z - \widehat{T}z - \lambda z = 0.$$

We thus eliminate the denominator and noticing once again that  $\widehat{T} + \widehat{B}_{12}\widehat{B}_{12}^\top + \widehat{B}_{22}\widehat{B}_{22}^\top = I$ , we obtain

$$-\lambda^3 z + \lambda^2 \widehat{B}_2 \widehat{B}_2^\top z + \lambda z - \widehat{B}_{12} \widehat{B}_{12}^\top z = 0. \quad (4.6)$$

Let us write  $z = z_0 + z_B$  with  $\widehat{B}_{12}^\top z_0 = 0$  and  $z_B \perp z_0$ . We multiply (4.6) from the left by  $z_0^\top$  and by  $z_B^\top$  in sequence, that is

$$(-\lambda^3 + \lambda)\|z_0\|^2 + \lambda^2 z_0^\top \widehat{B}_{22} \widehat{B}_{22}^\top z_0 + \lambda^2 z_0^\top \widehat{B}_{22} \widehat{B}_{22}^\top z_B = 0 \quad (4.7)$$

and

$$(-\lambda^3 + \lambda)\|z_B\|^2 + \lambda^2 z_B^\top \widehat{B}_{22} \widehat{B}_{22}^\top z_0 \quad (4.8)$$

$$+ \lambda^2 z_B^\top \widehat{B}_{12} \widehat{B}_{12}^\top z_B + \lambda^2 z_B^\top \widehat{B}_{22} \widehat{B}_{22}^\top z_B - z_B^\top \widehat{B}_{12} \widehat{B}_{12}^\top z_B = 0. \quad (4.9)$$

We note that for  $z_B = 0$  (and thus  $z_0 \neq 0$ ), and using  $\lambda > 0$ , equation (4.7) gives  $(-\lambda^2 + 1)\|z_0\|^2 + \lambda\|\widehat{B}_{22}^\top z_0\|^2 = 0$ . Using  $\lambda\|\widehat{B}_{22}^\top z_0\|^2 \geq 0$  we get  $(-\lambda^2 + 1)\|z_0\|^2 \leq 0$ , so that  $\lambda \geq 1$ . In the following we can assume that  $z_B \neq 0$ . Moreover, we can assume that  $\lambda < 1$ , otherwise 1 would be the sought after extreme. From (4.7) we get

$$\begin{aligned} \lambda^2 z_0^\top \widehat{B}_{22} \widehat{B}_{22}^\top z_B &= -\lambda\|z_0\|^2 + \lambda^3\|z_0\|^2 - \lambda^2\|z_0^\top \widehat{B}_{22}\|^2 \\ &= -\lambda(1 - \lambda^2)\|z_0\|^2 - \lambda^2\|z_0^\top \widehat{B}_{22}\|^2 \leq 0, \end{aligned}$$

so that from (4.8)-(4.9) we obtain

$$0 \leq (-\lambda^3 + \lambda)\|z_B\|^2 + \lambda^2 z_B^\top \widehat{B}_{12} \widehat{B}_{12}^\top z_B + \lambda^2 z_B^\top \widehat{B}_{22} \widehat{B}_{22}^\top z_B - z_B^\top \widehat{B}_{12} \widehat{B}_{12}^\top z_B,$$

that is,  $\lambda^3\|z_B\|^2 - \lambda^2 z_B^\top \widehat{B}_2 \widehat{B}_2^\top z_B - \lambda\|z_B\|^2 + \|\widehat{B}_{12}^\top z_B\|^2 \leq 0$ . Therefore we obtain the inequality (for  $z_B \neq 0$ )

$$\lambda^3 - \lambda^2\|\widehat{B}_2\|^2 - \lambda + \sigma_{\min}^2 \leq 0,$$

where  $\sigma_{\min}$  is the smallest nonzero singular value of  $\widehat{B}_{12}$ . Finally, using  $\lambda^3 \geq 0$  and  $\|\widehat{B}_2\| \leq 1$  (cf. Proposition 2.2) we obtain  $\lambda^2 + \lambda - \sigma_{\min}^2 \geq 0$  from which the lower extreme of  $\mathcal{I}_+$  follows.

Let  $\lambda < 0$ . We multiply (4.6) from the left by  $z^\top$  and notice that  $z^\top \widehat{B}_2 \widehat{B}_2^\top z \leq \|z\|^2$ , and  $-z^\top \widehat{B}_{12} \widehat{B}_{12}^\top z \leq 0$ . Therefore,

$$0 \leq \lambda(-\lambda^2 + \lambda + 1)\|z\|^2.$$

Since  $\lambda < 0$ , it must hold that  $-\lambda^2 + \lambda + 1 \leq 0$ , from which the bound  $\lambda \leq \frac{1}{2}(1 - \sqrt{5})$  follows. To obtain the left-hand bound, we observe that in (4.6) it holds that  $z^\top \widehat{B}_2 \widehat{B}_2^\top z \geq z^\top \widehat{B}_{12} \widehat{B}_{12}^\top z$ . Therefore, multiplying (4.6) from the left by  $z^\top$  we obtain

