# Analysis of block matrix preconditioners for elliptic optimal control problems

T. P. Mathew[1], M. Sarkis[1, 2, *, †] and C. E. Schaerer[1]

[1]*Instituto Nacional de Matemática Pura e Aplicada-IMPA, Estrada Dona Castorina 110, Rio de Janeiro, RJ 22460-320, Brazil*
[2]*Department of Mathematical Sciences, Worcester Polytechnic Institute, Worcester, MA 01609, U.S.A.*

## SUMMARY

In this paper, we describe and analyse several block matrix iterative algorithms for solving a *saddle point* linear system arising from the discretization of a linear-quadratic *elliptic control* problem with Neumann boundary conditions. To ensure that the problem is well posed, a *regularization* term with a parameter $\alpha$ is included. The first algorithm reduces the saddle point system to a symmetric positive definite Schur complement system for the control variable and employs conjugate gradient (CG) acceleration, however, double iteration is required (except in special cases). A preconditioner yielding a rate of convergence independent of the mesh size $h$ is described for $\Omega \subset R^2$ or $R^3$, and a preconditioner independent of $h$ and $\alpha$ when $\Omega \subset R^2$. Next, two algorithms avoiding double iteration are described using an *augmented Lagrangian* formulation. One of these algorithms solves the augmented saddle point system employing MINRES acceleration, while the other solves a symmetric positive definite reformulation of the augmented saddle point system employing CG acceleration. For both algorithms, a symmetric positive definite preconditioner is described yielding a rate of convergence independent of $h$. In addition to the above algorithms, two *heuristic* algorithms are described, one a projected CG algorithm, and the other an indefinite block matrix preconditioner employing GMRES acceleration. Rigorous convergence results, however, are not known for the heuristic algorithms. Copyright © 2007 John Wiley & Sons, Ltd.

## 1. INTRODUCTION

In this paper, we study the convergence of several iterative methods for solving a linear-quadratic elliptic optimal control problem with Neumann boundary conditions [1–5]. Such problems seek to

---

*Correspondence to: M. Sarkis, Instituto Nacional de Matemática Pura e Aplicada-IMPA, Estrada Dona Castorina 110, Rio de Janeiro, RJ 22460-320, Brazil.
†E-mail: msarkis@impa.br

WILEY
InterScience®
DISCOVER SOMETHING GREAT

determine a control function $u(\cdot)$ defined on the boundary $\partial\Omega$ of a domain $\Omega$, to minimize some performance functional $J(y, u)$ of the form

$$J(y, u) \equiv \tfrac{1}{2}(\|y - \hat{y}\|_{L^2(\Omega_0)}^2 + \alpha_1 \|u\|_{L^2(\partial\Omega)}^2 + \alpha_2 \|u\|_{H^{-1/2}(\partial\Omega)}^2) \qquad (1)$$

where $\hat{y}(.)$ is a given target function that we seek to match on $\Omega_0 \subset \Omega$ with the solution $y(\cdot) \in \mathscr{V}_f$ to the elliptic problem in (2) with Neumann data $u(\cdot)$. Here, $\|u\|_{H^{-1/2}(\partial\Omega)}^2$ denotes a dual Sobolev norm that will be defined later, and $\mathscr{V}_f$ denotes the affine space

$$\mathscr{V}_f \equiv \left\{(y, u) : -\Delta y(x) + \sigma y(x) = f(x) \text{ in } \Omega \text{ and } \frac{\partial y(x)}{\partial n} = u(x) \text{ on } \partial\Omega \right\} \qquad (2)$$

defined in terms of a forcing $f(\cdot)$ and parameter $\sigma > 0$. The parameters $\alpha_1$, $\alpha_2$ are chosen to *regularize* the functional $J(y, u)$ to yield a well posed problem. A typical choice is $\alpha_1 > 0$ and small, with $\alpha_2 = 0$. However, we shall also consider $\alpha_1 = 0$ with $\alpha_2 > 0$, which yields a weaker regularization term, and show that it has attractive computational properties.

The finite element discretization of an elliptic optimal control problem yields a saddle point linear system with a coefficient matrix that is symmetric indefinite. There is extensive literature on saddle point iterative methods, see [6], while specific preconditioners have been studied for discretizations of optimal control problems [1, 2, 4, 7–12]. Our discussion focuses on the analysis of block matrix algorithms based on conjugate gradient (CG) or MINRES acceleration [6, 11, 13–18, 20, 21]. The first algorithm we consider requires double iteration and is based on the solution of a reduced Schur complement system for the control variable $u$. We describe a preconditioner which yields a well-conditioned system with respect to $h$, but dependent on $\alpha$, for $\Omega \subset R^2$ or $R^3$, and a preconditioner which yields a well-conditioned system with respect to $h$ and $\alpha$ when $\Omega \subset R^2$. The second family of algorithms we study avoids double iteration, and employs an *augmented Lagrangian* reformulation of the original saddle point system [22]. Motivated by [18–20], we describe a symmetric positive definite preconditioner for the augmented system, employing MINRES acceleration, and a similar preconditioner, motivated by [15, 21, 23], for a symmetric positive definite reformulation of the augmented system, employing CG acceleration. In both the cases, the preconditioners yield a rate of convergence independent of the mesh size $h$, but dependent on the regularization parameters. We also describe a *heuristic* algorithm based on the projected gradient method (motivated by [24]) and a non-symmetric block matrix preconditioner based on GMRES acceleration.

This paper is organized as follows. In Section 2, we formulate the linear-quadratic elliptic control problem with Neumann boundary conditions. Its weak formulation and finite element discretization are described, and the block matrix form of the resulting saddle point system (with Lagrange multiplier $p(\cdot)$). In Section 3, we describe a reduced Schur complement system for the control variable $u$ (obtained by formal elimination of $y$ and the Lagrange multiplier variable $p$). The reduced system is symmetric positive definite, and we describe suitable preconditioners for it, requiring double iteration. In Section 4, we describe an augmented Lagrangian reformulation of the original saddle point system [22] to regularize the system without altering its solution. We describe a symmetric positive definite block diagonal preconditioner for the augmented saddle point system, for use with MINRES acceleration, and a similar preconditioner for a symmetric positive definite reformulation of the augmented saddle point system, for use with CG acceleration. The rates of convergence are shown to be independent of $h$, but dependent on the regularization parameters. In Section 5, we outline alternative algorithms, one based on the projected gradient

method (motivated by [24]), and another based on block matrix preconditioning of the original saddle point system (using GMRES acceleration).

