# ITERATIVE METHODS FOR DOUBLE SADDLE POINT SYSTEMS

FATEMEH PANJEH ALI BEIK[1] AND MICHELE BENZI[2]

**Abstract.** We consider several iterative methods for solving a class of linear systems with double saddle point structure. Both Uzawa-type stationary methods and block preconditioned Krylov subspace methods are discussed. We present convergence results and eigenvalue bounds together with illustrative numerical experiments using test problems from two different applications.

**1. Introduction.** In this paper we consider iterative methods for solving large, sparse linear systems of equations of the form

$$
(1.1) \qquad \mathscr{A}u \equiv \begin{bmatrix} A & B^T & C^T \\ B & 0 & 0 \\ C & 0 & -D \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} \equiv b,
$$

where $A \in \mathbb{R}^{n \times n}$ is symmetric positive definite (SPD), $B \in \mathbb{R}^{m \times n}$, $C \in \mathbb{R}^{p \times n}$, and $D \in \mathbb{R}^{p \times p}$ is symmetric positive semidefinite (SPS) and possibly zero. Throughout the paper we assume that $n \geq m + p$.

Linear systems of the form (1.1) arise frequently from mixed and mixed-hybrid formulations of second-order elliptic equations [5, Sect. 7.2],[10] and elasticity [5, Sect. 9.3.1] problems. Numerical methods in constrained optimization [11, 12] and liquid crystal modeling [15] also lead to sequences of linear systems of the type (1.1). We further mention that finite element models of certain incompressible flow problems arising in the analysis of non-Newtonian fluids and in geophysics lead to large linear systems with coefficient matrices of the form

$$
\mathscr{B} = \begin{bmatrix} A & C^T & B^T \\ C & -D & 0 \\ B & 0 & 0 \end{bmatrix} \quad \text{and} \quad \mathscr{C} = \begin{bmatrix} -D & C & 0 \\ C^T & A & B^T \\ 0 & B & 0 \end{bmatrix};
$$

see, e.g., [1] and [6], respectively. It is easy to see that both $\mathscr{B}$ and $\mathscr{C}$ can be brought into the same form as matrix $\mathscr{A}$ in (1.1) by means of symmetric permutations (row and column interchanges).

It is important to observe that matrix $\mathscr{A}$ can be regarded as a $2 \times 2$ block matrix in two different ways, according to which of the following partitioning strategies is used:

$$
(1.2) \qquad \mathscr{A} = \left[ \begin{array}{c|cc} A & B^T & C^T \\ \hline B & 0 & 0 \\ C & 0 & -D \end{array} \right] \quad \text{or} \quad \mathscr{A} = \left[ \begin{array}{cc|c} A & B^T & C^T \\ B & 0 & 0 \\ \hline C & 0 & -D \end{array} \right].
$$

The first partitioning highlights the fact that problem (1.1) can in principle be treated as a "standard" saddle point problem, possibly stabilized (or regularized) when $D \neq 0$;

---

[1] Department of Mathematics, Vali-e-Asr University of Rafsanjan, PO Box 518, Rafsanjan, Iran (f.beik@vru.ac.ir).

[2] Department of Mathematics and Computer Science, Emory University, Atlanta, Georgia 30322, USA (benzi@mathcs.emory.edu). Work supported in part by NSF grant DMS-1418889.

see, e.g., [4]. On the other hand, the second partitioning shows that (1.1) can also be regarded as having a *double* saddle point structure, since the (1,1) block is itself the coefficient matrix of a saddle point problem; see, e.g., [15]. While in this paper we will make use of both partitionings, we are especially interested in studying solvers and preconditioners that make explicit use of the $3 \times 3$ block structure of $\mathscr{A}$.

The paper is orgnaized as follows. In section 2 we give a detailed discussion of conditions that ensure the unique solvability of (1.1). Section 3 is devoted to analyzing Uzawa-type stationary iterations for problem (1.1) based on the two partitionings (1.2). Block preconditioners for Krylov-type methods are discussed and analyzed in section 4. Illustrative numerical experiments are presented in section 5. Section 6 contains brief concluding remarks.

**2. Solvability conditions.** In this section we investigate the solvability of (1.1) under various assumptions on the blocks $A$, $B$, $C$ and $D$. Invertibility conditions for the coefficient matrix $\mathscr{A}$ in (1.1) under different assumptions on the blocks can be found scattered in the literature; see, for instance, [4], [5, Chapter 3], as well as [2] and [8] for eigenvalue bounds. While our results overlap in part with known ones, we find it useful to collect all the needed statements with complete proofs here, also in order to make the paper self-contained. In the following, for a real square matrix $A$ we write $A \succ 0$ ($A \succcurlyeq 0$) if $A$ is SPD (respectively, SPS) and $A \succ B$ ($A \succeq B$) if $A$ and $B$ are real symmetric matrices such that $A - B$ is SPD (respectively, SPS). Moreover, we write $(x; y; z)$ to denote the vector $(x^T, y^T, z^T)^T$.

The following theorem provides a necessary and sufficient condition for the invertibility of the matrix $\mathscr{A}$ in the case that the $(1,1)$ and $(3,3)$ blocks are both SPD.

PROPOSITION 2.1. *Assume that $A \succ 0$ and $D \succ 0$. Then matrix $\mathscr{A}$ is invertible if and only if $B^T$ has full column rank.*

*Proof.* Let $B^T$ have full column rank and assume that $\mathscr{A}u = 0$ for $u = (x; y; z)$, i.e.,

$$(2.1) \qquad Ax + B^T y + C^T z = 0,$$

$$(2.2) \qquad Bx \qquad\qquad = 0,$$

$$(2.3) \qquad Cx \qquad - Dz = 0.$$

If $x = 0$, then (2.3) implies $z = 0$ (since $D \succ 0$) and thus from (2.1) we conclude that $y = 0$, since $B^T$ has full column rank. Hence, $u = 0$. If $z = 0$, then from (2.1) and (2.2) we obtain $0 = Bx = -BA^{-1}B^T y$ and thus $y = 0$ since $BA^{-1}B^T$ is SPD. Hence, $x = 0$ and thus again it must be $u = 0$. Let us assume now that both of the vectors $x$ and $z$ are nonzero. Multiplying (2.1) by $x^T$ from the left, we find

$$(2.4) \qquad x^T Ax + x^T B^T y + x^T C^T z = 0.$$

From (2.3), it can be seen that $z^T Cx = z^T Dz$. Substituting $z^T Cx = z^T Dz$ and (2.2) into (2.4), we have

$$(2.5) \qquad x^T Ax = -z^T Dz.$$

In view of the positive definiteness of $A$ and $D$, the preceding equality implies that $x = 0$ and $z = 0$ which shows that $u = 0$.

Conversely, suppose that $\mathscr{A}$ is nonsingular. Let $y \in \mathbb{R}^m$ be such that $B^T y = 0$. Setting $u = (0; y; 0)$, we obtain $\mathscr{A}u = 0$. In view of the invertibility of $\mathscr{A}$, we conclude

that $y = 0$. This completes the proof. $\square$

Next, we consider relaxing the assumptions of Proposition 2.1 so that either $A \succeq 0$ or $D \succeq 0$. In the following theorem we establish sufficient conditions which guarantee the nonsingularity of $\mathscr{A}$. We further show that some of these conditions are also necessary.

THEOREM 2.2. *Let $A$ and $D \neq 0$ be SPS matrices. Assume that at least one of them is positive definite and $B^T$ has full column rank. Then the following statements hold:*

**Case 1.** *Suppose that $A \succ 0$ and $D \succeq 0$.*
- *If $\ker(C^T) \cap \ker(D) = \{0\}$ and $\mathrm{range}(B^T) \cap \mathrm{range}(C^T) = \{0\}$, then $\mathscr{A}$ is nonsingular.*
- *If $\mathscr{A}$ is nonsingular then $\ker(C^T) \cap \ker(D) = \{0\}$.*

**Case 2.** *Suppose that $A \succeq 0$ and $D \succ 0$.*
- *If $\ker(A) \cap \ker(B) \cap \ker(C) = \{0\}$ and $\mathrm{range}(B^T) \cap \mathrm{range}(C^T) = \{0\}$, then $\mathscr{A}$ is nonsingular.*
- *If $\mathscr{A}$ is nonsingular then $\ker(A) \cap \ker(B) \cap \ker(C) = \{0\}$.*

*Proof.* For clarity we divide the proof into two steps. In the first step we show the validity of the stated sufficient conditions for the invertibility of $\mathscr{A}$ for both cases. In the second step, it is proved that in each case one of the conditions is also necessary.
**Step I.** Let $u = (x; y; z)$ be an arbitrary vector such that $\mathscr{A}u = 0$. We recall from the proof of Proposition 2.1 that relation (2.5) must hold true.

Let us first consider the case that $A \succ 0$. From (2.5), it can be seen that $x = 0$, hence $Dz = 0$ from (2.3). Note that $B^T y + C^T z = 0$ together with the assumption $\mathrm{range}(B^T) \cap \mathrm{range}(C^T) = \{0\}$ imply that $C^T z = 0$ and $B^T y = 0$. Since $B^T$ has full column rank, $B^T y = 0$ implies $y = 0$. From $z \in \ker(C^T)$ and $Dz = 0$, we may immediately conclude from the assumption that $z = 0$, hence $u = 0$ and thus $\mathscr{A}$ is nonsingular.

For the second case, assume that $D \succ 0$. From (2.5), we can see that $z = 0$ since $A \succeq 0$. In addition $x^T A x = 0$ which implies that $Ax = 0$, i.e., $x \in \ker(A)$. Since $\mathscr{A}u = 0$, we have $Bx = 0$ and $Cx = 0$, i.e., $x \in \ker(B)$ and $x \in \ker(C)$. Consequently, we deduce that $x = 0$ and therefore $y = 0$ in view of the fact that $B^T$ has full column rank. Hence, $u = (x; y; z)$ is the zero vector, which shows the invertibility of $\mathscr{A}$.
**Step II.** Suppose that $\mathscr{A}$ is a nonsingular matrix.

Consider the case that $A \succ 0$. Assume there exists a nonzero vector $z \in \ker(C^T) \cap \ker(D)$. Then letting $u = (0; 0; z)$, we get $\mathscr{A}u = 0$, which is a contradiction. Hence $\ker(C^T) \cap \ker(D) = \{0\}$ is a necessary condition for the invertibility of $\mathscr{A}$.

Finally, let us consider Case 2 and show that $\ker(A) \cap \ker(B) \cap \ker(C) = \{0\}$ is a necessary condition for the invertibility of $\mathscr{A}$. If there exists a nonzero vector $x \in \ker(A) \cap \ker(B) \cap \ker(C)$, then for $u = (x; 0; 0)$, we have $\mathscr{A}u = 0$, which is again a contradiction. Therefore $\ker(A) \cap \ker(B) \cap \ker(C) = \{0\}$. $\square$

It is worth noting that the sufficient condition $\mathrm{range}(B^T) \cap \mathrm{range}(C^T) = \{0\}$ given in Theorem 2.2 is not a necessary condition for $\mathscr{A}$ to be invertible in the case that $D \neq 0$. We illustrate this fact with the following two simple examples in which $\mathscr{A}$ is nonsingular and $\mathrm{range}(B^T) \cap \mathrm{range}(C^T) \neq \{0\}$:

- For the case that $A \succ 0$ and $D \succcurlyeq 0$ is a nonzero singular matrix, consider the $8 \times 8$ matrix

$$\mathscr{A} = \begin{bmatrix} I_4 & B^T & C^T \\ B & 0 & 0 \\ C & 0 & -D \end{bmatrix},$$

  where $I_n$ stands for the $n \times n$ identity matrix, and the matrices $B$, $C$ and $D$ are given as follows:
  (2.6)
$$B = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 0 & 1 & 0 \end{bmatrix}, \quad C = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 \end{bmatrix}, \quad \text{and} \quad D = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}.$$

- For the case that $A \succcurlyeq 0$ is a singular matrix and $D \succ 0$, consider the $8 \times 8$ matrix $\mathscr{A}$ where $B$ and $C$ are defined as (2.6),

$$A = \begin{bmatrix} 0 & 0 \\ 0 & I_3 \end{bmatrix} \quad \text{and} \quad D = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}.$$

In both examples the matrix $\mathscr{A}$ is invertible and $\operatorname{range}(B^T) \cap \operatorname{range}(C^T) \neq \{0\}$.

The following proposition addresses the case where $D$ is a zero matrix. We beging by noting that in this case, a necessary condition for $\mathscr{A}$ to be invertible is that $C^T$ has full column rank. Indeed, if there exists a nonzero vector $z$ such that $C^T z = 0$ then $\mathscr{A} u = 0$ for $u = (0; 0; z) \neq 0$ and thus $\mathscr{A}$ cannot be invertible.