$$0 \geq \lambda(1 - \lambda^2)\|z\|^2 + (\lambda^2 - 1)z^\top \widehat{B}_{12} \widehat{B}_{12}^\top z,$$

that is,  $(1 - \lambda^2)(\lambda\|z\|^2 - z^\top \widehat{B}_{12} \widehat{B}_{12}^\top z) \leq 0$ . Since  $\lambda\|z\|^2 - z^\top \widehat{B}_{12} \widehat{B}_{12}^\top z \leq 0$  for all  $z \neq 0$ , it must be  $1 - \lambda^2 \geq 0$ , from which the bound  $-1 \leq \lambda$  follows.  $\square$

REMARK 4.2. Interval bounds of the type as in Proposition 4.1 could be obtained either directly, or upon matrix row and column permutations, from Proposition 2.1. However, they are either loose, or they require further estimates for some of the involved quantities.

We note that all bounds of Proposition 4.1 except for the left extreme of  $\mathcal{I}^+$  are independent of the spectral properties of the blocks and are thus mesh independent. Therefore, the performance of the *exact*  $P_{ad}$  (i.e., with  $\widetilde{C}(D) = C(D)$ ) will solely depend on  $\sigma_{\min}$ , and to a large extent on the choice of  $\gamma$ .

The following example shows that the bounds of Proposition 4.1 are sharp.

EXAMPLE 4.3. We consider the following data:

$$S = \left[ \begin{array}{cc|cc} 0 & 0 & -1 & \delta \\ 0 & 2 & 10 & 1 \\ \hline -1 & 10 & -20 & 0 \\ \delta & 1 & 0 & 0 \end{array} \right]$$

and  $\delta$  is freely chosen. The preconditioner  $\mathcal{P}_{ad}$  is determined with  $D = \text{diag}([10, 2])$  ( $\gamma = 10$ ) and  $C = B_1 A_1^{-1} B_1^\top + B_2 D^{-1} B_2^\top$ . The table below shows the true eigenvalues of  $SP_{ad}^{-1}$  and the estimates of Proposition 4.1 for various selections of  $\delta$ .

$\delta$		spectral values			
$10^{-2}$	$\text{spec}(SP_{ad}^{-1})$	-1.0000	-0.61805	0.00603	1.6180
	$\mathcal{I}^-, \mathcal{I}^+$	-1.0000	-0.61803	0.00599	1.6180
$10^0$	$\text{spec}(SP_{ad}^{-1})$	-1.0000	-0.66646	0.42423	1.5774
	$\mathcal{I}^-, \mathcal{I}^+$	-1.0000	-0.61803	0.33426	1.6180
$10^2$	$\text{spec}(SP_{ad}^{-1})$	-1.0000	-0.70458	0.99980	1.4194
	$\mathcal{I}^-, \mathcal{I}^+$	-1.0000	-0.61803	0.61797	1.6180

We note that the accuracy of all bounds can vary somewhat with  $\delta$ , but that for  $\delta = 10^{-2}$  all bounds are accurate.

If the exact diagonal block  $C$  of  $P_{ad}$  is replaced by an accurate approximation  $\widetilde{C}$ , then the spectrum of  $SP_{ad}^{-1}$  does not change significantly, and the four extremes of

Proposition 4.1 are each multiplied by a modest constant; cf., e.g., [23]. On the other hand, if  $\tilde{C}$  does not fully capture all spectral features of an ill-conditioned matrix  $C$ , then the matrix  $C\tilde{C}^{-1}$  may have an isolated, small cluster of eigenvalues close to zero. This is indeed the situation we encounter with the data in the application problem we discuss in the next section. Theorem 2.4 in [22] ensures that in this case, all eigenvalues of  $SP_{ad}^{-1}$  with  $\tilde{C}$  remain quite close to those obtained when using the exact  $C$  (cf. Proposition 4.1), except for a small cluster of eigenvalues that scale like the cluster near zero in  $C\tilde{C}^{-1}$ . The number of these eigenvalues is strictly related to the number of small eigenvalues of  $C\tilde{C}^{-1}$ . We will explore in detail this phenomenon in the next section.

**4.1. Numerical experiments.** In this section we report on our experience with the reduced dimension linear system stemming from the Monge-Kantorovich problem. The data stem from an image registration problem, where  $u$  represents the image density, which is known at time zero (initial image) and at time  $T = 1/8$  (final image). The data and discretization procedure by means of finite differences are as described in [3]. We stress here that since discretization in time is performed as for the other variables, once the number of nodes  $n_x, n_t$  is selected, each block of the reduced linear systems will have dimension  $n_x^2 n_t$ . Since our main interest is in the solution of the system, we will not report results for the whole simulation, corresponding to the numerical solution of the nonlinear problem. Instead, in most experiments we select a Newton cycle, typically the first one, and analyze the algebraic linear system to be solved at that cycle. It is important to realize that the numerical results do change at different cycles, although the qualitative performance is the same when comparing different strategies.