## 2. OPTIMAL CONTROL PROBLEM

Let $\Omega \subset R^d$ denote a polygonal domain with boundary $\partial\Omega$. We consider the problem of determining a control function $u(\cdot)$ denoting the Neumann data on $\partial\Omega$, so that the solution $y(\cdot)$ to the following Neumann problem with forcing term $f(\cdot)$:

$$-\Delta y(x) + \sigma\, y(x) = f(x) \quad \text{in } \Omega$$

$$\frac{\partial y(x)}{\partial n} = u(x) \quad \text{on } \partial\Omega \tag{3}$$

minimizes the following performance functional $J(y, u)$:

$$J(y, u) = \frac{1}{2}\left( \|y - \hat{y}\|^2_{L^2(\Omega_0)} + \alpha_1 \|u\|^2_{L^2(\partial\Omega)} + \alpha_2 \|u\|^2_{H^{-1/2}(\partial\Omega)} \right) \tag{4}$$

where $\hat{y}(\cdot)$ is a given target, and $\alpha_1$, $\alpha_2 \geqslant 0$ are *regularization* parameters. For simplicity, we shall assume that $\sigma > 0$, and as a result our theoretical bounds will depend on $\sigma$. The term $\|u\|_{H^{-1/2}(\partial\Omega)}$ denotes a dual Sobolev norm

$$\|u\|_{H^{-1/2}(\partial\Omega)} \equiv \sup_{v \in H^{1/2}(\partial\Omega)} \frac{\int_{\partial\Omega} u\, v\, \mathrm{d}s_x}{\|v\|_{H^{1/2}(\partial\Omega)}}$$

where $H^{1/2}(\partial\Omega) = [L^2(\partial\Omega), H^1(\partial\Omega)]_{1/2}$ is a fractional index Sobolev space defined using Hilbert scales, see [25]. An integral expression for $\|v\|_{H^{1/2}(\partial\Omega)}$ can be found in [25].

To obtain a weak formulation of the minimization of (4) within the constraint set (3), we employ the function space $H^1(\Omega)$ for $y(\cdot)$ and $H^{-1/2}(\partial\Omega)$ for $u(\cdot)$. Given $f(\cdot) \in L^2(\Omega)$, define the constraint set $\mathcal{V}_f \subset \mathcal{V} \equiv H^1(\Omega) \times H^{-1/2}(\partial\Omega)$ as follows:

$$\mathcal{V}_f \equiv \{(y, u) \in \mathcal{V} : \mathcal{A}(y, w) = (f, w) + \langle u, w \rangle, \ \forall w \in H^1(\Omega)\} \tag{5}$$

where the forms are defined by

$$\mathcal{A}(u, w) \equiv \int_\Omega (\nabla u \cdot \nabla w + \sigma u w)\, \mathrm{d}x \quad \text{for } u, w \in H^1(\Omega)$$

$$(f, w) \equiv \int_\Omega f(x)\, w(x)\, \mathrm{d}x \quad \text{for } w \in H^1(\Omega) \tag{6}$$

$$\langle u, w \rangle \equiv \int_{\partial\Omega} u(x) w(x)\, \mathrm{d}s_x \quad \text{for } u \in H^{-1/2}(\partial\Omega), \ w \in H^{1/2}(\partial\Omega)$$

The constrained minimization problem then seeks $(y_*, u_*) \in \mathcal{V}_f$ satisfying

$$J(y_*, u_*) = \min_{(y,u) \in \mathcal{V}_f} J(y, u) \tag{7}$$

*Remark 1*

The regularization terms $(\alpha_1/2) \int_{\partial\Omega} u^2(x)\,ds_x + (\alpha_2/2)\|u\|^2_{H^{-1/2}(\partial\Omega)}$ must be chosen to modify $J(y, u)$ so that the minimization problem (7) is well posed. When $f(\cdot) = 0$, the constraint set $\mathcal{V}_0$ will be a *closed subspace* of $H^1(\Omega) \times H^{-1/2}(\partial\Omega)$, yet the minimization of $J(y, u)$ within $\mathcal{V}_0$ will not be well posed for $\alpha_1 = 0$ and $\alpha_2 = 0$ (due to the $L^2(\Omega_0)$ term). To ensure well posedness of (7), saddle point theory [26, 27] requires the functional $J(., .)$ to be *coercive* within $\mathcal{V}_0$. When $\alpha_1 > 0$ and $\alpha_2 = 0$, it can be shown that $J(y, u)$ is coercive in $\mathcal{V}_0$ (though $\|u\|^2_{L^2(\partial\Omega)}$ is not defined for $u \in H^{-1/2}(\partial\Omega)$, it will be defined for finite element approximations). When $\alpha_1 = 0$ and $\alpha_2 > 0$, elliptic regularity theory shows that $J(y, u)$ is coercive within $\mathcal{V}_0$. This regularization term has the advantage of involving a weaker norm.

To obtain a saddle point formulation of (7), introduce $p(\cdot) \in H^1(\Omega)$ as a Lagrange multiplier function to enforce the constraints. Define the following Lagrangian functional $\mathcal{L}(\cdot, \cdot, \cdot)$:

$$\mathcal{L}(y, u, p) \equiv J(y, u) + (\mathcal{A}(y, p) - (f, p) - \langle u, p \rangle) \tag{8}$$

for $(y, u, p) \in H^1(\Omega) \times H^{-1/2}(\partial\Omega) \times H^1(\Omega)$. Then, the constrained minimum $(y_*, u_*)$ can be obtained from the saddle point $(y_*, u_*, p_*) \in H^1(\Omega) \times H^{-1/2}(\partial\Omega) \times H^1(\Omega)$ of $\mathcal{L}(\cdot, \cdot, \cdot)$

$$\sup_q \mathcal{L}(y_*, u_*, q) = \mathcal{L}(y_*, u_*, p_*) = \inf_{(y,u)} \mathcal{L}(y, u, p_*) \tag{9}$$

Saddle point problem (9) will be *well posed* if an *inf–sup* condition holds (it will hold trivially for this problem), and if $J(., .)$ is *coercive* within the subspace $\mathcal{V}_0$, see [26, 27]. As mentioned before, the coercivity of $J(y, u)$ can be proved within the subspace $\mathcal{V}_0$ when $\alpha_1 = 0$ and $\alpha_2 > 0$, or when $\alpha_1 > 0$ and $\alpha_2 = 0$, but not when $\alpha_1 = 0$ and $\alpha_2 = 0$. In a strict sense, the term $\|u\|_{L^2(\partial\Omega)}$ will not be defined for $u \in H^{-1/2}(\partial\Omega)$. However, this term will be well defined for finite element functions.