PROPOSITION 2.3. *Let $A \succ 0$ and assume that $B^T$ and $C^T$ have full column rank. Consider the linear system (1.1) with $D = 0$. Then $\operatorname{range}(B^T) \cap \operatorname{range}(C^T) = \{0\}$ is a necessary and sufficient condition for the coefficient matrix $\mathscr{A}$ to be invertible.*

*Proof.* As seen in the proof of Theorem 2.2, $\operatorname{range}(B^T) \cap \operatorname{range}(C^T) = \{0\}$ is a sufficient condition for invertibility of $\mathscr{A}$. Therefore we only need to show that it is also a necessary condition when $D = 0$ in (1.1). To this end, suppose that there exists a nonzero vector $v \in \operatorname{range}(B^T) \cap \operatorname{range}(C^T)$. As a result, $v = B^T y$ and $v = C^T z$ for some nonzero vectors $y$ and $z$ and letting $u = (0; y; -z)$, we get $\mathscr{A} u = 0$, contradicting the invertibility of $\mathscr{A}$. Hence it must be $\operatorname{range}(B^T) \cap \operatorname{range}(C^T) = \{0\}$. □

REMARK 2.4. We stress that in the case $D = 0$, both $B^T$ and $C^T$ must have full column rank for $\mathscr{A}$ to be invertible. In contrast, in the case that $D \succeq 0$ and $D \neq 0$, only the matrix $B^T$ is required to have full column rank while the matrix $C^T$ can be rank deficient.

In the remainder of the paper we will always assume that $\mathscr{A}$ is nonsingular.

**3. Uzawa-like iterative methods.** Uzawa-type methods have long been among the most popular algorithms for solving linear systems in saddle point form [4, Sect. 8.1]. In this section we study two variants of Uzawa's algorithm, motivated by the two possible block partitionings (1.2). We discuss first the case where the matrix $D$ in (1.1) is zero, and then the case $D \neq 0$.

**3.1. Uzawa-like iterative methods of the first type.** In this subsection we present two Uzawa-like iterative methods for solving (1.1) when $D = 0$. To this end we first consider the following two splittings for $\mathscr{A}$,

$$\mathscr{A} = \mathscr{M}_1 - \mathscr{N}_1 = \mathscr{M}_2 - \mathscr{N}_2,$$

where

$$\mathcal{M}_1 = \begin{bmatrix} A & 0 & 0 \\ B & -\frac{1}{\alpha}I & 0 \\ C & 0 & -\frac{1}{\beta}I \end{bmatrix}, \quad \mathcal{N}_1 = \begin{bmatrix} 0 & -B^T & -C^T \\ 0 & -\frac{1}{\alpha}I & 0 \\ 0 & 0 & -\frac{1}{\beta}I \end{bmatrix},$$

$$\mathcal{M}_2 = \begin{bmatrix} A & B^T & 0 \\ B & 0 & 0 \\ C & 0 & -\frac{1}{\alpha}I \end{bmatrix}, \quad \text{and} \quad \mathcal{N}_2 = \begin{bmatrix} 0 & 0 & -C^T \\ 0 & 0 & 0 \\ 0 & 0 & -\frac{1}{\alpha}I \end{bmatrix},$$

where $\alpha$ and $\beta$ are two given nonzero parameters.

The corresponding iterative schemes for solving (1.1) are given by

$$(3.1) \qquad u_{k+1} = \mathcal{G}_1 u_k + \mathcal{M}_1^{-1} b, \qquad k = 0, 1, 2 \ldots,$$

and

$$(3.2) \qquad u_{k+1} = \mathcal{G}_2 u_k + \mathcal{M}_2^{-1} b, \qquad k = 0, 1, 2 \ldots,$$

respectively, where $u_0$ is arbitrary,

$$(3.3) \qquad \mathcal{G}_1 = \mathcal{M}_1^{-1} \mathcal{N}_1 = \begin{bmatrix} 0 & -A^{-1}B^T & -A^{-1}C^T \\ 0 & I - \alpha B A^{-1} B^T & -\alpha B A^{-1} C^T \\ 0 & -\beta C A^{-1} B^T & I - \beta C A^{-1} C^T \end{bmatrix},$$

and

$$(3.4) \qquad \mathcal{G}_2 = \mathcal{M}_2^{-1} \mathcal{N}_2 = \begin{bmatrix} 0 & 0 & -\tilde{A}C^T \\ 0 & 0 & -S_B^{-1} B A^{-1} C^T \\ 0 & 0 & I - \alpha C \tilde{A} C^T \end{bmatrix},$$

with $\tilde{A} = A^{-1} - A^{-1}B^T S_B^{-1} B A^{-1}$ and $S_B = B A^{-1} B^T$.

In the rest of this subsection, we analyze the convergence properties of iterative methods (3.1) and (3.2).

PROPOSITION 3.1. *Assume that $A \succ 0$, $B^T$ and $C^T$ have full column rank, and* $\text{range}(B^T) \cap \text{range}(C^T) = \{0\}$. *Then all of the eigenvalues of the following matrix are real and positive for positive parameters $\alpha$ and $\beta$:*

$$\mathcal{S}_{\alpha,\beta} = \begin{bmatrix} \alpha B A^{-1} B^T & \alpha B A^{-1} C^T \\ \beta C A^{-1} B^T & \beta C A^{-1} C^T \end{bmatrix}.$$

*Proof.* Evidently, we have

$$\mathcal{S}_{\alpha,\beta} = \begin{bmatrix} \alpha I & 0 \\ 0 & \beta I \end{bmatrix} \begin{bmatrix} B A^{-1} B^T & B A^{-1} C^T \\ C A^{-1} B^T & C A^{-1} C^T \end{bmatrix}.$$

On the other hand, we can write

$$(3.5) \qquad \mathcal{S} = \begin{bmatrix} B A^{-1} B^T & B A^{-1} C^T \\ C A^{-1} B^T & C A^{-1} C^T \end{bmatrix} = \begin{bmatrix} B \\ C \end{bmatrix} A^{-1} \begin{bmatrix} B^T & C^T \end{bmatrix}.$$

Given the following partitioning for the matrix $\mathscr{A}$,

$$(3.6) \qquad \mathscr{A} = \left[ \begin{array}{c|cc} A & B^T & C^T \\ \hline B & 0 & 0 \\ C & 0 & 0 \end{array} \right] = \left[ \begin{array}{c|c} A & A_{12} \\ \hline A_{12}^T & 0 \end{array} \right],$$

we see that $\mathscr{S}$ is just the Schur complement $\mathscr{S} = A_{12}^T A^{-1} A_{12}$. Under our assumptions, $\mathscr{A}$ is invertible by Proposition 2.3, and therefore so is $\mathscr{S}$. Moreover, the positive definiteness of $A^{-1}$ implies the positive definiteness of $\mathscr{S}$. This shows that $\mathscr{S}_{\alpha,\beta}$ is the product of two SPD matrices and thus its eigenvalues must be real and positive. $\square$

In the sequel, we first discuss the convergence properties of the iterative method (3.1) and then conclude this subsection with a necessary and sufficient condition for the convergence of the iterative method (3.2). For the convergence analysis of iterative scheme (3.1) we need the following useful lemma, which is a special case of Weyl's Theorem [9, Theorem 4.3.1].

LEMMA 3.2. *Let $A$ and $B$ be two Hermitian matrices. Then,*

$$\lambda_{\max}(A + B) \leq \lambda_{\max}(A) + \lambda_{\max}(B),$$
$$\lambda_{\min}(A + B) \geq \lambda_{\min}(A) + \lambda_{\min}(B).$$

PROPOSITION 3.3. *Let $\mathscr{A}$ in (1.1) be nonsingular with $A \succ 0$ and $D = 0$. If the parameters $\alpha > 0$ and $\beta > 0$ satisfy*

$$(3.7) \qquad \alpha \lambda_{\max}(B A^{-1} B^T) + \beta \lambda_{\max}(C A^{-1} C^T) < 2,$$

*then the iterative scheme (3.1) is convergent for any initial guess, i.e., $\rho(\mathscr{G}_1) < 1$.*

*Proof.* Note that if $\mathscr{A}$ is nonsingular with $A \succ 0$, then in view of Proposition 2.3 and Remark 2.4 all of the assumptions in Proposition 3.1 are satisfied. Next, observe that the nonzero eigenvalues of

$$\mathscr{S}_{\alpha,\beta} = \left[ \begin{array}{c} \alpha B \\ \beta C \end{array} \right] A^{-1} \left[ \begin{array}{cc} B^T & C^T \end{array} \right]$$

are the same as those of

$$\mathscr{S}_1 = A^{-1} \left[ \begin{array}{cc} B^T & C^T \end{array} \right] \left[ \begin{array}{c} \alpha B \\ \beta C \end{array} \right] = \alpha A^{-1} B^T B + \beta A^{-1} C^T C.$$

From (3.3) it is clear that the iterative scheme (3.1) is convergent if and only if $\rho(I - \mathscr{S}_{\alpha,\beta}) < 1$. From Proposition 3.1, it is known that the eigenvalues of $\mathscr{S}_{\alpha,\beta}$ are all positive. As a result,

$$\bar{\lambda}_{\min}(\mathscr{S}_1) \leq \lambda(\mathscr{S}_{\alpha,\beta}) \leq \lambda_{\max}(\mathscr{S}_1),$$

where $\lambda(Q)$, $\lambda_{\max}(Q)$ and $\bar{\lambda}_{\min}(Q)$ respectively denote an arbitrary eigenvalue, the maximum eigenvalue and the minimum nonzero eigenvalue of a given matrix $Q$ having real and nonnegative spectrum. Writing again $S_B = B A^{-1} B^T$, $S_C = C A^{-1} C^T$, and using Lemma 3.2, it is easy to see that

$$(3.8) \qquad 1 - (\alpha \lambda_{\max}(S_B) + \beta \lambda_{\max}(S_C)) \leq \lambda(\mathscr{G}_1) \leq 1 - (\alpha \lambda_{\min}(S_B) + \beta \lambda_{\min}(S_C)).$$

From the above relation and invoking the fact that $S_B \succ 0$ and $S_C \succ 0$, we can see that a sufficient condition for the convergence of (3.1) is

$$-1 < 1 - (\alpha \lambda_{\max}(S_B) + \beta \lambda_{\max}(S_C)),$$

which is precisely (3.7). This completes the proof. □

REMARK 3.4. Under the assumptions of Proposition 3.3, it follows from (3.8) that

(3.9) $$\rho(\mathscr{G}_1) \leq f(\alpha, \beta),$$

where

$$f(\alpha, \beta) = \max\{\ |1 - (\alpha \lambda_{\max}(S_B) + \beta \lambda_{\max}(S_C))|, |1 - (\alpha \lambda_{\min}(S_B) + \beta \lambda_{\min}(S_C))|\ \}.$$

It can be verified that the upper bound $f(\alpha, \beta)$ is minimized (and less than 1) for any pair $(\alpha^*, \beta^*)$ with $\alpha^*, \beta^* > 0$ such that

(3.10) $$\alpha^*(\lambda_{\max}(S_B) + \lambda_{\min}(S_B)) + \beta^*(\lambda_{\max}(S_C) + \lambda_{\min}(S_C)) = 2,$$

or, equivalently,

$$\beta^* = \frac{2 - \alpha^*(\lambda_{\max}(S_B) + \lambda_{\min}(S_B))}{\lambda_{\max}(S_C) + \lambda_{\min}(S_C)}.$$

REMARK 3.5. In the special case when $\alpha = \beta$, (3.1) reduces to the standard Uzawa method based on the partitioning (3.6) of $\mathscr{A}$. In this case the conditions (3.7) and (3.10) simply become

$$0 < \alpha < \frac{2}{\lambda_{\max}(S_B) + \lambda_{\max}(S_C)},$$

and the upper bound on the spectral radius of the iteration matrix $\mathscr{G}_1$ is minimized for

$$\alpha^* = \frac{2}{\lambda_{\max}(S_B) + \lambda_{\max}(S_C) + \lambda_{\min}(S_B) + \lambda_{\min}(S_C)}.$$

We further recall that in the case that $\alpha = \beta$, the asymptotic convergence rate of Uzawa's method (3.1) is the same as that of the stationary Richardson iteration applied to the Schur complement system obtained by the eliminating the $x$ variable from (1.1), with the coefficient matrix $\mathscr{S}$ defined by (3.5). Under the assumptions of Proposition 3.1, we have $\mathscr{S} \succ 0$. From a well-known result on the convergence of the Richardson iteration (e.g., [16, Chapter 4]) we may conclude that a necessary and sufficient condition for the convergence of the iterative method (3.1) is given by

$$0 < \alpha < \frac{2}{\lambda_{\max}(\mathscr{S})},$$

and the optimum values of $\alpha$ is given by

$$\alpha^* = \frac{2}{\lambda_{\max}(\mathscr{S}) + \lambda_{\min}(\mathscr{S})},$$

which leads to the smallest possible spectral radius of the iteration matrix $\mathscr{G}_1$.

We conclude this subsection with a brief discussion of the convergence properties of the iterative method (3.2). To this end we first need the following two propositions. We recall that $S_B = BA^{-1}B^T$.