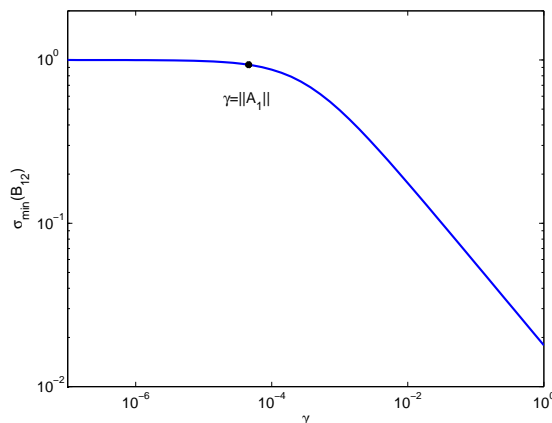


FIG. 4.1. Value of  $\sigma_{\min}(\hat{B}_{12})$  as  $\gamma$  varies (cf. Proposition 4.1).

We first consider solving the reduced system (4.1) by means of MINRES, preconditioned by  $P_{ad}$ , where the matrix  $D$  is obtained by replacing zero diagonal entries of  $Q^\top Q$  with  $\gamma$  (cf. (4.3) and (4.4)), and the second diagonal block of  $P_{ad}$  is given by  $C = B_1 A_1^{-1} B_1^\top + B_2 D^{-1} B_2^\top$ . For the proposed discretization  $A_1$  is diagonal, therefore its Euclidean norm can be computed very cheaply to get  $\gamma$ . To fully appreciate our choice of  $\gamma$ , in Figure 4.1 we show the value of  $\sigma_{\min}(\hat{B}_{12})$  in Proposition 4.1, for  $\gamma \in [10^{-7}, 1]$  and  $n = 1000$  (see below for more explanation on the data). The thick dot refers to the value  $\gamma = \|A_1\|$  in (4.4). Recalling that in general  $\sigma_{\min} \leq 1$ , the

estimates of Proposition 4.1 suggest that  $\sigma_{\min}$  should be as close as possible to one to have a good clustering, and this is clearly achieved by tiny values of  $\gamma$ . On the other hand, since  $\gamma$  affects the condition number of  $C$  and  $D$ , it should not be chosen to be too small, otherwise dealing with an inexact version of  $P_{ad}$  may become numerically difficult. The plot shows that our choice of  $\gamma = \|A_1\|$  appears to be very close to the trade-off value of  $\gamma$ : the largest possible value yielding  $\sigma_{\min} \approx 1$ . The choice of  $\gamma$  proposed in [3],  $\gamma = \text{mean}(\text{diag}(B_2 B_2^\top) / \text{diag}(B_1 A_1^{-1} B_1)) = 0.00249$  clearly gives a smaller value of  $\sigma_{\min}(\widehat{B}_{12})$ .

The exact application of the preconditioner yields the results shown in Table 4.1, for different discretization meshes:  $n_x \times n_x$  in 2D space, and  $n_t$  in time, yielding diagonal blocks both of size  $n_x^2 n_t \times n_x^2 n_t$ . In addition,  $VG$  is of size  $2n_x^2 n_t$ . Various combinations of  $n_x, n_t$  are reported.

TABLE 4.1

*Monge-Kantorovich problem. Performance (CPU Time in seconds and # iterations) of MINRES preconditioned by the exact version of  $\mathcal{P}_{ad}$ .  $n$  is the size of each of the (1,1) and (2,2) blocks. A dash stands for excessive memory requirements.*

$n_x$	$n_t$	$n$	# it	time	$n_x$	$n_t$	$n$	# it	time
10	10	1000	12	0.19	10	10	1000	12	0.19
15	15	3375	13	1.28	20	10	4000	16	3.44
20	20	8000	14	8.17	30	10	9000	21	34.16
25	25	15625	15	33.31	40	20	32000	21	376.54
30	30	27000	16	109.25	50	20	50000	-	-
35	35	42875	16	309.97					
40	40	64000	-	-					

The convergence tolerance, based on the preconditioned residual Euclidean norm, was set to  $10^{-10}$ . Note that this is a quite small tolerance, as significantly larger values are used within the Newton iteration; For instance,  $10^{-4}$  was used in [3]. We decided to set a strict tolerance because it allowed us to uncover some phenomena that were not visible otherwise, at least for small size problems and during the first few Newton iterations. In particular, the stagnation phases we will discuss later in this section are visible at the first Newton cycle only for a low residual norm (cf. left plot of Figure 4.2), as they also depend on the spectral distribution of the right-hand side. However, for finer discretizations or for different data, these phenomena were visible also at an earlier convergence stage, as shown in the right plot of Figure 4.2), where data from the last Newton step were used.