*Remark 2*

If $\hat{y}(\cdot)$ is sufficiently smooth, the functional $J(y, u) = \frac{1}{2}\|y - \hat{y}\|^2_{H^1(\Omega_0)}$ can also be employed. When $\Omega_0 = \Omega$, the functional $J(y, u) = \frac{1}{2}\|y - \hat{y}\|^2_{H^1(\Omega)}$ can easily be shown to be coercive within $\mathcal{V}_0$ without additional regularization terms, and the saddle point problem will be well posed. Efficient computational algorithms considered in this paper can easily be adapted to this case.

To obtain a finite element discretization of the constrained minimization problem, choose a quasi-uniform triangulation $\tau_h(\Omega)$ of $\Omega$. Let $V_h(\Omega) \subset H^1(\Omega)$ denote a finite element space [28–30] defined on $\tau_h(\Omega)$, and let $V_h(\partial\Omega) \subset L^2(\partial\Omega)$ denote its restriction to $\partial\Omega$. A finite element discretization of (7) will seek $(y_h^*, u_h^*) \in V_h(\Omega) \times V_h(\partial\Omega)$ such that

$$J(y_*, u_*) = \min_{(y_h, u_h) \in \mathcal{V}_{h,f}} J(y_h, u_h) \tag{10}$$

where the discrete constraint space $\mathcal{V}_{h,f} \subset \mathcal{V}_h \equiv V_h(\Omega) \times V_h(\partial\Omega)$ is defined by

$$\mathcal{V}_{h,f} = \{(y_h, u_h) \in \mathcal{V}_h : \mathcal{A}(y_h, w_h) = (f, w_h) + \langle u_h, w_h \rangle, \ \forall w_h \in V_h(\Omega)\} \tag{11}$$

Let $p_h \in V_h(\Omega)$ denote discrete Lagrange multiplier variables, and let $\{\phi_1(x), \ldots, \phi_n(x)\}$ and $\{\psi_1(x), \ldots, \psi_m(x)\}$ denote standard nodal finite element basis for $V_h(\Omega)$ and $V_h(\partial\Omega)$, respectively.

Then, expanding each unknown $y_h$, $u_h$ and $p_h$ with respect to the basis for each finite element space

$$y_h(x) = \sum_{i=1}^{n} \mathbf{y}_i \phi_i(x), \quad u_h(x) = \sum_{j=1}^{m} \mathbf{u}_j \psi_j(x), \quad p_h(x) = \sum_{l=1}^{n} \mathbf{p}_l \phi_l(x) \tag{12}$$

and substituting into the weak saddle point formulation yields the system:

$$\begin{bmatrix} M & 0 & A^{\mathrm{T}} \\ 0 & G & B^{\mathrm{T}} \\ A & B & 0 \end{bmatrix} \begin{bmatrix} \mathbf{y} \\ \mathbf{u} \\ \mathbf{p} \end{bmatrix} = \begin{bmatrix} \mathbf{f}_1 \\ \mathbf{f}_2 \\ \mathbf{f}_3 \end{bmatrix} \tag{13}$$

where the block submatrices $M$ and $A$, and the matrix $Q$ to be used later, are defined by

$$M_{ij} \equiv \int_{\Omega_0} \phi_i(x)\phi_j(x)\,\mathrm{d}x \quad \text{for } 1 \leqslant i, j \leqslant n$$

$$A_{ij} \equiv \int_{\Omega} (\nabla\phi_i(x) \cdot \nabla\phi_j(x) + \sigma\phi_i(x)\phi_j(x))\,\mathrm{d}x \quad \text{for } 1 \leqslant i, j \leqslant n \tag{14}$$

$$Q_{ij} \equiv \int_{\partial\Omega} \psi_i(x)\psi_j(x)\,\mathrm{d}s_x \quad \text{for } 1 \leqslant i, j \leqslant m$$

and the discrete forcing are defined by $(\mathbf{f}_1)_i = \int_{\Omega_0} \hat{y}(x)\phi_i(x)\,\mathrm{d}x$, for $1 \leqslant i \leqslant n$ with $\mathbf{f}_2 = \mathbf{0}$, and $(\mathbf{f}_3)_i = \int_{\Omega} f(x)\phi_i(x)\,\mathrm{d}x$ for $1 \leqslant i \leqslant n$. The matrix $M$ of dimension $n$ corresponds to a mass matrix on $\Omega_0$, and the matrix $A$ to the Neumann stiffness matrix. The matrix $Q$ of dimension $m$ corresponds to a lower dimensional mass matrix on $\partial\Omega$. The matrix $B$ will be defined in terms of $Q$, based on the following ordering of nodal unknowns in $\mathbf{y}$ and $\mathbf{p}$. Order the nodes in the *interior* of $\Omega$ prior to the nodes on $\partial\Omega$. Denote such block partitioned vectors as $\mathbf{y} = (\mathbf{y}_I^{\mathrm{T}}, \mathbf{y}_B^{\mathrm{T}})^{\mathrm{T}}$ and $\mathbf{p} = (\mathbf{p}_I^{\mathrm{T}}, \mathbf{p}_B^{\mathrm{T}})^{\mathrm{T}}$, and define $B$ of dimension $n \times m$ as

$$B = \begin{bmatrix} 0 \\ -Q \end{bmatrix} \quad \text{and} \quad B^{\mathrm{T}} = [0 \ -Q^{\mathrm{T}}] \tag{15}$$

and define matrix $G$ of dimension $m$, representing the regularizing terms as

$$G \equiv \alpha_1 Q + \alpha_2(B^{\mathrm{T}}A^{-1}B) \tag{16}$$

It will be shown later, that $\mathbf{u}^{\mathrm{T}}(B^{\mathrm{T}}A^{-1}B)\mathbf{u}$ is spectrally equivalent to $\|u_h\|^2_{H^{-1/2}(\partial\Omega)}$, when $\mathbf{u}$ is the nodal vector associated with a finite element function $u_h(\cdot)$.