PROPOSITION 3.6. *Assume that $A \succ 0$ and $B^T$ has full column rank. Then* $\tilde{A} = A^{-1} - A^{-1}B^T S_B^{-1} BA^{-1} \succcurlyeq 0$.
*Proof.* Since $A$ is SPD, we can write

$$\tilde{A} = A^{-1/2}(I - A^{-1/2}B^T S_B^{-1} BA^{-1/2})A^{-1/2}.$$

The nonzero eigenvalues of

$$A^{-1/2}B^T S_B^{-1} BA^{-1/2}$$

are the same as those of

$$BA^{-1/2}A^{-1/2}B^T S_B^{-1} = S_B S_B^{-1} = I$$

and therefore they are all equal to 1. Hence, $I - A^{-1/2}B^T S_B^{-1} BA^{-1/2} \succcurlyeq 0$ as claimed. ∎

PROPOSITION 3.7. *Suppose that $A \succ 0$, $B^T, C^T$ have full column rank, and* $\tilde{A} = A^{-1} - A^{-1}B^T S_B^{-1} BA^{-1}$. *If $z^T(C\tilde{A}C^T)z = 0$ for some nonzero vector $z$, then* $\mathrm{range}(B^T) \cap \mathrm{range}(C^T) \neq \{0\}$.
*Proof.* Suppose that $z$ is a nonzero vector such that $z^T(C\tilde{A}C^T)z = 0$. Setting $y = C^T z$ and invoking Proposition 3.6, we obtain that $y^T \tilde{A} y = 0$ where $\tilde{A} \succcurlyeq 0$. Note that $C^T$ has full column rank, hence $y \neq 0$. From [9, Page 400], we obtain that $\tilde{A}y = 0$, or $y \in \ker(\tilde{A})$. On the other hand, $\tilde{A}y = 0$ implies that $y = B^T S_B^{-1} BA^{-1}y$, which shows that $y \in \mathrm{range}(B^T)$. Consequently, in view of the definition of $y$, we have that that $y \in \mathrm{range}(B^T) \cap \mathrm{range}(C^T)$ as claimed. ∎

The following proposition provides a necessary and sufficient condition under which $\rho(\mathscr{G}_2) < 1$.
PROPOSITION 3.8. *Assume that $\mathscr{A}$ is invertible, with $A \succ 0$ and $D = 0$. A necessary and sufficient condition for the iterative scheme* (3.2) *to be convergent is*

$$0 < \alpha < \frac{2}{\lambda_{\max}(\hat{S}_C)},$$

*where $\hat{S}_C = C\tilde{A}C^T$ and $\tilde{A}$ is defined as before. The minimum value of the spectral radius $\rho(\mathscr{G}_2)$ is attained for*

$$\alpha^* = \frac{2}{\lambda_{\max}(\hat{S}_C) + \lambda_{\min}(\hat{S}_C)}.$$

*Proof.* Since $\mathscr{A}$ is assumed to be nonsingular, by Remark 2.4, the matrices $B^T$ and $C^T$ have full column rank and Proposition 2.3 guarantees $\mathrm{range}(B^T) \cap \mathrm{range}(C^T) \neq \{0\}$. Therefore, Proposition 3.7 implies that $\hat{S}_C \succ 0$. From the structure of the matrix $\mathscr{G}_2$, given by (3.4), it is clear that a necessary and sufficient condition for the

convergence of (3.2) is that $\rho(I - \alpha \hat{S}_C) < 1$; moreover, the asymptotic convergence rate is the same as that of Richardson's method for solving a linear system of equations with coefficient matrix $\hat{S}_C$. Now the conclusions follow from the results in [16, Chapter 4] on the convergence properties of Richardson's method applied to a linear system with an SPD coefficient matrix. □

**3.2. Uzawa-like iterative methods of the second type.** In this subsection we focus primarily on the case $D \neq 0$ and present two Uzawa-like iterative schemes for solving (1.1). Nevertheless, we stress that these iterative schemes can also be applied in the case that $D = 0$ (with $\mathscr{A}$ invertible). In what follows we assume that a splitting $D = M - N$ is given.

Similarly to the previous subsection, we consider two splittings for $\mathscr{A}$,

$$\mathscr{A} = \bar{\mathscr{M}}_1 - \bar{\mathscr{N}}_1 = \bar{\mathscr{M}}_2 - \bar{\mathscr{N}}_2,$$

where now

$$\bar{\mathscr{M}}_1 = \begin{bmatrix} A & 0 & 0 \\ B & -\frac{1}{\alpha}I & 0 \\ C & 0 & -M \end{bmatrix}, \quad \bar{\mathscr{N}}_1 = \begin{bmatrix} 0 & -B^T & -C^T \\ 0 & -\frac{1}{\alpha}I & 0 \\ 0 & 0 & -N \end{bmatrix},$$

$$\bar{\mathscr{M}}_2 = \begin{bmatrix} A & B^T & 0 \\ B & 0 & 0 \\ C & 0 & -M \end{bmatrix}, \quad \text{and} \quad \bar{\mathscr{N}}_2 = \begin{bmatrix} 0 & 0 & -C^T \\ 0 & 0 & 0 \\ 0 & 0 & -N \end{bmatrix}.$$

Consequently we may define the following iterative methods for solving (1.1),

$$(3.11) \qquad u_{k+1} = \bar{\mathscr{G}}_1 u_k + \bar{\mathscr{M}}_1^{-1} b, \qquad k = 0, 1, 2 \ldots,$$

and

$$(3.12) \qquad u_{k+1} = \bar{\mathscr{G}}_2 u_k + \bar{\mathscr{M}}_2^{-1} b, \qquad k = 0, 1, 2 \ldots,$$

where $u_0$ is arbitrary,

$$(3.13) \qquad \bar{\mathscr{G}}_1 = \begin{bmatrix} 0 & -A^{-1}B^T & -A^{-1}C^T \\ 0 & I - \alpha BA^{-1}B^T & -\alpha BA^{-1}C^T \\ 0 & -M^{-1}CA^{-1}B^T & M^{-1}(N - CA^{-1}C^T) \end{bmatrix},$$

and

$$(3.14) \qquad \bar{\mathscr{G}}_2 = \begin{bmatrix} 0 & 0 & -\tilde{A}C^T \\ 0 & 0 & -S_B^{-1}BA^{-1}C^T \\ 0 & 0 & M^{-1}(N - C\tilde{A}C^T) \end{bmatrix},$$

where $\tilde{A}$ is defined in (3.4).

We recall next the following theorem, which plays a key role in the convergence analysis of both iterative methods (3.11) and (3.12).

THEOREM 3.9. *[9, Theorem 7.7.3] Let $A$ and $B$ be two $n \times n$ real symmetric matrices such that $A$ is positive definite and $B$ is positive semidefinite. Then $A \succcurlyeq B$*

*if and only if $\rho(A^{-1}B) \leq 1$, and $A \succ B$ if and only if $\rho(A^{-1}B) < 1$.*

THEOREM 3.10. *Let $A \succ 0$ and let $B^T$ have full column rank. Furthermore, assume that $D = M - N$ is a convergent splitting with $M \succ 0$ and $N \succcurlyeq 0$. If $M \succ S_C$, then for*

$$(3.15) \qquad 0 < \alpha < \frac{1 + \lambda_{\min}(M^{-1}N) - \lambda_{\max}(M^{-1}S_C)}{\lambda_{\max}(S_B)}$$

*the iterative scheme* (3.11) *converges to the solution of* (1.1).

*Proof.* From the structure of $\bar{\mathscr{G}}_1$ it is clear that to prove the convergence of iterative method (3.11) we only need to show that the spectral radius of the following matrix is less than one:

$$\mathscr{H} = \begin{bmatrix} I - \alpha B A^{-1} B^T & -\alpha B A^{-1} C^T \\ -M^{-1}CA^{-1}B^T & M^{-1}N - M^{-1}CA^{-1}C^T \end{bmatrix}$$

$$= \begin{bmatrix} I & 0 \\ 0 & M^{-1}N \end{bmatrix} - \begin{bmatrix} \alpha B \\ M^{-1}C \end{bmatrix} A^{-1} \begin{bmatrix} B^T & C^T \end{bmatrix}.$$

It is easy to see that this matrix is similar to the symmetric matrix

$$\begin{bmatrix} I & 0 \\ 0 & M^{-\frac{1}{2}}NM^{-\frac{1}{2}} \end{bmatrix} - \begin{bmatrix} \sqrt{\alpha}B \\ M^{-\frac{1}{2}}C \end{bmatrix} A^{-1} \begin{bmatrix} \sqrt{\alpha}B^T & C^T M^{-\frac{1}{2}} \end{bmatrix}.$$

Using Lemma 3.2 and straightforward computations, we get the bounds

$$(3.16) \qquad \beta_1 \leq \lambda(\mathscr{H}) \leq \beta_2,$$

for the eigenvalues of $\mathscr{H}$, where

$$(3.17) \qquad \beta_1 = \lambda_{\min}(M^{-1}N) - \alpha\lambda_{\max}(S_B) - \lambda_{\max}(M^{-1}S_C)$$

and

$$(3.18) \qquad \beta_2 = 1 - \alpha\lambda_{\min}(S_B) - \lambda_{\min}(M^{-1}S_C).$$

Evidently, the positive definiteness of $S_B$ implies $\beta_2 < 1$. On the other hand, by Theorem 3.9, the assumption $M \succ S_C$ is equivalent to $\lambda_{\max}(M^{-1}S_C) < 1$ and therefore the set of values of $\alpha$ satisfying (3.15) is not empty and for these values of $\alpha$, we have $\beta_1 > -1$. Therefore, $\rho(\mathscr{H}) < 1$. $\square$

REMARK 3.11. In addition to the assumptions of Theorem 3.10, let us assume that

$$(3.19) \qquad \lambda_{\min}(M^{-1}N) \geq \lambda_{\min}(M^{-1}S_C).$$

From the proof of Theorem 3.10, we can conclude that

$$\rho(\bar{\mathscr{G}}_1) < h(\alpha),$$

where $h(\alpha) = \max\{|\beta_1|, |\beta_2|\}$ in which $\beta_1$ and $\beta_2$ are respectively defined by (3.17) and (3.18). Consequently, it can be observed that $\alpha^* = \text{argmin } h(\alpha)$ is given by

$$(3.20) \qquad \alpha^* = \frac{1 - (\lambda_{\max}(M^{-1}S_C) + \lambda_{\min}(M^{-1}S_C) - \lambda_{\min}(M^{-1}N))}{\lambda_{\max}(S_B) + \lambda_{\min}(S_B)}.$$

Note that if $C^T$ is rank deficient or $N$ is singular then $\lambda_{\min}(M^{-1}S_C) = 0$ and $\lambda_{\min}(M^{-1}N) = 0$, respectively. We observe here that the condition (3.19) guarantees that the value of $\alpha^*$ given by (3.20) is positive.

PROPOSITION 3.12. *Assume that $A \succ 0$ and $D \succ 0$. Let the splitting $D = M - N$ be such that $M \succ 0$. If $N - C\tilde{A}C^T \succcurlyeq 0$, then the iterative scheme (3.12) is convergent for any initial guess.*

*Proof.* Notice that by Proposition 3.7, we have $C\tilde{A}C^T \succcurlyeq 0$. The assumptions imply that $M \succ N - C\tilde{A}^T C^T$. Therefore, using Theorem 3.9 we immediately obtain $\rho(\bar{\mathscr{G}}_2) = \rho(M^{-1}(N - C\tilde{A}C^T)) < 1$. □

PROPOSITION 3.13. *Assume that $A \succ 0$, $D \succcurlyeq 0$, $\operatorname{range}(B^T) \cap \operatorname{range}(C^T) = \{0\}$ and $C^T$ has full column rank. Let the splitting $D = M - N$ be such that $M \succ 0$. If $N - C\tilde{A}C^T \succcurlyeq 0$, then the iterative scheme (3.12) is convergent for any initial guess.*

*Proof.* Notice that the assumptions, together with Propositions 2.3 and 3.7, imply that $C\tilde{A}C^T \succ 0$. Therefore $D + C\tilde{A}C^T \succ 0$, which is equivalent to $M \succ N - C\tilde{A}C^T$. The result follows immediatly from Theorem 3.9. □

REMARK 3.14. Consider the case that $D \succ 0$. From the structure of $\bar{\mathscr{G}}_2$ in (3.14), we can see that a suitable choice for the splitting $D = M - N$ is given by

$$M = D + CA^{-1}C^T \quad \text{and} \quad N = CA^{-1}C^T.$$

Notice that the above splitting satisfies the conditions on the splitting $D = M - N$ required in Proposition 3.12. However, both $M$ and $N$ would be dense matrices in general. As a result, this splitting may not be a practical one for large problems in general situations. Nevertheless, this observation suggests that approximations to such choices of $M$ and $N$ may lead to effective preconditioners for Krylov subspace methods.

REMARK 3.15. Consider the following splitting for $D \succ 0$,

$$M = D + \omega I \quad \text{and} \quad N = \omega I,$$

where $\omega$ is a given nonnegative parameter. In view of Proposition 3.12, the iterative method (3.12) is convergent to the exact solution of (1.1) for any initial guess if

$$\omega \geq \lambda_{\max}(C\tilde{A}C^T).$$

Indeed, if $\omega$ satisfies the above inequality, then Lemma 3.2 implies that $\lambda_{\min}(N - C\tilde{A}C^T) \geq 0$, which is equivalent to say that $N - C\tilde{A}C^T \succcurlyeq 0$. Hence, all the assumptions of Proposition 3.12 are satisfied.