According to Table 4.1, both discretizations show that this preconditioner is “optimal”, in the sense that the number of iterations does not grow with the problem size. On the other hand, the computational costs of this “exact” preconditioner become quickly prohibitive, mostly due to the factorization and solution with the matrix  $C = C(D) = B_1 A_1^{-1} B_1^\top + B_2 D^{-1} B_2^\top$ . As in the previous application problem, we thus substitute the use of the exact  $C$  and its inverse with that of an Algebraic Multigrid operator; in the following experiments we selected once again the code HSL\_MI20 ([6]), with the parameters `control.one_pass_coarsen=1`, `control.v_iterations=5`, `control.st_parameter=.85`. In the following, the cost for generating the Multigrid operator will always be included in the reported timing. We note that  $C$  does not stem from the discretization of a standard second order elliptic operator, therefore we

do not expect optimality of HSL\_MI20. We also remark that the matrix  $C$  also arises in the numerical solution of the original problem in [3]: in there, the authors opted for an inner-outer procedure, in which a system with  $C$  is solved by means of a Krylov subspace method at each (outer) iteration. We do not explore this possibility here, as this involves the introduction of an extra parameter, the inner stopping tolerance; moreover, it was shown in [3] that this approach still displays some mesh dependence.

Table 4.4 shows that the computational costs decrease significantly when using the multigrid operator, although mesh independence seems to be lost, as the number of iterations grows with the grid refinement. This phenomenon will be analyzed more closely in the following. But first we would like to notice that our numerical experiments seem to support our choice of  $\gamma$ . Indeed, for  $\gamma$  as in [3] and  $n = 1000$ , MINRES with  $P_{ad}$  and AMG converged in 69 iterations and 0.47 seconds; for  $n = 27,000$  the method converged in 139 iterations and 38 seconds. These figures should be compared with the corresponding ones in Table 4.4. We notice, however, that the value of  $\gamma$  proposed by the authors of [3] was tailored towards their solution method, which differs from ours, therefore different performance may be expected in the two cases.

A close look at the convergence history of MINRES reveals that the convergence behavior is a little more subtle than the numbers in Table 4.4 suggest, and that phases of (almost) stagnation alternate with much steeper convergence rates. This pattern can be appreciated in Figure 4.2, where the MINRES convergence history for three mesh refinements is reported. The left plot uses data from the first Newton iteration, while the right plot uses those from the last Newton iteration, at convergence. We recall that both the coefficient matrix and the right-hand side change at each Newton iteration. Especially for the data in the first Newton iteration, we see that for a large portion of the convergence history, all curves look alike, suggesting an almost mesh independence. Moreover, only at some stage stagnation occurs in each of the histories. Such a pattern seems to indicate that the spectral distribution of  $SP_{ad}^{-1}$  is good, except perhaps for a few badly behaved eigenvalues, which delay convergence. Note also that stagnation occurs at different convergence stages in the two plots, and this mainly depends on the right-hand side eigendecomposition. Therefore, to completely avoid the stagnation phase, it would not be sufficient to stop the MINRES iteration earlier by setting a looser stopping tolerance, as stagnation may still occur.

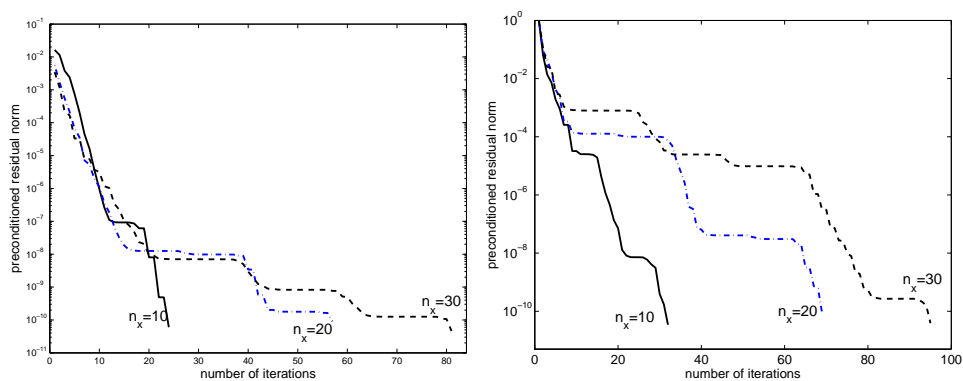


FIG. 4.2. Convergence history of MINRES with AMG-based  $P_{ad}$ , for three mesh refinements (here  $n_x = n_t$ ). Left: System data from first Newton iteration. Right: System data from last Newton iteration (at convergence).