*Remark 3*
The discrete performance functional has the representation

$$J(\mathbf{y}, \mathbf{u}) = \frac{1}{2} \begin{bmatrix} \mathbf{y} \\ \mathbf{u} \end{bmatrix}^{\mathrm{T}} \begin{bmatrix} M & 0 \\ 0 & G \end{bmatrix} \begin{bmatrix} \mathbf{y} \\ \mathbf{u} \end{bmatrix} - \begin{bmatrix} \mathbf{y} \\ \mathbf{u} \end{bmatrix}^{\mathrm{T}} \begin{bmatrix} \mathbf{f}_1 \\ \mathbf{f}_2 \end{bmatrix} \tag{17}$$

after omission of a constant term, while the discretized constraints have the form

$$A\mathbf{y} + B\mathbf{u} = \mathbf{f}_3 \qquad (18)$$

The following properties can be easily verified for the matrices $M$ and $A$. The matrix $M$ will be singular when $\Omega_0 \neq \Omega$. The matrix $A$ will be symmetric and positive definite when $\sigma > 0$. However, when $\sigma = 0$, matrix $A$ will be *singular* with $A\mathbf{1} = \mathbf{0}$ for $\mathbf{1} = (1, \ldots, 1)^{\mathrm{T}}$. In this case, we shall require $\mathbf{1}^{\mathrm{T}}\mathbf{f}_3 = 0$ for solvability of (13). Theory for saddle systems [27] requires the quadratic form $\mathbf{y}^{\mathrm{T}}M\mathbf{y} + \mathbf{u}^{\mathrm{T}}G\mathbf{u}$ to be positive definite for $A\mathbf{y} + B\mathbf{u} = \mathbf{0}$ and $\mathbf{y} \neq \mathbf{0}$ and $\mathbf{u} \neq \mathbf{0}$. This will ensure solvability of (13) (provided $\mathbf{1}^{\mathrm{T}}\mathbf{f}_3 = 0$ when $A\mathbf{1} = \mathbf{0}$).

*Remark 4*

The role of regularization and the role of the parameters $\alpha_i$ can be heuristically understood by considering the following least-squares problem. Let $H$ be a rectangular or singular matrix of dimension $m \times n$ with singular value decomposition $H = U\Sigma V^{\mathrm{T}}$. Then, a minimum of the least-squares functional $F(\mathbf{x}) = \frac{1}{2}\|H\mathbf{x} - \mathbf{b}\|^2$ will be given by $\mathbf{x}_* = H^{\dagger}\mathbf{b} = V\Sigma^{\dagger}U^{\mathrm{T}}\mathbf{b}$. If $H$ has a non-trivial null space, there will be an affine space of minima. Indeed, if $N$ is a matrix of dimension $n \times k$ whose columns form a basis for the null space of $H$, with $\mathrm{Range}(N) = \mathrm{Kernel}(H)$, then a general minimum of $F(\mathbf{x})$ will be $\mathbf{x}_* + N\boldsymbol{\beta}$ for any vector $\boldsymbol{\beta} \in R^k$. If we employ partial regularization and define $\tilde{F}(\mathbf{x}) = F(\mathbf{x}) + (\alpha/2)\|P_N\mathbf{x}\|^2$ where $P_N$ denotes the Euclidean orthogonal projection onto the null space of $H$, then the minimum of $\tilde{F}(\cdot)$ can be verified to be unique and occur at $\mathbf{x}_* = H^{\dagger}\mathbf{b}$ for any $\alpha > 0$. If, however, a regularization term of the form $(\alpha/2)\|\mathbf{x}\|^2$ is employed, and the minimum of $\hat{F}(\mathbf{x}) = F(\mathbf{x}) + (\alpha/2)\|\mathbf{x}\|^2$ is sought, this will yield the linear system $(H^{\mathrm{T}}H + \alpha I)\mathbf{x} = H^{\mathrm{T}}\mathbf{b}$. Using, the singular value decomposition of $H$, we may obtain the following representation of the unique solution to the regularized problem $\mathbf{x} = V(\Sigma^{\mathrm{T}}\Sigma + \alpha I)^{-1}\Sigma^{\mathrm{T}}U^{\mathrm{T}}\mathbf{b}$. The $i$th diagonal entry of $(\Sigma^{\mathrm{T}}\Sigma + \alpha I)^{-1}\Sigma^{\mathrm{T}}$ will be $\sigma_i/(\sigma_i^2 + \alpha)$, so that if $\sigma_i > 0$, then $\sigma_i/(\sigma_i^2 + \alpha) \to 1/\sigma_i$ as $\alpha \to 0^+$, while if $\sigma_i = 0$, then $\sigma_i/(\sigma_i^2 + \alpha) = 0$. Thus, $\mathbf{x} \to \mathbf{x}_* = H^{\dagger}\mathbf{b}$ as $\alpha \to 0^+$. In our applications, matrix $H$ will correspond to $(B^{\mathrm{T}}A^{-\mathrm{T}}MA^{-1}B)$, while $\mathbf{x}$ will correspond to $\mathbf{u}$ and $F(.)$ to $J(.)$.

*Remark 5*

The choice of parameter $\alpha > 0$ will typically be problem dependent. When matrix $H$ arises from the discretization of a *well-posed* problem, the singular values of $H$ will be bounded away from 0. In this case, if $\alpha$ is chosen appropriately smaller than the smallest non-zero singular value of $H$, the regularized solution will approach the pseudo-inverse solution. However, when matrix $H$ arises from the discretization of an *ill-posed* problem, its singular values will *cluster* around 0. In this case, care must be exercised in the choice of regularization parameter $\alpha > 0$, to balance the accuracy of the modes associated with the large singular values, and to dampen the modes associated with the very small singular values.

*Remark 6*

In applications, alternate performance functionals may be employed, which measure the difference between $y(\cdot)$ and $\hat{y}(\cdot)$ at different subregions of $\Omega$. For instance, given nodes $z_1, \ldots, z_r \in \Omega$, we may minimize the distance between $y(\cdot)$ and $\hat{y}(\cdot)$ at these points:

$$J(y, u) = \frac{1}{2}\left(\sum_{l=1}^{r} |y(z_l) - \hat{y}(z_l)|^2 + \alpha_1\|u\|_{L^2(\partial\Omega)}^2 + \alpha_2\|u\|_{H^{-1/2}(\partial\Omega)}^2\right)$$

This performance functional requires the measurement of $y(x) - \hat{y}(x)$ at the $r$ discrete nodes. We must choose either $\alpha_1 > 0$ or $\alpha_2 > 0$ to regularize the problem. Another performance functional, described below, requires the measurement of $y(x) - \hat{y}(x)$ only on $\partial \Omega$

$$J(y, u) = \frac{1}{2} \left( \int_{\partial \Omega} |y(x) - \hat{y}(x)|^2 \, ds_x + \alpha_1 \|u\|^2_{L^2(\partial \Omega)} + \alpha_2 \|u\|^2_{H^{-1/2}(\partial \Omega)} \right) \tag{19}$$

We shall obtain $M = \text{blockdiag}(0, Q)$ and require $\alpha_1 > 0$ or $\alpha_2 > 0$ to regularize the problem.