REMARK 3.16. Consider the case that $\mathscr{A}$ is nonsingular with $A \succ 0$ and $D = 0$. Let the splitting $D = M - N$ be such that $M = C\tilde{A}C^T$ and $N = C\tilde{A}C^T$. From Proposition 3.13 one may deduce that the iterative scheme (3.12) is convergent for any initial guess. While in general it is not practical to form the matrix $C\tilde{A}C^T$ explicitly, there are cases where this may be possible (see, for instance, the experiments in section 5.1, where good results are reported using this approach).

In the case that $D = M - N \succ 0$ and $N - C\tilde{A}C^T \succcurlyeq 0$, the convergence of iterative method (3.12) has been proved by Proposition 3.12. Next, we study the convergence

of the iterative method in the case that $N - C\tilde{A}C^T \preccurlyeq 0$. It is interesting to note that in this case we can drop the assumption that $D$ is positive semidefinite, as long as $\mathscr{A}$ is nonsingular. We first recall a useful lemma.

LEMMA 3.17. *[17] Suppose that $A$ and $B$ are $n \times n$ Hermitian matrices with $A$ negative definite and $B$ positive semidefinite. Then*

$$\lambda_{\min}(A)\lambda_{\min}(B) \leq \lambda_{\max}(AB) \leq \lambda_{\max}(A)\lambda_{\min}(B),$$

$$\lambda_{\min}(A)\lambda_{\max}(B) \leq \lambda_{\min}(AB) \leq \lambda_{\max}(A)\lambda_{\max}(B).$$

THEOREM 3.18. *Assume $\mathscr{A}$ be nonsingular, with $A \succ 0$, and consider a splitting $D = M - N$ such that $M \succ 0$. If $N - C\tilde{A}C^T \preccurlyeq 0$ and*

(3.21)                    $$\lambda_{\min}(N - C\tilde{A}C^T) + \lambda_{\min}(M) > 0,$$

*then the iterative scheme (3.12) is convergent to the solution of (1.1).*

*Proof.* From Lemma 3.17 we have that

$$\lambda_{\min}(N - C\tilde{A}C^T)\lambda_{\min}(M^{-1}) \leq \lambda_{\max}(\bar{\mathscr{G}}_2) \leq \lambda_{\max}(N - C\tilde{A}C^T)\lambda_{\min}(M^{-1})$$

and

$$\lambda_{\min}(N - C\tilde{A}C^T)\lambda_{\max}(M^{-1}) \leq \lambda_{\min}(\bar{\mathscr{G}}_2) \leq \lambda_{\max}(N - C\tilde{A}C^T)\lambda_{\max}(M^{-1}).$$

In view of the fact that $N - C\tilde{A}C^T$ is negative semidefinite and $M \succ 0$ we may immediately conclude that $\lambda_{\max}(N - C\tilde{A}C^T)\lambda_{\min}(M^{-1}) < 1$ and thus $\lambda_{\max}(\bar{\mathscr{G}}_2) < 1$. On the other hand, condition (3.21) implies that $\lambda_{\min}(N - C\tilde{A}C^T)\lambda_{\max}(M^{-1}) > -1$, which completes the proof. □

REMARK 3.19. In the case that $D \succ 0$, we may simply use the splitting $D = M - N$ where $M = D$ and $N = 0$, provided $D$ is easily invertible and

(3.22)                    $$\lambda_{\max}(C\tilde{A}C^T) < \lambda_{\min}(D),$$

which is equivalent to (3.21). It is not difficult to verify that (3.22) is satisfied if the following relation holds:

(3.23)                    $$\lambda_{\max}(CA^{-1}C^T) < \lambda_{\min}(D).$$

Numerical results for iterative methods (3.12) corresponding to this splitting and the splitting given in Remark 3.14 will be discussed in section 5.

**4. Preconditioning techniques.** In this section we develop and analyze several block preconditioners to be used in conjunction with Krylov subspace methods to solve linear system of equations of the form (1.1). The section is divided into two subsections which correspond again to the two main cases $D = 0$ and $D \neq 0$, respectively.

**4.1. Block preconditioners of the first type.** In this part we discuss the eigenvalue distribution of the preconditioned matrices corresponding to the following block diagonal and block triangular preconditioners for solving systems of the form (1.1) with $D = 0$:

$$\mathscr{P}_{D} = \begin{bmatrix} A & 0 & 0 \\ 0 & BA^{-1}B^T & 0 \\ 0 & 0 & CA^{-1}C^T \end{bmatrix}, \quad \mathscr{P}_{T} = \begin{bmatrix} A & B^T & C^T \\ 0 & -BA^{-1}B^T & 0 \\ 0 & 0 & -CA^{-1}C^T \end{bmatrix},$$

(4.1)
$$\mathscr{P}_{GD} = \begin{bmatrix} A & 0 & 0 \\ 0 & BA^{-1}B^T & BA^{-1}C^T \\ 0 & CA^{-1}B^T & CA^{-1}C^T \end{bmatrix}, \quad \mathscr{P}_{GT,1} = \begin{bmatrix} A & 0 & 0 \\ B & -BA^{-1}B^T & -BA^{-1}C^T \\ C & -CA^{-1}B^T & -CA^{-1}C^T \end{bmatrix},$$

and

(4.2)
$$\mathscr{P}_{GT,2} = \begin{bmatrix} A & B^T & 0 \\ B & 0 & 0 \\ C & 0 & -\bar{S} \end{bmatrix},$$

where

(4.3)
$$\bar{S} = C(A^{-1} + A^{-1}B^T S_B^{-1} BA^{-1})C^T.$$

These preconditioners can be regarded as extensions or generalizations of "standard" block diagonal and block triangular preconditioners for saddle point problems (see, e.g., [4] and [7] for extensive treatments). We note that the two block triangular preconditioners $\mathscr{P}_{GT,1}$ and $\mathscr{P}_{GT,2}$ correspond to the two natural possible partitioning of the matrix $\mathscr{A}$ shown in (1.2). We also remark that all these preconditioners are examples of "ideal" preconditioners, in the sense that in general the matrices $S_B = BA^{-1}B^T$, $S_C = CA^{-1}C^T$, $BA^{-1}C^T$ (or $CA^{-1}B^T$), and $\bar{S}$ will be full and therefore cannot be formed explicitly. In practice, they (or their inverses) will have to be approximated, possibly by some iterative process; the same applies to the action of $A^{-1}$ when solving the systems associated with the preconditioners.[1] Hence, in practice, the preconditioners will have to be applied "inexactly", possibly necessitating the use of a flexible Krylov subspace method. Nevertheless, the spectral analysis for the ideal case is still useful as it provides insight on the performance of the inexact preconditioners, at least for "sufficiently accurate" inexact solves.

We also mention that one can just as well adopt block upper triangular variants of the preconditioners $\mathscr{P}_{GT,1}$ and $\mathscr{P}_{GT,2}$. It has been shown in [14] that the difference between employing block lower and upper preconditioners should not be very significant, with the block upper triangular versions often working slightly better in practice. Nevertheless, in our numerical experiments we opted for $\mathscr{P}_{GT,1}$ instead of the block upper triangular version as the subsystem corresponding to $\mathscr{S}$ (see (3.5)) is solved inexactly by an inner iteration, while the subsystem associated with coefficient matrix $A$ is solved "exactly." Hence, using forward substitution leads to a more accurate application of the preconditioner. For consistency we also chose to adopt the lower triangular form for $\mathscr{P}_{GT,2}$.

---

[1]See section 5.1, however, for an example in which some of these matrices remain sparse and can be formed explicitly.

Our first result concerns the block diagonal preconditioner $\mathscr{P}_D$. It is obvious that $\mathscr{P}_D$ is invertible (indeed, SPD) if and only if $A \succ 0$ and $B$, $C$ have full rank. Under these assumptions, $\mathscr{P}_D$ can be used to precondition the Minimal Residual (MINRES) method [13].

THEOREM 4.1. *Suppose that $A \succ 0$, $B^T$ and $C^T$ have full column rank, and that $D = 0$ in (1.1). Then*

$$(4.4) \qquad \sigma(\mathscr{P}_D^{-1}\mathscr{A}) \subset \left(-1, \frac{1 - \sqrt{1 + 4\gamma_*}}{2}\right) \cup \{1\} \cup \left(\frac{1 + \sqrt{1 + 4\gamma_*}}{2}, 2\right),$$

*with*

$$(4.5) \qquad \gamma_* = \min \frac{x^T(B^T S_B^{-1} B + C^T S_C^{-1} C)x}{x^T A x} > 0,$$

*where the minimum is taken over all $x \in \mathbb{R}^n$, $x \notin \ker(B) \cap \ker(C)$, such that $(x; y; z)$ is an eigenvector of $\mathscr{P}_D^{-1}\mathscr{A}$. In particular, the set $\{1\} \cup \left(\frac{1+\sqrt{1+4\gamma_*}}{2}, 2\right)$ contains $n$ eigenvalues and the negative interval $\left(-1, \frac{1-\sqrt{1+4\gamma_*}}{2}\right)$ contains $m + p$ eigenvalues. Furthermore, if $\lambda \neq 1$ is an eigenvalue of $\mathscr{P}_D^{-1}\mathscr{A}$, then $1 - \lambda$ is also an eigenvalue.*

*Proof.* Since $\mathscr{A}$ is symmetric and $\mathscr{P}_D$ is SPD, all the eigenvalues and corresponding eigenvectors are real. Let $\lambda$ be an arbitrary eigenvalue of $\mathscr{P}_D^{-1}\mathscr{A}_D$, then there exists a vector $(x; y; z) \neq (0; 0; 0)$ such that

$$(4.6) \qquad Ax + B^T y + C^T z = \lambda A x,$$
$$(4.7) \qquad Bx \qquad\qquad = \lambda B A^{-1} B^T y,$$
$$(4.8) \qquad Cx \qquad\qquad = \lambda C A^{-1} C^T z.$$

Note that it must be $x \neq 0$, otherwise $y = 0$ and $z = 0$ by (4.7)-(4.8). If $\ker(B) \cap \ker(C) \neq \{0\}$ then $\lambda = 1$ is an eigenvalue, since any vector $(x; 0; 0)$ with $x \neq 0$, $x \in \ker(B) \cap \ker(C)$ will be a corresponding eigenvector of $\mathscr{A}$. Conversely, any eigenvector corresponding to $\lambda = 1$ is necessarily of this form.

Assume now that $\lambda \neq 1$. We compute $y = \frac{1}{\lambda}(BA^{-1}B^T)^{-1}Bx \equiv \frac{1}{\lambda}S_B^{-1}Bx$ and $z = \frac{1}{\lambda}(CA^{-1}C^T)^{-1}Cx \equiv \frac{1}{\lambda}S_C^{-1}Cx$ from (4.7) and (4.8), respectively. Substituting the computed $y$ and $z$ into (4.6) and premultiplying by $x^T$, we obtain the following quadratic equation:

$$(4.9) \qquad \lambda^2 - \lambda - \gamma = 0,$$

where

$$\gamma = \frac{x^T \left(B^T S_B^{-1} B + C^T S_C^{-1} C\right) x}{x^T A x} > 0.$$

The roots of (4.9) are given by

$$(4.10) \qquad \lambda_+ = \frac{1 + \sqrt{1 + 4\gamma}}{2} \quad \text{and} \quad \lambda_- = \frac{1 - \sqrt{1 + 4\gamma}}{2},$$

which shows that $\lambda_\pm = 1 - \lambda_\mp$. Since

$$\frac{x^T B^T S_B^{-1} Bx}{x^T A x} \leq \lambda_{\max}(A^{-1} B^T S_B^{-1} B) = 1,$$

and, in a similar way,

$$\frac{x^T C^T S_C^{-1} C x}{x^T A x} \leq 1,$$

we obtain that $\gamma \in (0, 2)$ and thus $-1 < \lambda_- < \frac{1 - \sqrt{1 + 4\gamma_*}}{2} < 0$ and $1 < \frac{1 + \sqrt{1 + 4\gamma_*}}{2} < \lambda_+ < 2$, proving (4.4).