TABLE 4.2  
*First approximate eigenvalues of the pencil  $(C, \tilde{C})$ .*

approx $\lambda_i$	$n_x = 10$	$n_x = 20$	$n_x = 30$
i=1	9.4632e-04	1.7676e-05	1.3247e-05
i=2	9.4999e-04	1.8274e-05	1.3937e-05
i=3	8.5302e-01	1.5449e-04	4.6116e-05
i=4	8.5400e-01	1.5484e-04	4.6200e-05
i=5	8.5722e-01	6.2973e-01	3.7707e-01
i=6	8.5752e-01	6.3079e-01	3.7791e-01
i=7	8.6410e-01	6.6089e-01	3.8408e-01
i=8	8.6825e-01	6.6091e-01	3.8409e-01
i=9	8.7875e-01	6.6557e-01	4.0425e-01
i=10	8.8097e-01	6.6674e-01	4.0514e-01

As already mentioned, this setting was recently analyzed in [22], where it is shown that the occurrence of such stagnation pattern may be associated with the non-full optimality of the preconditioning matrix. More precisely, assume that  $\tilde{C}$  is an approximation to  $C$  with both matrices symmetric and positive definite. It is shown in [22] that if the pencil  $(C, \tilde{C})$  has a small group of *positive* eigenvalues close to zero, well separated from the other eigenvalues, then the preconditioned matrix  $SP_{ad}^{-1}$  will inherit such a cluster: an isolated group of *negative* eigenvalues close to zero arises. Such a cluster of “badly” behaved eigenvalues causes temporary stagnation of MINRES, while the method tries to resolve the cluster.

Table 4.2 reports the approximation by means of 50 iterations of the Lanczos process (see, e.g., [1]) of the smallest 10 eigenvalues of  $(C, \tilde{C})$  for the three mesh refinements considered in Figure 4.2. In all cases the small group of well separated eigenvalues is readily visible. It is also interesting to notice that the cluster moves towards zero with the mesh refinements, causing a more severe stagnation phase for finer meshes. This observation is in agreement with the curves shown in Figure 4.2.

TABLE 4.3  
 *$j_*$ th Newton iteration. Solution via MINRES with the augmentation procedure when the approximate eigenvectors are computed either at iteration  $j_*$  or at the first iteration. Reported are number of iterations and CPU Time in seconds.*

$n_x$	$n_t$	$j_*$	defl.vec. at $j_*$ th it.		defl.vec. at first it.	
			# it	time	# it	time
20	10	4	21	0.80	21	0.79
30	10	5	30	2.29	30	2.26
40	20	5	38	15.81	38	16.55

We also notice that the occurrence of the small cluster does not seem to be due to some anomalous behavior of the AMG preconditioner. Indeed, the original matrix  $C$  does have a few small eigenvalues, which are apparently *not* captured by the approximation  $\tilde{C}$ , and thus they are still visible in the pencil  $(C, \tilde{C})$ . We observed a similar behavior when using an incomplete LU decomposition as  $\tilde{C}$ . The occurrence of this small cluster in  $C$  is problem and discretization dependent, therefore such a difficulty might not arise when analyzing slightly different application problems.

To cure the delay due to small tiny clusters, a deflation strategy is proposed in [22], which can be easily incorporated in the usual preconditioned MINRES iteration.

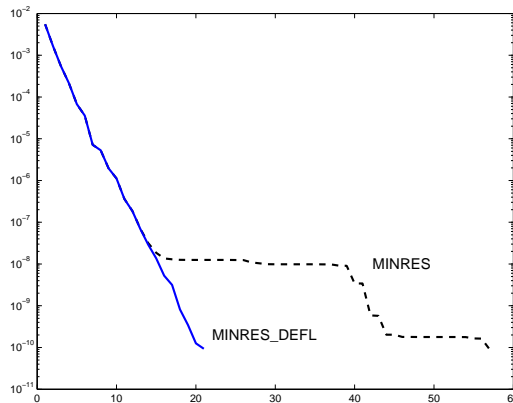


FIG. 4.3. Convergence history of MINRES with AMG-based  $P_{ad}$  without and with approximate deflation strategy;  $n_x = n_t = 20$ .

The approach consists of employing *approximate* eigenvectors of the negative cluster of  $SP_{ad}^{-1}$  to enrich the method recurrence, so that no iterations are used to approximate the corresponding quantities. It is important to realize that the approximate eigenvector does not have to be accurate, so that cheap devices can be used: if  $y$  is an eigen-approximation to an exact eigenvector of the pencil  $(C, \tilde{C})$ , then the vector  $[-D^{-1}B_1^T y; y]$  may be a good approximation to an eigenvector of the pencil  $(S, P_{ad})$  [22]. We adopted this strategy in the following experiments, where the vectors  $y$  were obtained with the Lanczos iteration [1], with a single cycle of dimension 50 on  $(C, \tilde{C})$ . The retained approximate eigenvectors were selected as those corresponding to eigenvalues of  $(C, \tilde{C})$  less than  $10^{-1}$ . A typical convergence history is depicted in Figure 4.3 where the convergence improvement can be fully appreciated. A more complete account of the results obtained with this approach is summarized in the right portion of Table 4.4. The rightmost values report the CPU time used to generate these approximations, and the number of eigenvectors retained. Although these timings are high with respect to the solution time, it should be kept in mind that this information can be recycled in subsequent Newton steps and thus can be fully amortized. As an example, Table 4.3 reports the results from solving the system at the  $j_*$ th Newton iteration (last iteration before convergence) when the augmenting approximate eigenvectors are computed: (i) at the same  $j_*$  iteration, or (ii) at the first iteration. No significant difference is observed between the two choices, suggesting that the approximate eigenvectors computed at the first Newton iteration can be effectively used throughout the whole process. If additional spectral information on the problem is available beforehand, then further savings can be obtained by avoiding the use of an eigenvalue routine.