## 3. PRECONDITIONED SCHUR COMPLEMENT ALGORITHMS

The first algorithm we consider for solving (13) is based on the solution of a reduced system for the discrete control $\mathbf{u}$. We shall assume that $\sigma > 0$ and that $G > 0$. Then, formally solving the third block row in (13) will yield $\mathbf{y} = A^{-1}(\mathbf{f}_3 - B\mathbf{u})$. Solving the first block row in (13) will yield $\mathbf{p} = A^{-T}(\mathbf{f}_1 - MA^{-1}\mathbf{f}_3 + MA^{-1}B\mathbf{u})$. Substituting these into the second block row of (13) will yield the following reduced Schur complement system for $\mathbf{u}$:

$$(G + B^T A^{-T} M A^{-1} B)\mathbf{u} = \mathbf{f}_2 - B^T A^{-T} \mathbf{f}_1 + B^T A^{-T} M A^{-1} \mathbf{f}_3 \tag{20}$$

The Schur complement matrix $(G + B^T A^{-T} M A^{-1} B)$ will be symmetric and positive definite of dimension $m$, and system (20) can be solved using a PCG algorithm. Each matrix vector product with $G + B^T A^{-T} M A^{-1} B$ will require the action of $A^{-1}$ twice per iteration (this can be computed iteratively, resulting in a *double iteration*). Once $\mathbf{u}$ has been determined by solution of (20), we obtain $\mathbf{y} = A^{-1}(\mathbf{f}_3 - B\mathbf{u})$ and $\mathbf{p} = A^{-T}(\mathbf{f}_1 - MA^{-1}\mathbf{f}_3 + MA^{-1}B\mathbf{u})$. The following result shows that if the parameters $\alpha_1 > 0$ or $\alpha_2 > 0$ are held fixed independent of $h$, then matrix $G$ will be spectrally equivalent to the Schur complement $(G + B^T A^{-T} M A^{-1} B)$, and can be used as a preconditioner. Unfortunately, in practice $\alpha_i$ may be small (and possibly dependent on $h$), and for such a case alternative preconditioners will be described later in Remark 7 and Subsection 3.1.

*Lemma 3.1*
Suppose that $\alpha_1 > 0$ and $\alpha_2 = 0$ or $\alpha_1 = 0$ and $\alpha_2 > 0$. Then, there exists $\gamma, \tilde{c} > 0$ independent of $h$, $\alpha_1$, and $\alpha_2$ such that:

$$(\mathbf{u}^T G \mathbf{u}) \leqslant \mathbf{u}^T (G + B^T A^{-T} M A^{-1} B)\mathbf{u} \leqslant (1 + c)(\mathbf{u}^T G \mathbf{u}) \quad \forall \mathbf{u} \in R^m \tag{21}$$

where $c = (\gamma/\alpha_1)\tilde{c}$ when $\alpha_2 = 0$, and $c = (\gamma/\alpha_2)$ when $\alpha_1 = 0$.

*Proof*
The lower bound follows trivially since $(B^T A^{-T} M A^{-1} B) \geqslant 0$. To obtain the upper bound, employ Poincare–Friedrichs' inequality which yields $\gamma > 0$ independent of $h$ such that $\mathbf{y}^T M \mathbf{y} \leqslant \gamma \mathbf{y}^T A \mathbf{y}$. Substituting this, yields

$$(\mathbf{u}^T B^T A^{-T} M A^{-1} B \mathbf{u}) \leqslant \gamma (\mathbf{u}^T B^T A^{-T} A A^{-1} B \mathbf{u}) = \gamma (\mathbf{u}^T B^T A^{-1} B \mathbf{u})$$

When $\alpha_1 = 0$, matrix $G = \alpha_2 (B^T A^{-1} B)$ and the desired bound will hold trivially for $c = (\gamma/\alpha_2)$. When $\alpha_2 = 0$, employ the block partition $\mathbf{y} = (\mathbf{y}_I^T, \mathbf{y}_B^T)^T$ to obtain

$$B^T A^{-1} B = \begin{bmatrix} 0 \\ -Q \end{bmatrix}^T \begin{bmatrix} A_{II} & A_{IB} \\ A_{IB}^T & A_{BB} \end{bmatrix}^{-1} \begin{bmatrix} 0 \\ -Q \end{bmatrix} = Q^T (A_{BB} - A_{IB}^T A_{II}^{-1} A_{IB})^{-1} Q$$

The matrix $S = (A_{BB} - A_{IB}^T A_{II}^{-1} A_{IB})$ denotes the discrete Dirichlet to Neumann map, and is known to be symmetric and positive (when $\sigma > 0$), see [31]. Let $u_h$, $w_h \in V_h(\partial\Omega)$ denote finite element functions associated with **u** and **w**, respectively. Then, using properties of $S$ yields:

$$\mathbf{u}^T B^T A^{-1} B\mathbf{u} = \mathbf{u}^T Q^T S^{-1} Q\mathbf{u} = \|S^{-1/2}Q\mathbf{u}\|^2$$

$$= \left( \sup_{\mathbf{v}\in R^m} \frac{(S^{-1/2}Q\mathbf{u}, \mathbf{v})}{\|\mathbf{v}\|} \right)^2$$

$$= \left( \sup_{\mathbf{v}\in R^m} \frac{(Q\mathbf{u}, S^{-1/2}\mathbf{v})}{\|\mathbf{v}\|} \right)^2$$

$$= \left( \sup_{\mathbf{w}\in R^m} \frac{(Q\mathbf{u}, \mathbf{w})}{\|S^{1/2}\mathbf{w}\|} \right)^2$$

$$\leqslant \tilde{c} \left( \sup_{w_h\in V_h(\partial\Omega)} \frac{\langle u_h, w_h \rangle}{\|w_h\|_{1/2,\partial\Omega}} \right)^2$$

$$\leqslant \tilde{c} \left( \sup_{w\in H^{1/2}(\partial\Omega)} \frac{\langle u_h, w \rangle}{\|w\|_{1/2,\partial\Omega}} \right)^2$$

$$= \tilde{c}(\|u_h\|_{-1/2,\partial\Omega})^2$$

$$\leqslant \tilde{c}(\|u_h\|_{0,\partial\Omega})^2$$

$$= \tilde{c}\mathbf{u}^T Q\mathbf{u} \tag{22}$$

where $\tilde{c}$ denotes a parameter independent of $h$, which bounds the energy associated with $S$ in terms of the fractional Sobolev norm $H^{1/2}(\partial\Omega)$. This equivalence between $\|S^{1/2}\mathbf{w}\|$ and $\|w_h\|_{1/2,\partial\Omega}$ is a standard result in domain decomposition literature [31]. We used $(\cdot, \cdot)$ to denote the Euclidean inner product with norm $\|\cdot\|$, and $\langle\cdot, \cdot\rangle$ to denote the duality pairing between $H^{1/2}(\partial\Omega)$ and $H^{-1/2}(\partial\Omega)$ (pivoted using the $L^2(\partial\Omega)$ inner product), and $u_h(\cdot)$ to denote the finite element function associated with a nodal vector **u**. We also employed the definition of dual norms of Sobolev spaces and the property that $\|u\|_{-1/2,\partial\Omega} \leqslant \|u\|_{0,\partial\Omega}$ when $u \in L^2(\partial\Omega)$. This yields the upper bound $c = (\gamma\tilde{c}/\alpha_1)$ when $G = \alpha_1 Q$. Importantly, under additional assumptions $\mathbf{u}^T Q^T S^{-1} Q\mathbf{u}$ will be equivalent to $\|u_h\|^2_{H^{-1/2}(\partial\Omega)}$, see Remark 7. $\qquad\square$