Finally, recalling that $\mathscr{A}$ has $n$ positive and $m + p$ negative eigenvalues (see, e.g., [4, Sect. 3.4]) and observing that $\mathscr{P}_D^{-1} \mathscr{A}$ is similar to $\mathscr{P}_D^{-\frac{1}{2}} \mathscr{A} \mathscr{P}_D^{-\frac{1}{2}}$, we conclude by Sylvester's Law of Inertia that there are exactly $n$ eigenvalues that are either 1 or lie in the positive interval in (4.4), and exactly $m + p$ eigenvalues lying in the negative interval, counted with their multiplicities. □

REMARK 4.2. It is clear from the foregoing proof that for any positive eigenvalue of the form $\lambda_+$, there must be a corresponding negative eigenvalue $\lambda_- = 1 - \lambda_+$; see (4.10). On the other hand, we also showed that $\mathscr{P}_D^{-1} \mathscr{A}_D$ must have $n$ positive and $m + p$ negative eigenvalues, and in general $n > m + p$. This is true whether $\lambda = 1$ is an eigenvalue or not. This apparent contradiction can be explained by observing that the multiplicity of $\lambda_+$ as an eigenvalue of $\mathscr{P}_D^{-1} \mathscr{A}_D$ will generally be different from that of the corresponding $\lambda_-$. Indeed, there may be a different number of eigenvectors of the form $(x; y; z)$ corresponding to $\lambda_+$ and to $\lambda_-$, all with the same $x$ (and thus the same $\gamma$) but different $y$ or $z$. Hence, while the negative and positive interval must contain the same number of *distinct* non-unit eigenvalues, the multiplicities of the positive and negative eigenvalues must add up to $n$ and $m + p$, respectively.

REMARK 4.3. While Theorem 4.1 shows that the positive eigenvalues are nicely bounded (between 1 and 2), it does not provide any useful information on the rightmost negative eigenvalue, since $\gamma_*$, while always strictly greater than zero, can in principle be arbitrarily small. Nevertheless, in special cases, given additional assumptions on the blocks $A$, $B$ and $C$, it may be possible to derive a positive lower bound for $\gamma$, and therefore an upper bound (away from zero) for the rightmost negative eigenvalue.

Next, we prove a result concerning the spectrum of of matrices preconditioned with the block triangular preconditioner $\mathscr{P}_T$. We note that since this preconditioner is nonsymmetric, it cannot be used with MINRES. Note that $\mathscr{P}_T$ is guaranteed to be nonsingular when $A \succ 0$ and $B$, $C$ have full rank.

THEOREM 4.4. *Under the assumptions of Theorem 4.1, $\sigma(\mathscr{P}_T^{-1} \mathscr{A}) \subset (0, 2)$, with $\lambda = 1$ being an eigenvalue of multiplicity at least $n$. Moreover, the spectrum of $\mathscr{P}_T^{-1} \mathscr{A}$ is symmetric with respect to $\lambda = 1$; i.e., if $\lambda_1 \neq 1$ and $\lambda_2 \neq 1$ are two eigenvalues of $\mathscr{P}_T^{-1} \mathscr{A}$, then $\lambda_1 + \lambda_2 = 2$.*

*Proof.* Suppose that $\lambda$ is an arbitrary eigenvalue of $\mathscr{P}_T^{-1} \mathscr{A}$ with the corresponding eigenvector $(x; y; z)$, i.e.,

$$\begin{align}
(4.11) \qquad Ax + B^T y + C^T z &= \lambda(Ax + B^T y + C^T z), \\
(4.12) \qquad Bx \phantom{+ B^T y + C^T z} &= -\lambda B A^{-1} B^T y, \\
(4.13) \qquad Cx \phantom{+ B^T y + C^T z} &= -\lambda C A^{-1} C^T z.
\end{align}$$

Notice that $x \neq 0$, otherwise, in view of the fact that $B^T$ and $C^T$ are full column

rank, $x = 0$ implies $(x; y; z) = (0; 0; 0)$ in contradiction with the fact that $(x; y; z)$ is an eigenvector.

Clearly, $\lambda = 1$ is an eigenvalue of $\mathscr{P}_T^{-1}\mathscr{A}$ with corresponding eigenvector of the form $(x; -S_B^{-1}Bx; -S_C^{-1}Cx)$. The multiplicity of this eigenvalue is therefore at least $n$. Assume now that $\lambda \neq 1$. From (4.11), we deduce that

$$(4.14) \qquad\qquad Ax + B^T y + C^T z = 0.$$

Similar to the proof of Theorem 4.1, we compute $y$ and $z$ from (4.12) and (4.13) in terms of $\lambda$ and $x$, respectively. Substituting the derived values of $y$ and $z$ into (4.14), we get

$$(4.15) \qquad\qquad \lambda = \frac{x^* \left( B^T S_B^{-1} Bx + C^T S_C^{-1} C \right) x}{x^* Ax}.$$

(Note that since $\lambda$ is real, the corresponding eigenvector can also be chosen to be real and therefore $x^*$ in (4.15) can be replaced by $x^T$.) Hence, $\lambda$ has the same expression as $\gamma$ in the proof of Theorem 4.1, therefore (4.15) shows that $\lambda \in (0, 2]$.

Next, recall that $\sigma(\mathscr{A}\mathscr{P}_T^{-1}) = \sigma(\mathscr{P}_T^{-1}\mathscr{A})$. Straightforward computations reveal that

$$\mathscr{A}\mathscr{P}_T^{-1} = \begin{bmatrix} A & B^T & C^T \\ B & 0 & 0 \\ C & 0 & 0 \end{bmatrix} \begin{bmatrix} A^{-1} & A^{-1}B^T S_B^{-1} & A^{-1}C^T S_C^{-1} \\ 0 & -S_B^{-1} & 0 \\ 0 & 0 & -S_C^{-1} \end{bmatrix}$$

$$= \begin{bmatrix} I & 0 & 0 \\ BA^{-1} & I & BA^{-1}C^T S_C^{-1} \\ CA^{-1} & CA^{-1}B^T S_B^{-1} & I \end{bmatrix}.$$

The above relation, incidentally, confirms that the number of eigenvalues which are equal to one cannot be less than $n$, the order of the $(1, 1)$-block. In addition, it can be seen that the remaining $m + p$ eigenvalues of $\mathscr{P}_T^{-1}\mathscr{A}$ are the eigenvalues of $I + \hat{\mathscr{S}}$, where

$$\hat{\mathscr{S}} = \begin{bmatrix} 0 & BA^{-1}C^T S_C^{-1} \\ CA^{-1}B^T S_B^{-1} & 0 \end{bmatrix}$$

$$= \begin{bmatrix} 0 & BA^{-1}C^T \\ CA^{-1}B^T & 0 \end{bmatrix} \begin{bmatrix} S_B^{-1} & 0 \\ 0 & S_C^{-1} \end{bmatrix}.$$

To conclude the proof, we only need to show that the distribution of the eigenvalues of $\hat{\mathscr{S}}$ is symmetric with respect to zero. Hence, all the eigenvalues of $\mathscr{P}_T^{-1}\mathscr{A}$ must lie in the interval $(0, 2)$. In view of the fact that $S_B \succ 0$ and $S_C \succ 0$, matrix $\hat{\mathscr{S}}$ is similar to

$$\check{\mathscr{S}} = \begin{bmatrix} S_B^{-1/2} & 0 \\ 0 & S_C^{-1/2} \end{bmatrix} \begin{bmatrix} 0 & BA^{-1}C^T \\ CA^{-1}B^T & 0 \end{bmatrix} \begin{bmatrix} S_B^{-1/2} & 0 \\ 0 & S_C^{-1/2} \end{bmatrix},$$

and therefore the two matrices have the same eigenvalues. Evidently,

$$\check{\mathscr{S}} = \begin{bmatrix} 0 & X \\ X^T & 0 \end{bmatrix},$$

with $X = S_B^{-1/2} B A^{-1} C^T S_C^{-1/2}$. It is well known that the eigenvalues of a matrix of the above form are given by $\pm \sigma_i(X)$, where $\sigma_i(X)$ stands for the $i$th singular value of $X$. This shows the symmetric distribution of the eigenvalues of $\hat{\mathscr{S}}$ with respect to zero. $\square$

REMARK 4.5. Similar to Remark 4.3, we note that without additional assumptions on the matrices $A$, $B$ and $C$ we cannot give a useful lower bound on the eigenvalues of $\mathscr{P}_T^{-1} \mathscr{A}$.

We conclude this section with a few brief remarks on the preconditioners $\mathscr{P}_{GD}$, $\mathscr{P}_{GT,1}$, and $\mathscr{P}_{GT,2}$. We observe that the first two are just special cases of the "ideal" block diagonal and block (lower) triangular preconditioners for saddle point problems based on the first of the two partitionings in (1.2); the third one is the ideal block (lower) triangular preconditioner based on the second partitioning of $\mathscr{A}$ in (1.2). The spectral properties of preconditioned saddle point matrices with any of these block preconditioners are well known; see, e.g., [4, Sec. 10.1.1–10.1.2]. In particular, $\mathscr{P}_{GD}^{-1} \mathscr{A}$ has only three distinct eigenvalues and is diagonalizable, while $\mathscr{P}_{GT,1}^{-1} \mathscr{A}$ and $\mathscr{P}_{GT,2}^{-1} \mathscr{A}$ have all the eigenvalues equal to 1 and are non-diagonalizable but have minimum polynomial of degree 2. Hence, MINRES and the Generalized Minimum Residual Method (GMRES) [16] will reach the exact solution in at most three and two steps, respectively. As before, these ideal block preconditioners may be prohibitively expensive to construct and apply; in practice, they are ususally replaced by inexact variants.

**4.2. Block preconditioners of the second type.** In this part the eigenvalue distributions of the preconditioned matrices are discussed for the case that the coefficient matrix $\mathscr{A}$ has nonzero $(3,3)$-block. We consider two following types of block triangular preconditioners:

$$(4.16) \qquad \tilde{\mathscr{P}}_T = \begin{bmatrix} A & B^T & C^T \\ 0 & -BA^{-1}B^T & 0 \\ 0 & 0 & -(D + CA^{-1}C^T) \end{bmatrix},$$

and

$$(4.17) \qquad \hat{\mathscr{P}}_T = \begin{bmatrix} A & B^T & C^T \\ 0 & -BA^{-1}B^T & -BA^{-1}C^T \\ 0 & 0 & -(D + CA^{-1}C^T) \end{bmatrix}.$$

We note that these preconditioners wil be nonsingular if $A \succ 0$, $B^T$ has full column rank, $D \succeq 0$ and $\ker(D) \cap \ker(C^T) = \{0\}$. From Theorem 2.2, these conditions also guarantee the invertibility of $\mathscr{A}$.

For ease of exposition, we present the analysis in several steps. Our first result is the following.

THEOREM 4.6. *Assume that $A \succ 0$, $B$ has full rank, $D \succeq 0$ and $\ker(D) \cap \ker(C^T) = \{0\}$. Then all the eigenvalues of $\mathscr{A} \tilde{\mathscr{P}}_T^{-1}$ are real and nonzero. Moreover, $\lambda = 1$ is an eigenvalue of algebraic multiplicity at least $n$.*

*Proof.* Under the stated assumptions, both $\mathscr{A}$ and $\tilde{\mathscr{P}}_T$ are nonsigular. We have

$$\mathscr{A}\,\tilde{\mathscr{P}}_T^{-1} = \begin{bmatrix} A & B^T & C^T \\ B & 0 & 0 \\ C & 0 & -D \end{bmatrix} \begin{bmatrix} A^{-1} & A^{-1}B^T S_B^{-1} & A^{-1}C^T \tilde{S}_C^{-1} \\ 0 & -S_B^{-1} & 0 \\ 0 & 0 & -\tilde{S}_C^{-1} \end{bmatrix}$$

$$= \begin{bmatrix} I & 0 & 0 \\ BA^{-1} & I & BA^{-1}C^T \tilde{S}_C^{-1} \\ CA^{-1} & CA^{-1}B^T S_B^{-1} & I + D\tilde{S}_C^{-1} \end{bmatrix},$$

where $\tilde{S}_C = D + CA^{-1}C^T$. Similar to the proof of Theorem 4.4, we find that the number of eigenvalues of $\mathscr{A}\,\tilde{\mathscr{P}}_T^{-1}$ which are equal to one is at least $n$, the order of the $(1,1)$-block, with the remaining eigenvalues being those of the matrix $I + \tilde{\mathscr{S}}_1$ where

$$\tilde{\mathscr{S}}_1 = \begin{bmatrix} 0 & BA^{-1}C^T \tilde{S}_C^{-1} \\ CA^{-1}B^T S_B^{-1} & D\tilde{S}_C^{-1} \end{bmatrix}$$

$$= \begin{bmatrix} 0 & BA^{-1}C^T \\ CA^{-1}B^T & D \end{bmatrix} \begin{bmatrix} S_B^{-1} & 0 \\ 0 & \tilde{S}_C^{-1} \end{bmatrix}.$$

Since $\tilde{\mathscr{S}}_1$ is the product of two symmetric matrices, one of which is positive definite, its eigenvalues are all real and the result is proved. $\square$

Next, we present bounds on the eigenvalues of the preconditioned matrices $\tilde{\mathscr{P}}_T^{-1}\mathscr{A}$ and $\hat{\mathscr{P}}_T^{-1}\mathscr{A}$. To this end, we make use of the Cholesky factorization of the $(1,1)$-block of $\mathscr{A}$, i.e., $A = LL^T$. Consider the lower triangular matrix $\mathscr{L}$ defined by

$$\mathscr{L} = \begin{bmatrix} L & 0 & 0 \\ 0 & I & 0 \\ 0 & 0 & I \end{bmatrix}.$$