For the sake of comparison, we also tried different parameter values in `hsl_mi20`, namely we selected `control.v_iterations=10` while all other values were selected as default (and not the ones chosen in previous experiments). The results are reported in Table 4.5, and show that a lower number of iterations can be achieved by allowing more V-cycles within the AMG operator; note however that the number of iterations fluctuates somewhat, although it does not seem to have an increasing pattern. We also notice that no clear stagnation phases could be detected during the convergence history in all cases. However, the costs of the AMG-based preconditioner are sig-

TABLE 4.4

Performance results of exact, inexact and deflated preconditioning variants. CPU time is in seconds. Data are taken from the first Newton iteration. The last column reports the time used to obtain the approximate eigenvectors of  $(C, \tilde{C})$  and the (automatically chosen) number of approximate eigenvectors retained.

$n_x$	$n_t$	$n$	$P_{ad}$		$P_{ad}$ w/AMG		$P_{ad}$ w/AMG+defl		
			# it	time	# it	time	# it	time	
10	10	1000	12	0.19	24	0.29	17	0.26	+ 0.52, 2 eigs
15	15	3375	13	1.28	52	1.43	25	0.80	+ 1.53, 5 eigs
20	20	8000	14	8.17	57	4.64	21	1.95	+ 4.51, 4 eigs
25	25	15625	15	33.31	85	14.64	36	6.68	+ 9.57, 5 eigs
30	30	27000	16	109.25	81	28.37	28	10.61	+19.03, 4 eigs
35	35	42875	16	309.97	122	65.82	47	27.33	+29.82, 5 eigs
40	40	64000	-		90	88.63	37	37.82	+52.23, 4 eigs
20	10	4000	16	3.44	46	1.49	18	0.70	+ 1.74, 4 eigs
30	10	9000	21	34.16	68	5.26	25	2.12	+ 4.36, 4 eigs
40	20	32000	21	376.54	76	31.63	32	12.55	+ 21.05, 4 eigs
50	20	50000	-		94	59.96	39	24.28	+ 33.21, 4 eigs

TABLE 4.5

Performance results of  $P_{ad}$  w/AMG, where the only non-default parameter was set to `control.v.iterations=10`. The digits report number of iterations and CPU time with data at the first Newton iteration.

$n_x$	$n_t$	$n$	$P_{ad}$ w/AMG( <code>v.IT=10</code> )	
			# it	time
10	10	1000	25	0.52
15	15	3375	29	1.82
20	20	8000	35	6.20
25	25	15625	43	16.81
30	30	27000	31	25.78
35	35	42875	56	70.61
40	40	64000	48	93.38
20	10	4000	33	2.54
30	10	9000	44	9.79
40	20	32000	40	45.22
50	20	50000	35	67.06

nificantly higher, compared with those obtained with the setting used in previous experiments. These results show the expected fact that by allowing a more accurate AMG preconditioner, performance in terms of number of iterations improves, at the price of higher costs. Nonetheless, other preconditioning techniques specifically adapted to this type of matrices may be better suited than the generic AMG operator we have used. With the current setting, the use of the deflation strategy appears to be more effective.

REMARK 4.4. We also implemented the indefinite preconditioner  $P_{indef}$  in (2.4) for this problem. However, the eigenvalue distribution of the preconditioned matrix, even in the exact case (that is, with  $\tilde{B}_2 = B_2$ ) was not favorable, as the (nonsymmetric) solver showed some stagnation phase. We could adopt a deflation-type strategy similar to what we proposed for the inexact  $P_{ad}$ , however the occurrence of a non-

symmetric matrix requires much more care in the implementation and accuracy of the deflation strategy, which is beyond the aim of this paper.