*Remark 7*
The first *inequality* in (22) will also be an equivalence [31]. The second inequality in (22) will be an equivalence too by considering $w_h$ as the $L_2(\partial\Omega)$ projection of $w$ on $V_h(\partial\Omega)$ and by using the stability of the $L_2(\partial\Omega)$ projection in the $H^{1/2}(\partial\Omega)$ norm [31]. Henceforth, let $\asymp$ denote an equivalence independent of $h$ and $\alpha_i$. Thus, $\mathbf{u}^T Q^T S^{-1} Q\mathbf{u} \asymp \|u_h\|^2_{H^{-1/2}(\partial\Omega)}$. The upper bound in Lemma 3.1 *deteriorates* as $\max\{\alpha_1, \alpha_2\} \to 0^+$. As a result, $G$ may not be a uniformly effective preconditioner for $(G + B^T A^{-T} M A^{-1} B)$ as $\alpha_i \to 0^+$. The spectral properties of the Schur complement $(G + B^T A^{-T} M A^{-1} B)$ can differ significantly from those of $G$ and $(B^T A^{-T} M A^{-1} B)$,

depending on the *weights* $\alpha_i$. For instance, when $\lambda_{\min}(G) \geqslant \lambda_{\max}(B^\mathrm{T} A^{-\mathrm{T}} M A^{-1} B)$ it can be verified that $\mathrm{cond}(G, G + B^\mathrm{T} A^{-\mathrm{T}} M A^{-1} B) \leqslant 2$, so that $G$ will be an effective preconditioner. On the other hand, if $M$ is non-singular and $\lambda_{\max}(G) \leqslant \lambda_{\min}(B^\mathrm{T} A^{-\mathrm{T}} M A^{-1} B)$ then

$$\mathrm{cond}(B^\mathrm{T} A^{-\mathrm{T}} M A^{-1} B, G + B^\mathrm{T} A^{-\mathrm{T}} M A^{-1} B) \leqslant 2$$

so that $(B^\mathrm{T} A^{-\mathrm{T}} M A^{-1} B)$ will be a more effective preconditioner. Some preconditioners which are uniformly effective with respect to $\alpha_1$, $\alpha_2$ and $h$ will be considered in Section 3.1.

*Remark 8*
When $M = 0$, computing the action of $G + B^\mathrm{T} A^{-\mathrm{T}} M A^{-1} B$ on a vector will be trivial. In this case, the Schur complement system for **u** can be solved *without* double iteration, retaining a convergence rate independent of $h$. If matrix $M$ is of *low rank l*, then matrix $B^\mathrm{T} A^{-\mathrm{T}} M A^{-1} B$ can be assembled explicitly (at the cost of $l$ matrix products with $A$), and the Sherman–Morrison–Woodbury formula can be employed to compute the solution to the perturbed system $G + B^\mathrm{T} A^{-\mathrm{T}} M A^{-1} B$. For instance, if $B^\mathrm{T} A^{-\mathrm{T}} M A^{-1} B = U U^\mathrm{T}$, then:

$$(G + U U^\mathrm{T})^{-1} = G^{-1} + G^{-1} U (I + U^\mathrm{T} G^{-1} U)^{-1} U^\mathrm{T} G^{-1}$$

where we use that $(I + U^\mathrm{T} G^{-1} U)$ is invertible since $G$ is symmetric positive definite. Such an approach will be efficient only if $l$ is small, since we will need to solve $l$ systems with coefficient matrix $A$ in a preprocessing step.

*Remark 9*
When matrix $M$ is non-singular and its inverse $M^{-1}$ is available, double iteration may also be avoided as follows. Suppose $\alpha_1 > 0$. Define $\boldsymbol{\mu} = -A^{-\mathrm{T}} M A^{-1} B \mathbf{u}$. Then, the following extended block matrix system is easily seen to be equivalent to (20):

$$\begin{bmatrix} A^\mathrm{T} M^{-1} A & B \\ B^\mathrm{T} & -G \end{bmatrix} \begin{bmatrix} \boldsymbol{\mu} \\ \mathbf{u} \end{bmatrix} = \begin{bmatrix} \mathbf{0} \\ \mathbf{g} \end{bmatrix} \tag{23}$$

where the right-hand side $\mathbf{g} = -\mathbf{f}_2 + B^\mathrm{T} A^{-\mathrm{T}} \mathbf{f}_1 - B^\mathrm{T} A^{-\mathrm{T}} M A^{-1} \mathbf{f}_3$ can be computed at an initial overhead cost (requiring the action of $A^{-1}$). The above symmetric *indefinite* system can be transformed into a *symmetric positive definite* system, using a technique described in [23] (without requiring the action of $A^{-1}$) as follows. Suppose $A_0$ is a matrix spectrally equivalent to $A$ (such as a domain decomposition preconditioner), and $M_0 = h^d I$ a suitably scaled matrix spectrally equivalent to $M$, and $G_0 = \alpha_1 h^{d-1} I$ also a suitably scaled matrix equivalent to $G$, such that $A_0^\mathrm{T} M_0^{-1} A_0 \leqslant A^\mathrm{T} M^{-1} A$ (in the sense of quadratic forms). Then, system (23) can be transformed into the following symmetric and positive definite system, see [15, 21, 23]:

$$\begin{bmatrix} K K_0^{-1} K - K & (K K_0^{-1} - I) B \\ B^\mathrm{T} (K_0^{-1} K - I) & G + B^\mathrm{T} K_0^{-1} B \end{bmatrix} \begin{bmatrix} \boldsymbol{\mu} \\ \mathbf{u} \end{bmatrix} = \begin{bmatrix} \mathbf{0} \\ \mathbf{g} \end{bmatrix} \tag{24}$$

where $K = A^\mathrm{T} M^{-1} A$ and $K_0 = A_0^\mathrm{T} M_0^{-1} A_0$. This system can be solved by a PCG algorithm, using a block preconditioner blockdiag($K_0, T_0$), where $K_0$ is as in the preceding, and $T_0$ is spectrally equivalent to $G_0 + B^\mathrm{T} K_0^{-1} B$. In the special case when $J(y, u) = \frac{1}{2} \mathscr{A}(y - \hat{y}, y - \hat{y})$, matrix $A$ will replace $M$ in (23) and we will obtain the simplification $A^\mathrm{T} M^{-1} A = A$, $K_0 = A_0$, and $T_0 = G_0$. We omit further details.