We define $\hat{\mathscr{A}} = \mathscr{L}^{-1}\mathscr{A}\mathscr{L}^{-T}$, which has the following structure:

$$\hat{\mathscr{A}} = \begin{bmatrix} I & \hat{B}^T & \hat{C}^T \\ \hat{B} & 0 & 0 \\ \hat{C} & 0 & D \end{bmatrix},$$

where $\hat{B} = BL^{-T}$ and $\hat{C} = CL^{-T}$. Now we consider the following two block triangular preconditioners for $\hat{\mathscr{A}}$,

$$(4.18) \qquad \tilde{\hat{\mathscr{P}}}_T = \begin{bmatrix} I & \hat{B}^T & \hat{C}^T \\ 0 & -\hat{B}\hat{B}^T & 0 \\ 0 & 0 & -(D + \hat{C}\hat{C}^T) \end{bmatrix}$$

and

$$(4.19) \qquad \hat{\hat{\mathscr{P}}}_T = \begin{bmatrix} I & \hat{B}^T & \hat{C}^T \\ 0 & -\hat{B}\hat{B}^T & -\hat{B}\hat{C}^T \\ 0 & 0 & -(D + \hat{C}\hat{C}^T) \end{bmatrix}.$$

It is not difficult to check that the following two relations hold:

$$\tilde{\mathscr{P}}_T^{-1}\mathscr{A} = \mathscr{L}^{-T}\tilde{\hat{\mathscr{P}}}_T^{-1}\hat{\mathscr{A}}\mathscr{L}^T \quad \text{and} \quad \hat{\mathscr{P}}_T^{-1}\mathscr{A} = \mathscr{L}^{-T}\hat{\hat{\mathscr{P}}}_T^{-1}\hat{\mathscr{A}}\mathscr{L}^T,$$

which reveal that $\sigma(\tilde{\mathscr{P}}_{T}^{-1}\mathscr{A}) = \sigma(\tilde{\hat{\mathscr{P}}}_{T}^{-1}\hat{\mathscr{A}})$ and $\sigma(\hat{\mathscr{P}}_{T}^{-1}\mathscr{A}) = \sigma(\hat{\hat{\mathscr{P}}}_{T}^{-1}\hat{\mathscr{A}})$.

THEOREM 4.7. *Under the same assumptions of Theorem 4.6*, $\sigma(\tilde{\mathscr{P}}_{T}^{-1}\mathscr{A}) = \sigma(\tilde{\hat{\mathscr{P}}}_{T}^{-1}\hat{\mathscr{A}}) \subset (0, 1 - \sqrt{\xi}] \cup \{1\} \cup [1 + \sqrt{\xi}, 2) \subset (0, 2)$, *where*

$$(4.20) \qquad \xi = \frac{\bar{\sigma}_{\min}^2(\hat{C})}{\lambda_{\max}(D) + \bar{\sigma}_{\min}^2(\hat{C})}.$$

*Here $\bar{\sigma}_{\min}(\hat{C})$ denotes the smallest nonzero singular value of $\hat{C}$.*

*Proof.* The equality $\sigma(\tilde{\mathscr{P}}_{T}^{-1}\mathscr{A}) = \sigma(\tilde{\hat{\mathscr{P}}}_{T}^{-1}\hat{\mathscr{A}})$ has already been noted. From Theorem 4.6 we already know that the spectrum is real and that $\lambda = 1$ is an eigenvalue of algebraic multiplicity at least $n$. Assume now that $\lambda \neq 1$ is an eigenvalue of $\tilde{\hat{\mathscr{P}}}_{T}^{-1}\hat{\mathscr{A}}$. There exists a (real) nonzero vector $(x; y; z)$ such that

$$(4.21) \qquad x + \hat{B}^T y + \hat{C}^T z = \lambda(x + \hat{B}^T y + \hat{C}^T z),$$
$$(4.22) \qquad \hat{B}x = -\lambda \hat{B}\hat{B}^T y,$$
$$(4.23) \qquad \hat{C}x - Dz = -\lambda(D + \hat{C}\hat{C}^T)z.$$

Notice that $x \neq 0$, otherwise $x = 0$ implies that $y$ and $z$ are both zero in contradiction with the fact that $(x; y; z)$ is an eigenvector.

From (4.21), we get

$$x + \hat{B}^T y + \hat{C}^T z = 0,$$

and therefore

$$\hat{B}x = -(\hat{B}\hat{B}^T y + \hat{B}\hat{C}^T z)$$

and

$$\hat{C}x = -(\hat{C}\hat{B}^T y + \hat{C}\hat{C}^T z).$$

Substituting the preceding two relations into (4.22) and (4.23), respectively, we get

$$(4.24) \qquad (\lambda - 1)\hat{B}\hat{B}^T y = \hat{B}\hat{C}^T z$$

and

$$(4.25) \qquad (\lambda - 1)(D + \hat{C}\hat{C}^T)z = \hat{C}\hat{B}^T y.$$

We observe that the vectors $y$ and $z$ must both be nonzero. Indeed, our assumptions imply that both $\hat{B}\hat{B}^T$ and $D + \hat{C}\hat{C}^T$ are positive definite, and this fact, together with (4.24) and (4.25), implies that $y = 0$ if and only if $z = 0$. Notice that $\hat{C}^T z \neq 0$, otherwise (4.24) implies that $\lambda = 1$ which is contrary to our assumption. By computing $y$ from (4.24) and then substituting it into (4.25), we obtain

$$(4.26) \qquad (\lambda - 1)^2 = \frac{z^T \hat{C} P \hat{C}^T z}{z^T (D + \hat{C}\hat{C}^T)z},$$

where $P = \hat{B}^T(\hat{B}\hat{B}^T)^{-1}\hat{B}$. Note that $P$ is an orthogonal projector, i.e., $P^2 = P$ and $P = P^T$. Using the fact that $\|Pv\|_2 \leq \|v\|_2$ for any vector $v$, we obtain as a consequence of (4.26) that $|\lambda - 1| < 1$, which is equivalent to say that $\lambda \in (0, 2)$.

Finally, we apply the Rayleigh–Ritz Theorem [9, Thm. 4.2.2] to obtain

$$(4.27) \quad \frac{1}{\lambda_{\max}(D)/\bar{\lambda}_{\min}(\hat{C}\hat{C}^T) + 1} \leq \frac{z^* \hat{C} P \hat{C}^T z}{z^*(D + \hat{C}\hat{C}^T)z} \leq \frac{1}{\lambda_{\min}(D)/\lambda_{\max}(\hat{C}\hat{C}^T) + 1},$$

where $\bar{\lambda}_{\min}(\hat{C}\hat{C}^T)$ denotes the smallest nonzero eigenvalue of $\hat{C}\hat{C}^T$. This shows that $|\lambda - 1| \geq \sqrt{\xi}$. The proof is complete. □

We conclude this section with a result on the preconditioner $\hat{\mathscr{P}}_{T}$.

THEOREM 4.8. *Assume that $A \succ 0$, $B^T$ has full column rank, and $D \succ 0$ in (1.1). Then $\sigma(\hat{\mathscr{P}}_{T}^{-1}\mathscr{A}) = \sigma(\hat{\hat{\mathscr{P}}}_{T}^{-1}\hat{\mathscr{A}}) \subseteq \{1\} \cup [1+\xi, 1+\eta] \subset [1,2)$, where $\xi$ is given by (4.20) and*

$$\eta = \frac{\sigma_{\max}^2(\hat{C})}{\lambda_{\min}(D) + \sigma_{\max}^2(\hat{C})}.$$

*Proof.* First, we note that $\hat{\mathscr{P}}_{T}^{-1}\mathscr{A}$ is invertible and $\sigma(\hat{\mathscr{P}}_{T}^{-1}\mathscr{A}) = \sigma(\hat{\hat{\mathscr{P}}}_{T}^{-1}\hat{\mathscr{A}})$, where $\hat{\hat{\mathscr{P}}}_{T}$ is given in (4.19). Let $\lambda$ be an eigenvalue of $\hat{\hat{\mathscr{P}}}_{T}^{-1}\hat{\mathscr{A}}$ with corresponding eigenvector $(x; y; z)$. We have

$$(4.28) \qquad x + \hat{B}^T y + \hat{C}^T z = \lambda(x + \hat{B}^T y + \hat{C}^T z),$$
$$(4.29) \qquad \hat{B}x \qquad\qquad = -\lambda(\hat{B}\hat{B}^T y + \hat{B}\hat{C}^T z),$$
$$(4.30) \qquad \hat{C}x \qquad - Dz = -\lambda(D + \hat{C}\hat{C}^T)z.$$

If $C^T$ (and therefore $\hat{C}^T$) does not have full column rank, we observe that $\lambda = 1$ is an eigenvalue with corresponding eigenvectors of the form $(0; 0; z)$, where $0 \neq z \in \ker(C^T)$. Hence, the multiplicity of $\lambda = 1$ is at least equal to $p - r$, where $r = \operatorname{rank}(C)$.

Let us now assume that $\hat{C}^T$ has full column rank, and let $x \in \mathbb{R}^n$ be any nonzero vector. It is then easy to see that $\lambda = 1$ is an eigenvalue of $\hat{\hat{\mathscr{P}}}_{T}^{-1}\hat{\mathscr{A}}$ with corresponding eigenvector

$$\left(x; -(\hat{B}\hat{B}^T)^{-1}(\hat{B} - \hat{B}\hat{C}^T(\hat{C}\hat{C}^T)^{-1}\hat{C})x; -(\hat{C}\hat{C}^T)^{-1}\hat{C}x\right).$$

Since there are $n$ linearly independent vectors of this form, $\lambda = 1$ is an eigenvalue of multiplicity at least $n$ of $\hat{\mathscr{P}}_{T}^{-1}\mathscr{A}$.

In the sequel we assume that $\lambda \neq 1$. From (4.28) we obtain

$$x + \hat{B}^T y + \hat{C}^T z = 0.$$

It follows that

$$(4.31) \qquad\qquad \hat{B}x = -(\hat{B}\hat{B}^T y + \hat{B}\hat{C}^T z),$$
$$(4.32) \qquad\qquad \hat{C}x = -(\hat{C}\hat{B}^T y + \hat{C}\hat{C}^T z).$$

Substituting (4.31) and (4.32) into (4.29) and (4.30), respectively, we get

$$(4.33) \qquad\qquad (\lambda - 1)(\hat{B}\hat{B}^T y + \hat{B}\hat{C}^T z) = 0,$$
$$(4.34) \qquad\qquad (\lambda - 1)(D + \hat{C}\hat{C}^T)z = -\hat{C}\hat{B}^T y.$$

From the above two relations it can be deduced that $y = 0$ if (and only if) $z = 0$, in which case $x = 0$ in contradiction with the assumption that $(x; y; z)$ is an eigenvector. Keeping in mind that $\lambda \neq 1$, the vector $y$ can be computed from (4.33) as $y = -(\hat{B}\hat{B}^T)^{-1}\hat{B}\hat{C}^T z$. In order to complete the proof, we first substitute $y$ in (4.34), and then multiply both sides of the resulting relation by $z^*$; note that we can actually use $z^T$ since the eigenvalues are necessarily real. Thus,

$$\lambda = 1 + \frac{z^T \hat{C} P \hat{C}^T z}{z^T (D + \hat{C}\hat{C}^T) z},$$

where $P = \hat{B}^T (\hat{B}\hat{B}^T)^{-1} \hat{B}$. As pointed before, the matrix $P$ is an orthogonal projector. The result immediately follows from (4.27). □

REMARK 4.9. We remark again that specific knowledge of the largest and smallest (nonzero) singular value of $\hat{C}$ (for instance, knowledge of their behavior as the meshsize $h \to 0$ in PDE problems) is required in order to make the foregoing spectral bounds explicit and useful. Also, it is well known that eigenvalue information alone does not suffice, in general, to predict the convergence behavior of nonsymmetric Krylov subspace methods like GMRES. Nevertheless, experience shows that in many cases of practical interest convergence can be expected to be fast when the spectrum is real, positive, and contained in an interval of modest length bounded away from zero. This behavior is also observed when the "ideal" preconditioners are replaced with inexact versions, as long as the preconditoner is applied with a reasonable degree of accuracy.

**5. Numerical experiments.** In this section, we present a selection of numerical tests aimed at illustrating the performance of some of the proposed solvers and preconditioners. Due to space limitations, we present detailed results only for some of the methods analyzed in the theoretical sections, and comment briefly on the remaining ones. We focus on two sets of problems of the type (1.1) arising from two very different applications, one with $D = 0$ and the other with $D \neq 0$. All of the reported numerical results were performed on a 64-bit 2.45 GHz core i7 processor and 8.00GB RAM using MATLAB version 8.3.0532. In all of the experiments we have used right-hand sides corresponding to random solution vectors, performing ten runs and then averaging the CPU-times. The iteration counts reported in the tables (under "Iter") are also averages (rounded to the nearest integer).

All of the methods require repeated solution (whether "exact" or inexact) of SPD linear systems as subtasks. These are either solved by sparse Cholesky factorization with symmetric approximate minimum degree (SYMAMD) reordering or by the preconditioned conjugate gradient (PCG) method. When using PCG, unless otherwise specified, the preconditioner used is a drop tolerance-based incomplete Cholesky factorization [3, 16] computed using the MATLAB function "ichol(.,opts)", where
  - opts.type = 'ict',
  - opts.droptol = 1e-2.