REMARK 4.5. Due to the more complex nature of the problem, we also experimented with nonsymmetric preconditioners, which require the use of a nonsymmetric solver. In particular, we explored a block (augmented) triangular preconditioner in the style of [3], namely

$$P_{ab} = \begin{bmatrix} D & B_2^\top \\ 0 & -C \end{bmatrix}, \quad C \approx B_1 A_1^{-1} B_1^\top + B_2 D^{-1} B_2^\top, \quad (4.10)$$

and  $D$  is a symmetric positive definite matrix completing the singular matrix  $A_2$ . In our tests,  $D$  and  $C$  were selected as in  $P_{ad}$ . However, for the realistic case where  $C$  is an (AMG) approximation to  $B_1 A_1^{-1} B_1^\top + B_2 D^{-1} B_2^\top$ , we observed convergence delays very similar to those described in Remark 4.4. For these reasons, we did not explore this preconditioner further. Note that in [3], where (4.10) was used to precondition (1.1), approximate solves with  $B_1 A_1^{-1} B_1^\top + B_2 D^{-1} B_2^\top$  were performed by an AMG-preconditioned Conjugate Gradient (PCG) iteration, giving rise to an inner-outer method, in which the inner PCG stopping tolerance was an extra parameter to be tuned; The outer method was necessarily a flexible solver [30].

**5. Conclusions.** We have shown that certain PDE-constrained optimization problems lead to structured algebraic linear systems, whose dimension can be significantly reduced without major computational overhead, thanks to the properties of the matrices defining the blocks. We have explored the numerical solution of these reduced problems by using different structured preconditioning techniques, reporting the effectiveness of the proposed strategies on two typical problems with data stemming from realistic applications.

**Acknowledgements.** We are grateful to Sue Thorne for providing us with the data on the distributed control problem of section 3.1. We are indebted with Michele Benzi, Eldad Haber and Lauren Taralli for making their full code available, which allowed us to extract both the data and important information for the experiments of section 4.1. We also thank Michele Benzi for helpful comments on an earlier version of this manuscript. All reported timings were obtained with Matlab ([21]), using the computer facilities of the SINCEM Laboratory at CIRSA, Ravenna.

#### REFERENCES

- [1] Z. Bai, J. Demmel, J. Dongarra, A. Ruhe, and H. van der Vorst, editors. *Templates for the solution of algebraic eigenvalue problems: a practical guide*. SIAM, Philadelphia, 2000.
- [2] Z.-Z. Bai, M. Benzi, F. Chen, and Z.-Q. Wang. Preconditioned MHSS iteration methods for a class of block two-by-two linear systems with applications to distributed control problems. Technical Report 2011-001, Math/CS Department, Emory University, January 2011.
- [3] M. Benzi, E. Haber, and L. Taralli. A preconditioning technique for a class of PDE-constrained optimization problems. Technical report, Department of Mathematics and Computer Science, Emory University, 2009. To appear in *Advances in Computational Mathematics*.
- [4] Michele Benzi, Gene H. Golub, and Jörg Liesen. Numerical solution of saddle point problems. *Acta Numerica*, 14:1–137, 2005.
- [5] L. Bergamaschi, J. Gondzio, M. Venturin, and G. Zilli. Inexact constraint preconditioners for linear systems arising in interior point methods. *Computational Optimization and Applications*, 36(2-3):137–147, 2007.
- [6] J. Boyle, M. D. Mihajlović, and J.A. Scott. HSL\_MI20: an efficient AMG preconditioner. *Int. J. Numer. Meth. Engrng.*, 82(1):64–98, 2010.
- [7] Achi Brandt. Algebraic multigrid theory: the symmetric case. *Applied Mathematics and Computation*, 19:23–56, 1986.