## 3.1. A uniformly effective preconditioner for $C \equiv (G + B^T A^{-T} M A^{-1} B)$

The task of finding an effective preconditioner for the Schur complement $C$ is complicated by the presence of the parameters $\alpha_1 \geqslant 0$ and $\alpha_2 \geqslant 0$. As noted before, when $\alpha_1$ or $\alpha_2$ is large (or equivalently, when $\lambda_{\min}(G)$ is sufficiently large), $G$ will be an effective preconditioner for $C$, while when both $\alpha_1$ and $\alpha_2$ are small (or equivalently, when $\lambda_{\max}(G)$ is sufficiently small), and when $M$ is non-singular, matrix $(B^T A^{-T} M A^{-1} B)$ will be an effective preconditioner for $C$. For intermediate values of $\alpha_i$, however, neither limiting approximation may be effective. In the special case when $\Omega \subset R^2$, we shall indicate a preconditioner uniformly effective with respect to $\alpha_1 > 0$ or $\alpha_2 > 0$. The general case will be considered in a subsequent paper.

The preconditioners that we shall formulate for $C$ will be based on spectrally equivalent representations of $G$ and $(B^T A^{-T} M A^{-1} B)$, for special choices of the matrix $M$. Lemma 3.2 below describes uniform spectral equivalences between $G$, $(B A^{-1} B)$, $(B^T A^{-T} M A^{-1} B)$ and one or more of the matrices $Q$, $S^{-1}$, $S^{-2}$ or $S^{-3}$, where $S = (A_{BB} - A_{IB}^T A_{II}^{-1} A_{IB})$ denotes the discrete Dirichlet to Neumann map. Properties of $S$ have been studied extensively in the domain decomposition literature [31].

*Lemma 3.2*
Let $\Omega \subset R^d$ be a smooth convex domain. Then, the following equivalences will hold:

$$Q \asymp h^{d-1} I$$
$$(B^T A^{-1} B) \asymp Q^T S^{-1} Q$$
$$(B^T A^{-T} M A^{-1} B) \asymp Q^T S^{-1} Q S^{-1} Q \quad \text{when } M = \text{blockdiag}(0, Q)$$
$$(B^T A^{-T} M A^{-1} B) \asymp Q^T S^{-1} Q^T S^{-1} Q S^{-1} Q \quad \text{when } M \asymp h^d I$$

(25)

with coefficients independent of $h$, $\alpha_1$ and $\alpha_2$, where $S = (A_{BB} - A_{IB}^T A_{II}^{-1} A_{IB})$.

*Proof*
The first equivalence is a Gram matrix property on $\partial \Omega$, while the second equivalence follows from $B^T A^{-1} B = Q S^{-1} Q$, proved in Lemma 3.1. To prove the third equivalence, use

$$A^{-1} = \begin{bmatrix} A_{II}^{-1} + A_{II}^{-1} A_{IB} S^{-1} A_{IB} A_{IB} A_{II}^{-1} & -A_{II}^{-1} A_{IB} S^{-1} \\ -S^{-1} A_{IB}^T A_{II}^{-1} & S^{-1} \end{bmatrix}$$

Employing this and using the block matrix structure of $B$ yields

$$A^{-1} B \mathbf{u} = \begin{bmatrix} -A_{II}^{-1} A_{IB} S^{-1} Q \mathbf{u} \\ S^{-1} Q \mathbf{u} \end{bmatrix}$$

Substituting this into $(B^T A^{-T} M A^{-1} B)$ with $M = \text{blockdiag}(0, Q)$ yields

$$B^T A^{-T} M A^{-1} B = Q S^{-1} Q S^{-1} Q$$

and the third equivalence follows. To prove the fourth equivalence, let $u_h$ denote a finite element control function defined on $\partial \Omega$ with associated nodal vector $\mathbf{u}$. Let $v_h$ denote the Dirichlet data associated with the Neumann data $u_h$, i.e. with associated nodal vector $\mathbf{v} = S^{-1} Q \mathbf{u}$. When $M \asymp h^d I$ is the mass matrix on $\Omega$, then $\mathbf{u}^T (B^T A^{-1} M A^{-1} B) \mathbf{u}$ will be equivalent to $\|E v_h\|^2_{L^2(\Omega)}$, where $E v_h$

denotes the *discrete harmonic* extension of the Dirichlet boundary data $v_h$ into $\Omega$ with associated nodal vector $A^{-1}B\mathbf{u}$. When $\Omega$ is *convex* and *smooth*, $H^2(\Omega)$ elliptic regularity will hold for (3) and a result from [32] shows that $\|Ev_h\|^2_{L^2(\Omega)}$ is spectrally equivalent to $\|v_h\|^2_{H^{-1/2}(\partial\Omega)}$. In matrix terms, the nodal vector associated with the discrete Dirichlet data $v_h$ will be $\mathbf{v} = S^{-1}Q\mathbf{u}$, given by the discrete Neumann to Dirichlet map. For $v_h \in H^{-1/2}(\partial\Omega)$, it will hold that $\|v_h\|^2_{H^{-1/2}(\partial\Omega)}$ is spectrally equivalent to $\mathbf{v}^T Q^T S^{-1} Q\mathbf{v}$, in turn equivalent to $\mathbf{u}^T Q^T S^{-1} Q^T S^{-1} Q S^{-1} Q\mathbf{u}$ and the fourth equivalence follows. □

As an immediate corollary, we obtain the following uniform equivalences.

*Corollary 3.3*
When $\Omega$ is a smooth convex domain, the following equivalences will hold:

$$C \asymp \alpha_1 h^{d-1} I + \alpha_2 Q^T S^{-1} Q + Q^T S^{-1} Q S^{-1} Q \quad \text{when } M = \text{blockdiag}(0, Q)$$

$$C \asymp \alpha_1 h^{d-1} I + \alpha_2 Q^T S^{-1} Q + Q^T S^{-1} Q^T S^{-1} Q S^{-1} Q \quad \text{when } M \asymp h^d I$$

(26)

*Proof*
Follows from Lemma 3.2. □

*Remark 10*
When $\partial\Omega$ is *non-smooth*, elliptic regularity results will be weaker and the bounds in Lemma 3.2 may involve poly-logarithmic terms in $h$. However, if the Neumann control is applied only on a smooth *subsegment* $\Gamma \subset \partial\Omega$, the bounds will be independent of $h$.