In all of the numerical tests below, the initial guess is taken to be the zero vector. For the Uzawa, MINRES, GMRES, and Flexible GMRES (FGMRES) methods the iterations are stopped once

$$\|b - \mathscr{A}(x^{(k)}; y^{(k)}; z^{(k)})\|_2 < 10^{-10} \|b\|_2.$$

For the inner PCG iterations (whenever applicable), the stopping tolerances used are specified below.

**5.1. Saddle point systems from potential fluid flow modeling.** Here we consider linear systems of equations of the form

$$
(5.1) \qquad
\begin{bmatrix} A & B^T & C^T \\ B & 0 & 0 \\ C & 0 & 0 \end{bmatrix}
\begin{bmatrix} x \\ y \\ z \end{bmatrix}
=
\begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix},
$$

arising from a low-order Raviart–Thomas mixed-hybrid finite element approximation [5] of Darcy's law and continuity equation describing the three-dimensional (3D) potential fluid flow problem in porous media. The continuous problem reads:

$$
\mathbf{A}\mathbf{u} = -\nabla p, \quad \nabla \cdot \mathbf{u} = q,
$$

where $\mathbf{u}$ is the fluid velocity, $p$ is the piezometric potential (fluid pressure), $\mathbf{A}$ is the symmetric and uniformly positive definite second-rank tensor of the hydraulic resistance of the medium with $[\mathbf{A}(\mathbf{x})]_{ij} \in L^\infty(\Omega)$ for $i, j = 1, 2, 3$, and $q$ represents the density of potential sources in the medium. The underlying spatial domain $\Omega$ is cubic, and the boundary conditions are given by

$$
p = p_D \quad \text{on} \quad \partial\Omega_D, \qquad \mathbf{u} \cdot \mathbf{n} = \mathbf{u}_N \quad \text{on} \quad \partial\Omega_N,
$$

where $\partial\Omega = \overline{\partial\Omega_D} \cup \overline{\partial\Omega_N}$ with $\partial\Omega_D \neq \emptyset$, $\partial\Omega_D \cap \partial\Omega_N = \emptyset$, and $\mathbf{n}$ is the outward normal vector defined (a.e.) on $\partial\Omega$. We refer to [10] for details of the problem and its discretization. The solution vectors $x$ and $y$ in (5.1) correspond to velocity and pressure degrees of freedom (respectively), while $z$ is a vector of Lagrange multipliers. For this problem we have that $A \succ 0$ and $B^T$, $C^T$ have full column rank, hence $\mathscr{A}$ is nonsingular. Details on the dimensions of sub-blocks $A$, $B$, and $C$ and further information can be found in [10, Table 1].

For this test problem, the SPD matrix $A$ is block diagonal with small blocks, and linear systems associated with it can be solved very cheaply by means of Cholesky factorization. Likewise, the Schur complements $S_B = BA^{-1}B^T$, $S_C = CA^{-1}C^T$, $\bar{S}$ (see (4.3)) and the matrix $BA^{-1}C^T$ are still relatively sparse matrices which can be formed explicitly at low expense.[2] For this problem, the Uzawa-type methods of the first kind were found to converge rather slowly, hence we do not report the results. On the other hand, as pointed out in Remark 3.16, the Uzawa-type method (3.12) converges fast. Concerning the block preconditioners, the best results were obtained with $\mathscr{P}_{GD}$ and $\mathscr{P}_{GT,1}$ in (4.1) and $\mathscr{P}_{GT,2}$ in (4.2). The block diagonal preconditioner $\mathscr{P}_{GD}$ was used with MINRES, while the two block triangular preconditioners $\mathscr{P}_{GT,1}$ and $\mathscr{P}_{GT,2}$ were used with both GMRES and FGMRES.

Apart from the inexpensive solves associated with $A$, the implementation of $\mathscr{P}_{GD}$ and $\mathscr{P}_{GT,1}$ requires solving linear systems associated with matrix $\mathscr{S}$ given in (3.5). In spite of the sparsity of $\mathscr{S}$, solution by sparse Cholesky factorization is expensive (recall that this is a 3D problem). Thus, we solve such systems with the PCG method with a very stringent stopping criterion (inner relative residual norm less than $tol = 10^{-15}$) for MINRES and GMRES and a looser one ($tol = 10^{-4}$) for FGMRES.

---

[2] The Schur complement $BA^{-1}B^T$ for this problem turns out to be a scalar multiple of the $m \times m$ identity matrix.

TABLE 1
*Numerical results for Uzawa's method (3.12), potential fluid flow problem.*

| size | $M = N = C\tilde{A}C^T$ | |
| --- | --- | --- |
| | Iter | CPU |
| 2125 | 2 | 0.0070 |
| 17000 | 2 | 0.0369 |
| 57375 | 2 | 0.1399 |
| 136000 | 2 | 0.4843 |
| 265625 | 2 | 1.7452 |
| 459000 | 2 | 4.1678 |

The application of the preconditioner $\mathscr{P}_{GT,2}$, on the other hand, requires solving at each step a linear system of the form $\mathscr{P}_{GT,2}(w_1; w_2; w_3) = (r_1; r_2; r_3)$. This amounts to solving a saddle point problem of size $(n + m) \times (n + m)$ of the form

$$(5.2) \qquad \begin{bmatrix} A & B^T \\ B & 0 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} r_1 \\ r_2 \end{bmatrix},$$

followed by solution of a linear systems with the coefficient matrix $\bar{S}$ (see (4.3)). The solution of (5.2) can be obtained in two steps as follows:

- **Step I.** Solve $S_B w_2 = BA^{-1}r_1 - r_2$, to find $w_2$.
- **Step II.** Set $w_1 = A^{-1}(r_1 - B^T w_2)$.

We recall that for this particular test problem, $A$ is block diagonal (with small blocks) and $S_B$ is just a scalar multiple of the identity, so the above solution process is extremely cheap and in our experiments we use it both within stationary iterative scheme (3.12) (with $M = N = C\tilde{A}C^T$), and within GMRES and FGMRES iterative methods. We stress that for this problem the matrix $C\tilde{A}C^T$ is sparse and it can be formed explicitly very cheaply. We also observed that $C\tilde{A}C^T$ is well-conditioned. Ib our tests, the linear systems with coefficient matrix $C\tilde{A}C^T$ which arise in applying iterative scheme (3.12) have been solved by CG with inner tolerance $10^{-12}$. In addition to solving (5.2), for applying $\mathscr{P}_{GT,2}$ we also need to solve $\bar{S}w_3 = -r_3 + Cw_1$ where $\bar{S}$ is defined by (4.3). As already mentioned, in this problem $\bar{S}$ can be formed explicitly as it is a sparse matrix. To solve $\bar{S}w_3 = -r_3 + Cw_1$, the PCG method was used where the inner stopping tolerances were chosen as before as $10^{-15}$ and $10^{-4}$ depending on whether GMRES or FGMRES is used, respectively.

In Tables 1, 2, 3, and 4 we report the results for Uzawa's method (3.12) and for the preconditioned MINRES, GMRES and FGMRES iterative methods. The total number $n + m + p$ of unknowns is reported under "size". As expected, MINRES/GMRES with the "ideal" bock diagonal/triangular preconditioners require exactly three and two steps to converge, independent of problem size. In Table 4, the cumulative number of inner PCG iterations required is reported under "Iter$_{pcg}$".

These results show that for this particular example, the best results are obtained with Uzawa's method (3.12) and the inexact block triangular preconditioners $\mathscr{P}_{GT,1}$ and $\mathscr{P}_{GT,2}$; of these last two, the latter one (based on the second of the two partitionings (1.2)) appears to be slighlty better in this particular case. We note the satisfactory scaling in terms of CPU time for sufficiently small $h$, especially for FGMRES with the inexact $\mathscr{P}_{GT,2}$ preconditioner. As for the other two preconditioners, $\mathscr{P}_D$

Table 2
*Results for MINRES with block diagonal preconditioner $\mathscr{P}_{GD}$, potential fluid flow problem.*

| size | Iter | CPU-time |
|------|------|----------|
| 2125 | 3 | 0.0125 |
| 17000 | 3 | 0.0947 |
| 57375 | 3 | 0.4829 |
| 136000 | 3 | 1.6226 |
| 265625 | 3 | 3.9002 |
| 459000 | 3 | 8.8899 |

Table 3
*Results for GMRES with block triangular preconditioners, potential fluid flow problem.*

| | Preconditioner | | | |
|------|------|------|------|------|
| | $\mathscr{P}_{GT,1}$ | | $\mathscr{P}_{GT,2}$ | |
| size | Iter | CPU-time | Iter | CPU-time |
| 2125 | 2 | 0.0191 | 2 | 0.0180 |
| 17000 | 2 | 0.1284 | 2 | 0.1180 |
| 57375 | 2 | 0.5247 | 2 | 0.4516 |
| 136000 | 2 | 1.5425 | 2 | 1.2936 |
| 265625 | 2 | 3.6811 | 2 | 3.1080 |
| 459000 | 2 | 7.9861 | 2 | 6.8368 |

Table 4
*Results for FGMRES with block triangular preconditioners, potential fluid flow problem.*

| | Preconditioner | | | | |
|------|------|------|------|------|------|
| | $\mathscr{P}_{GT,1}$ | | | $\mathscr{P}_{GT,2}$ | |
| size | Iter | $\text{Iter}_{pcg}$ | CPU-time | Iter | $\text{Iter}_{pcg}$ | CPU-time |
| 2125 | 5 | 25 | 0.0085 | 5 | 25 | 0.0073 |
| 17000 | 6 | 47 | 0.0575 | 6 | 53 | 0.0534 |
| 57375 | 6 | 66 | 0.2361 | 6 | 72 | 0.2265 |
| 136000 | 6 | 87 | 0.7480 | 6 | 95 | 0.6563 |
| 265625 | 6 | 108 | 1.8190 | 6 | 112 | 1.5220 |
| 459000 | 6 | 134 | 4.2658 | 5 | 117 | 3.0442 |

and $\mathscr{P}_T$, their performance was generally inferior, with worsening iteration counts for increasing problem sizes. The observed behavior appears to be due to the fact that for this problem, some of the eigenvalues of the preconditioned matrices corresponding to $\mathscr{P}_D$ and $\mathscr{P}_T$ approach zero as the mesh is refined. Still, these preconditioners, as well as the two Uzawa methods discussed in section 3.1, may well be useful in solving saddle point systems arising from other applications.

**5.2. Saddle point systems from liquid crystal directors modeling.** Continuum models for the orientational properties of liquid crystals require the minimization of free energy functionals of the form

$$(5.3) \qquad \mathscr{F}[u,v,w,U] = \frac{1}{2} \int_0^1 \left[ (u_z^2 + v_z^2 + w_z^2) - \alpha^2(\beta + w^2)U_z^2 \right] dz,$$

where $u, v, w$ and $U$ are functions of $z \in [0, 1]$ subject to suitable end-point conditions, $u_z = \frac{du}{dz}$ (etc.), and $\alpha, \beta$ are positive prescribed parameters. Approximation via a uniform piecewise-linear finite element scheme with $k+1$ cells using nodal quadrature and using the prescribed boundary conditions leads to replacing the functional $\mathscr{F}$ with a function $f$ of $4k$ variables:

$$\mathscr{F}[u, v, w, U] \approx f(u_1, \ldots, u_k, v_1, \ldots, v_k, w_1, \ldots, w_k, U_1, \ldots, U_k),$$

see [15, Eq. (2.4)] for the precise form of $f$.

Minimization of the free energy (5.3) must be carried out under the so-called *unit vector constraint*, which at the discrete level can be expresses by imposing that the solution components $u_j$, $v_j$ and $w_j$ satisfy

$$u_j^2 + v_j^2 + w_j^2 = 1, \quad j = 1, \ldots, k.$$

Introducing Lagrange multipliers $\lambda_1, \ldots, \lambda_k$, the problem reduces to finding the critical points of the Lagrangian function

$$L = f + \frac{1}{2} \sum_{j=1}^{k} \lambda_j (u_j^2 + v_j^2 + w_j^2 - 1).$$

Imposing the first-order conditions results in the system of $5k$ nonlinear equations $\nabla L(\mathbf{x}) = \mathbf{0}$ where the unknown vector $\mathbf{x} \in \mathbb{R}^{5k}$ collects the values $(u_j, v_j, w_j)$ ($j = 1, \ldots, k$), $(\lambda_1, \ldots, \lambda_k)$, and $(U_1, \ldots, U_k)$ (in this order). Solving this nonlinear system with Newton's method leads to a linear system of the form

$$(5.4) \qquad \nabla^2 L(\mathbf{x}^{(\ell)}) \, \delta \mathbf{x}^{(\ell)} = -\nabla L(\mathbf{x}^{(\ell)})$$

at each step $\ell$, where $\nabla^2 L(\mathbf{x}^{(\ell)})$ denotes the Hessian of $L$ evaluated at $\mathbf{x}^{(\ell)}$. As shown in [15], the Hessian has the following structure:

$$\nabla^2 L = \begin{bmatrix} A & B^T & C^T \\ B & 0 & 0 \\ C & 0 & -D \end{bmatrix},$$

where $A$ is $n \times n$, $B$ is $m \times n$, $C$ is $p \times n$ and $D \neq 0$ is $p \times p$ with $n = 3k$ and $m = p = k$. Therefore, it is necessary to solve a system of the form (1.1) within each Newton step. Details on the structure of the blocks $A$, $B$, $C$ and $D$ can be found in [15]. Here we note that $A$ is SPD, $B^T$ has full column rank and is such that $BB^T$ is diagonal (and indeed $BB^T = I_m$ if the unit vector constraints are satisfied exactly), $C$ is rank deficient, and $D$ is tridiagonal and SPD. Hence, $\mathscr{A}$ is nonsingular.[3] We also mention that in our experiments we used the following values of the parameters $\alpha$ and $\beta$ appearing in (5.3): $\alpha = 0.5\alpha_c$ and $\beta = 0.5$ where $\alpha_c \approx 2.721$ is known as the *critical switching value*. For further details we refer the reader to [15].