- [8] Zhi-Hao Cao. Augmentation block preconditioners for saddle point-type matrices with singular  $(1, 1)$  blocks. *Numerical Linear Algebra with Applications*, 15:515–533, 2008.
- [9] H. Sue Dollar, Nicholas I. M. Gould, Wil H. A. Schilders, and Andrew J. Wathen. Implicit-factorization preconditioning and iterative solvers for regularized saddle-point systems. *SIAM J. Matrix Anal. Appl.*, 28:170–189, 2006.
- [10] C. Durazzi and V. Ruggiero. Indefinitely preconditioned conjugate gradient method for large sparse equality and inequality constrained quadratic problems. *Num. Linear Algebra with Appl.*, pages 1–24, 2003.
- [11] Howard C. Elman, David J. Silvester, and Andrew J. Wathen. *Finite Elements and Fast Iterative Solvers, with applications in incompressible fluid dynamics*, volume 21 of *Numerical Mathematics and Scientific Computation*. Oxford University Press, 2005.
- [12] N. Gould and V. Simoncini. Spectral analysis of saddle point matrices with indefinite leading blocks. *SIAM J. Matrix Anal. Appl.*, 31:1152–1171, 2009.
- [13] N. I. M. Gould, M. E. Hribar, and J. Nocedal. On the solution of equality constrained quadratic programming problems arising in optimization. *SIAM J. Sci. Comput.*, 23:1376–1395, 2001.
- [14] C. Greif and D. Schötzau. Preconditioners for saddle point linear systems with highly singular  $(1, 1)$  blocks. *Electronic Transactions on Numerical Analysis*, 22:114121, 2006.
- [15] Chen Greif and Michael Overton. An analysis of low-rank modifications of preconditioners for saddle point systems. *Electronic Transactions on Numerical Analysis*, 37:307–320, 2010.
- [16] E. Haber and U. Ascher. Preconditioned all-at-once methods for large, sparse parameter estimation problems. *Inverse Problems*, 17:1847–1864, 2001.
- [17] R. Herzog and E. Sachs. Preconditioned conjugate gradient method for optimal control problems with control and state constraints. *SIAM J. Matrix Anal. Appl.*, 31(5):2291–2317, 2010.
- [18] P. Krzyzanowski. On block preconditioners for saddle point problems with singular or indefinite  $(1,1)$  block. *Numerical Linear Algebra with Applications*, 18(1):123–140, 2011.
- [19] L. Lukšan and J. Vlček. Indefinitely preconditioned inexact Newton method for large sparse equality constrained non-linear programming problems. *Num. Linear Algebra and Appl.*, 5:219–247, 1998.
- [20] Kent-Andre Mardal and Ragnar Winther. Preconditioning discretizations of systems of partial differential equations. *Numerical Linear Algebra with Applications*, 18(1):1–40, 2011.
- [21] The MathWorks, Inc. *MATLAB 7*, September 2004.
- [22] M. Olshanskii and V. Simoncini. Acquired clustering properties and solution of certain saddle point systems. *SIAM J. Matrix Anal. Appl.*, 31:2754–2768, 2010.
- [23] I. Perugia and V. Simoncini. Block-diagonal and indefinite symmetric preconditioners for mixed finite element formulations. *Numerical Linear Algebra with Applications*, 7:585–616, 2000.
- [24] T. Rees, M. Stoll, and A. Wathen. All-at-once preconditioning in PDE-constrained optimization. *Kybernetika*, 46(2):341–360, 2010.
- [25] T. Rees and A. Wathen. Preconditioning iterative methods for the optimal control of the Stokes equations. Technical report, 2010.
- [26] Tim Rees and Chen Greif. A preconditioner for linear systems arising from interior point optimization methods. *SIAM J. Sci. Comput.*, 29(5):1992–2007, 2007.
- [27] Tyrone Rees. *Preconditioning Iterative Methods for PDE Constrained Optimization*. PhD thesis, University of Oxford, 2010.
- [28] Tyrone Rees and Martin Stoll. Block-triangular preconditioners for PDE-constrained optimization. *Numerical Linear Algebra with Applications*, 17(6):977–996, Dec. 2010.
- [29] J. W. Ruge and K. Stüben. Algebraic multigrid. In *Multigrid Methods, Frontiers Appl. Math.* SIAM, Philadelphia, 1987.
- [30] Y. Saad. A flexible inner-outer preconditioned GMRES. *SIAM J. Sci. Comput.*, 14:461–469, 1993.
- [31] J. Schöberl and W. Zulehner. Symmetric indefinite preconditioners for saddle point problems with applications to PDE-constrained optimization problems. *SIAM J. Matrix Anal. Appl.*, 29(3):752–773, 2007.
- [32] D. Sesana and V. Simoncini. Spectral analysis of inexact constraint preconditioning for symmetric saddle point matrices. Technical report, Dipartimento di Matematica, Università di Bologna, 2010.
- [33] D. Silvester and M. Mihajlovic. A black-box multigrid preconditioner for the biharmonic equation. *BIT Numerical Mathematics*, 44:151–163, 2004.
- [34] D. Silvester and A. Wathen. Fast iterative solution of stabilized Stokes systems part II: using general block preconditioners. *SIAM J. Numer. Anal.*, 31:1352–1367, 1994.
- [35] V. Simoncini. Solution of structured algebraic linear systems in pde-constrained optimization problems. available at <http://www.dm.unibo.it/~simoncin/>, July 2010. Slides of the talk

- given at the “Erice 2010 Workshop on Nonlinear Optimization, Variational Inequalities and Equilibrium Problems, July 2-10, 2010”.
- [36] Valeria Simoncini and Daniel B. Szyld. Recent computational developments in Krylov subspace methods for linear systems. *nlaem*, 14(1):1–59, 2007.
  - [37] M. Stoll and A. Wathen. All-at-once solution of time-dependent PDE-constrained optimization problems. Technical Report 1017, The Mathematical Institute, University of Oxford, 2010.
  - [38] H. S. Thorne. Distributed control and constraint preconditioners. Technical Report RAL-TR-2010-016, Rutherford Appleton Laboratory, 2010.
  - [39] Fredi Tröltzsch. *Optimal control of partial differential equations : theory, methods and applications*. American Mathematical Society, Providence, RI, 2010. translated by J. Sprekels.
  - [40] Andy Wathen and Tyrone Rees. Chebyshev semi-iteration in preconditioning for problems including the mass matrix. *Electronic Transactions on Numerical Analysis*, 34:125–135, 2008-2009.
  - [41] Walter Zulehner. Non-standard norms and robust estimates for saddle point problems. NuMa-Report 2010-07, Institute of Computational Mathematics, Johannes Kepler University, November 2010.