First we consider the use of the Uzawa-type methods (3.11) and (3.12). While the first of these (which necessitates the selection of the parameter $\alpha$) was found to converge very slowly, method (3.12) converged extremely fast for suitable (parameter-free) choices of the splitting $D = M - N$. At each step of (3.12) we need to solve a

---

[3]We are assuming here that the Hessian is being evaluated away from bifurcation points and turning points, see again [15].

linear system of equations with coefficient matrix

$$\bar{\mathscr{M}}_2 = \left[\begin{array}{cc|c} A & B^T & 0 \\ B & 0 & 0 \\ \hline C & 0 & -M \end{array}\right].$$

This requires solving at each step of (3.12) a system of the form (5.2), followed by the solution of a system with coefficient matrix $M$. The first task requires similar steps to those (Steps I and II) described in Example 5.1. Within these steps, we used sparse Cholesky factorization with SYMAMD reordering for solving systems with coefficient matrix $A$, and the PCG method to solve the linear systems with coefficient matrix $S_B = BA^{-1}B^T$. Note that $S_B$ would be full and is not formed explicitly. As for the preconditioner used, the observation that $B$ has (nearly) orthogonal rows suggests that the matrix $BAB^T$ would be a good approximate inverse of $BA^{-1}B^T$, and indeed it was found to be a very effective preconditioner. Note that only sparse matrix-vector products are required for its application, and there is no construction overhead.

Concerning the choice of $M$, we considered the following two options:

$$(5.5) \qquad\qquad M = D + CA^{-1}C^T \quad \text{and} \quad M = D.$$

We note that in view of Remark 3.14, the first choice results in a convergent iteration. For $M = D$, we were able to check for some of the smaller problem sizes that condition (3.23) was satisfied. We conjecture that this is also true for the larger problem sizes. The second choice is especially easy to implement since $D$ is tridiagonal and SPD. The first one requires solving linear systems with the matrix $\tilde{S}_C = M = D + CA^{-1}C^T$. To this end we used again the PCG method, which does not require forming $\tilde{S}_C$ explicitly (it is a full matrix). As a preconditioner we used the (sparse) matrix $D + CD_A^{-1}C^T$, where $D_A$ is the diagonal part of $A$. Application of the preconditioner is accomplished via a sparse Cholesky factorization.

We experimented with different convergence tolerances for the inner PCG method and we found that for the first choice of $M$ in (5.5), a very stringent tolerance (of the order of machine precision) is needed for the outer Uzawa-type method to converge, particulalry for larger problems. For the second choice of $M$ the inner PCG converge tolerance could be relaxed (resulting in an inexact Uzawa method), but the total timings were not significantly affected. We report the performance of the two variants of the iterative method (3.12) in Table 5. As can be seen, convergence is extremely fast, and for larger problems only two steps are required. Clearly, using $M = D$ results in faster solution times.

Finally, we present a few results obtained with the block triangular preconditioners $\tilde{\mathscr{P}}_T$ and $\hat{\mathscr{P}}_T$ given in (4.16)–(4.17). The application of these two preconditioners can be performed "exactly" or inexactly, via block backsubstitution. Both versions of the preconditioners require the solution ("exact" and approximate) of linear systems with SPD coefficient matrices $D + CA^{-1}C^T$, $BA^{-1}B^T$, and $A$ at each (outer) GMRES or FGMRES iteration. The first two systems are solved via the PCG method, with the already described preconditioners for $D + CA^{-1}C^T$ and $BA^{-1}B^T$. The systems with matrix $A$ are solved via sparse Cholesky factorization with SYMAMD reordering when GMRES is used, and with PCG preconditioned with the incomplete Cholesky factorization described earlier in the case of FGMRES. In Table 6, we report the inner tolerances in the PCG method utilized inside the preconditioners $\tilde{\mathscr{P}}_T$ and $\hat{\mathscr{P}}_T$. The symbol "$\star$" in this table means that for solving linear systems with coefficient $A$, we used the sparse Cholesky factorization with SYMAMD reordering.

TABLE 5
*Numerical results for Uzawa's method* (3.12), *liquid crystal problem.*

| size | $M = D + CA^{-1}C^T$ | | $M = D$ | |
|---|---|---|---|---|
| | Iter | CPU | Iter | CPU |
| 5115 | 4 | 0.25794 | 4 | 0.12869 |
| 10235 | 4 | 0.40680 | 3 | 0.16013 |
| 20475 | 3 | 0.67560 | 3 | 0.28280 |
| 40955 | 2 | 0.90899 | 2 | 0.38881 |
| 81915 | 2 | 1.9106 | 2 | 0.82650 |
| 163835 | 2 | 4.2573 | 2 | 1.9536 |
| 327675 | 2 | 10.141 | 2 | 4.6240 |

The results for preconditioners $\tilde{\mathscr{P}}_T$ and $\hat{\mathscr{P}}_T$ are shown in Tables 7 and 8, respectively. In all cases we observe mesh-independent convergence rates, with no deterioration when using the inexact variants of the block preconditioners in place of the exact ones; indeed, in several cases FGMRES even requires one less iteration than GMRES with the "exact" preconditioner. The CPU timings are clearly much better for the inexact variants, especially for larger problems. Overall, the fastest solution times are obtained with FGMRES preconditioned by the inexact variant of the block preconditioner $\hat{\mathscr{P}}_T$. With this method, solution times exhibit almost linear scaling behavior.

**6. Conclusions.** In this paper we have introduced and analyzed several iterative methods and block preconditioners for the solution of sparse linear systems with double saddle point structure. While "standard" techniques for saddle point problems are certainly applicable to systems of the form (1.1), several of the methods investigated in this paper and their analysis make specific use of the $3 \times 3$ block structure of the coefficient matrix. Furthermore, different block partitionings (see (1.2)) lead to different solvers with distinct theoretical and practical properties.

Numerical experiments on test problems arising from two distinct application domains show that some of the proposed solvers can be very efficient in situations of practical interest, resulting in rapid convergence (independent of problem size) and scalable behavior. Of course, the performance of each method is highly problem-dependent, and specific information on the spectral properties of the problem at hand may be needed in order to make a good choice. We stress that it is quite possible that some of the methods that were found to be not competitive for the two test problems considered here may well turn out to be useful on other problems.

In our analysis we assumed that the various solvers and preconditioners were implemented exactly. Numerical experiments, however, showed that the rates of convergence do not necessarily deteriorate when inexact solves are used instead, often leading to significantly faster solution times relative to the "exact" versions. This is consistent with previous experience for block preconditioners; see, e.g., [4] or [7].

Table 6

*Inner tolerance in PCG method used inside the preconditioned methods, liquid crystal problem.*

| Preconditioner: | $\tilde{\mathscr{P}}_T$ | | | $\hat{\mathscr{P}}_T$ | | |
|---|---|---|---|---|---|---|
| Coefficient matrix: | $A$ | $S_B$ | $\tilde{S}_C$ | $A$ | $S_B$ | $\tilde{S}_C$ |
| GMRES | $\star$ | 1e–12 | 1e–12 | $\star$ | 1e–12 | 1e–12 |
| FGMRES | 1e–02 | 1e–02 | 1e–01 | 1e–03 | 1e–03 | 1e–01 |

Table 7

*Numerical results for preconditioner $\tilde{\mathscr{P}}_T$, liquid crystal problem.*

| | Method | | | |
|---|---|---|---|---|
| | GMRES | | FGMRES | |
| size | Iter | CPU | Iter | CPU |
| 5115 | 10 | 0.53144 | 9 | 0.08461 |
| 10235 | 9 | 0.98055 | 9 | 0.14981 |
| 20475 | 9 | 1.9727 | 8 | 0.27168 |
| 40955 | 9 | 3.6642 | 8 | 0.53854 |
| 81915 | 9 | 8.4782 | 8 | 1.2848 |
| 163835 | 9 | 17.947 | 8 | 3.1012 |
| 327675 | 9 | 42.743 | 8 | 7.4957 |

Table 8

*Numerical results for preconditioner $\hat{\mathscr{P}}_T$, liquid crystal problem.*

| | Method | | | |
|---|---|---|---|---|
| | GMRES | | FGMRES | |
| size | Iter | CPU | Iter | CPU |
| 5115 | 6 | 0.32469 | 6 | 0.02981 |
| 10235 | 6 | 0.62408 | 6 | 0.05179 |
| 20475 | 6 | 1.2614 | 6 | 0.08967 |
| 40955 | 6 | 2.4815 | 6 | 0.17196 |
| 81915 | 6 | 5.4721 | 6 | 0.33430 |
| 163835 | 6 | 11.879 | 6 | 0.69814 |
| 327675 | 6 | 28.196 | 6 | 1.5901 |

## REFERENCES

[1] A. Aposporidis, E. Haber, M. A. Olshanskii, and A. Veneziani, *A mixed formulation of the Bingham fluid flow problem: Analysis and numerical solution*, Comput. Methods Appl. Mech. Engrg., 200 (2011), pp. 2434–2446.

[2] Z.-Z. Bai, *Eigenvalue estimates for saddle point matrices of Hermitian and indefinite leading blocks*, J. Comput. Appl. Math., 237 (2013), pp. 295–306.

[3] M. Benzi, *Preconditioning techniques for large linear systems: a survey*, J. Comput. Phys., 182 (2002), pp. 417–477.

[4] M. Benzi, G. H. Golub, and J. Liesen, *Numerical solution of saddle point problems*, Acta Numer., 14 (2005), pp. 1–137.

[5] D. Boffi, F. Brezzi, and M. Fortin, *Mixed Finite Element Methods and Applications*, Springer Ser. Comput. Math., Springer-Verlag, New York, 2013.

[6] P. Chidyagway, S. Ladenheim, and D. B. Szyld, *Constraint preconditioning for the coupled Stokes–Darcy system*, SIAM J. Sci. Comput., 38 (2016), pp. A668–A690.

[7] H. C. Elman, D. J. Silvester, and A. J. Wathen, *Finite Elements and Fast Iterative Solvers. With Applications in Incompressible Fluid Dynamics*, Second Edition, Oxford University Press, Oxford, UK, 2014.

[8] N. I. M. Gould and V. Simoncini, *Spectral analysis of saddle point matrices with indefinite leading blocks*, SIAM J. Matrix Anal. Appl., 31 (2009), pp. 1152–1171.

[9] R. A. Horn and C. R. Johnson, *Matrix Analysis*, Cambridge University Press, Cambridge, UK, 1985.

[10] J. Maryška, M. Rozložník, and M. Tůma, *Schur complement systems in the mixed-hybrid finite element approximation of the potential fluid flow problem*, SIAM J. Sci. Comput., 22 (2005), pp. 704–723.

[11] B. Morini, V. Simoncini, and M. Tani, *Spectral estimates for unreduced symmetric KKT systems arising from Interior Point methods*, Numer. Linear Algebra Appl., 23 (2016), pp. 776–800.

[12] B. Morini, V. Simoncini, and M. Tani, *A comparison of reduced and unredued KKT systems arising from interior point methods*, Preprint, University of Bologna, February 2016.

[13] C. C. Paige and M. A. Saunders, *Solution of sparse indefinite systems of linear equations*, SIAM J. Numer. Anal., 12 (1975), pp. 617–629.

[14] J. Pestana, *On the eigenvalues and eigenvectors of block triangular preconditioned block matrices*, SIAM J. Matrix Anal. Appl., 35 (2014), pp. 517–525.

[15] A. Ramage and E. C. Gartland, Jr., *A preconditioned nullspace method for liquid crystal director modeling*, SIAM J. Sci. Comput., 35 (2013), pp. B226–B247.

[16] Y. Saad, *Iterative Methods for Sparse Linear Systems. Second Edition*, Society for Industrial and Applied Mathematics, Philadelphia, 2003.

[17] F. Zhang and Q. Zhang, *Eigenvalue inequalities for matrix product*, IEEE Trans. Automat. Control, 51 (2006), pp. 1506–1509.