*3.1.1. A fast Fourier transform (FFT)-based preconditioner for C.* When $\Omega \subset R^2$, matrix $S$ (hence, $Q$, $S^{-1}$, etc.) will have an approximate spectral representation involving the discrete Fourier transform $U$. The Dirichlet to Neumann map $S$ will be spectrally equivalent to an appropriately *scaled* square root of the discretization of the *Laplace–Beltrami* operator $L_B = -d^2/ds_x^2$ on $\partial\Omega$, see [31]. On $\partial\Omega$ and for quasi-uniform triangulation on $\partial\Omega$, the discretization of the Laplace–Beltrami operator with periodic boundary conditions will yield a matrix spectrally equivalent to the circulant matrix $H_0 = h^{-1}\text{circ}(-1, 2, -1)$, since $\partial\Omega$ is a loop. Matrix $H_0$ will be diagonalized by the discrete Fourier transform $U$, yielding $H_0 = U\Lambda_{H_0}U^T$, where $\Lambda_{H_0}$ is a diagonal matrix whose entries can be computed analytically [31]. If $Q$ denotes the mass matrix on $\partial\Omega$, then it will be spectrally equivalent to the circulant matrix $Q_0 \equiv (h/6)\text{circ}(1, 4, 1)$ and diagonalized by the discrete Fourier transform, with $Q_0 = U\Lambda_{Q_0}U^T$. An analytical expression can be derived for the eigenvalues $\Lambda_{Q_0}$, where $(h/3) \leqslant (\Lambda_{Q_0})_i \leqslant h$. Based on the above expressions, we may employ the representations

$$S_0 \asymp Q_0^{1/2}(Q_0^{-1/2}H_0 Q_0^{-1/2})^{1/2}Q_0^{1/2} \equiv U\Lambda_{S_0}U^T = U(\Lambda_{Q_0}^{1/4}\Lambda_{H_0}^{1/2}\Lambda_{Q_0}^{1/4})U^T$$

$$S_0^r \asymp U\Lambda_{S_0}^r U^T \asymp U(\Lambda_{Q_0}^{r/2}\Lambda_{H_0}^{r/2})U^T$$

$$Q \asymp U\Lambda_{Q_0}U^T \asymp hUIU^T$$

The following approximate representations will hold for $C \asymp C_0$:

$$C_0 \asymp U(\alpha_1 \Lambda_{Q_0} + \alpha_2 \Lambda_{Q_0}^2 \Lambda_{S_0}^{-1} + \Lambda_{Q_0}^3 \Lambda_{S_0}^{-2}) U^{\mathrm{T}} \quad \text{when } M = \text{blockdiag}(0, Q)$$

$$C_0 \asymp U(\alpha_1 \Lambda_{Q_0} + \alpha_2 \Lambda_{Q_0}^2 \Lambda_{S_0}^{-1} + \Lambda_{Q_0}^4 \Lambda_{S_0}^{-3}) U^{\mathrm{T}} \quad \text{when } M \asymp h^2 I$$

(27)

The eigenvalues of $C_0^{-1}$ can be found analytically, and the action of $C_0^{-1}$ can be computed at low cost using FFTs. Such preconditioners, however, are not easily generalized to $\Omega \subset R^3$.

*3.1.2. An algebraic preconditioner for $S^{-1}$.* We also describe an algebraic preconditioner $\tilde{S}^{-1}$ for $S^{-1}$, applicable when $\Omega \subset R^2$ or $R^3$. It can precondition $G = \alpha_2 (Q^{\mathrm{T}} S^{-1} Q)$. If $G \leqslant (B^{\mathrm{T}} A^{-\mathrm{T}} M A^{-1} B)$, we may also apply it repeatedly to precondition $C \asymp Q^{\mathrm{T}} S^{-1} Q S^{-1} Q$ or $C \asymp Q^{\mathrm{T}} S^{-1} Q^{\mathrm{T}} S^{-1} Q S^{-1} Q$, depending on whether $M = \text{blockdiag}(0, Q)$ or $M \asymp h^d I$. This preconditioner for $S^{-1}$ is based on a *subregion* $(\Omega \backslash \overline{D}) \subset \Omega$ surrounding $\partial \Omega$. Let $\overline{D} \subset \Omega$ be a subregion with $\text{dist}(\partial D, \partial \Omega) \geqslant \beta > 0$, independent of $h$. Let $\tilde{A}$ denote the submatrix of $A$

$$\tilde{A} = \begin{bmatrix} \tilde{A}_{II} & A_{IB} \\ A_{IB}^{\mathrm{T}} & A_{BB} \end{bmatrix}$$

corresponding to a discretization of the elliptic equation on $\Omega \backslash \overline{D}$ with Neumann boundary conditions on $\partial \Omega$ and zero Dirichlet boundary conditions on $\partial D$. By construction, the matrix $\tilde{S} = (A_{BB} - A_{IB}^{\mathrm{T}} \tilde{A}_{II}^{-1} A_{IB})$ will be spectrally equivalent to $S = (A_{BB} - A_{IB}^{\mathrm{T}} A_{II}^{-1} A_{IB})$, since the Schur complement energy of the discrete harmonic extension into $\Omega \backslash \overline{D}$ will be equivalent to the Schur complement energy of the discrete harmonic extension into $\Omega$ (as both will be equivalent to the $H^{1/2}(\partial \Omega)$ norm square of its boundary data). Applying $\tilde{S}$ will require an exact solver for $\tilde{A}_{II}$ (such as a band solver, if $\beta > 0$ is small).

# 4. PRECONDITIONED-AUGMENTED LAGRANGIAN ALGORITHMS

The second category of algorithms we consider for solving system (13) will avoid double iteration, and correspond to saddle point preconditioners for an *augmented Lagrangian* reformulation [22] of the original system. Traditional saddle point algorithms, such as Uzawa and block preconditioners [14, 15, 18–21, 23], may not be directly applicable to system (13) since matrix $M$ can possibly be *singular*. Instead, in the augmented Lagrangian system, the block submatrix blockdiag$(M, G)$ is transformed into a symmetric positive definite submatrix, so that traditional saddle point methods can be applied. We shall describe preconditioners employing MINRES [14, 18–20] and CG acceleration [15, 21, 23].

Augmenting the Lagrangian [22] is a method suitable for regularizing a saddle point system without altering its solution. Formally, the augmented Lagrangian method seeks the minimum of an augmented energy functional $J_{\text{aug}}(\mathbf{y}, \mathbf{u})$

$$J_{\text{aug}}(\mathbf{y}, \mathbf{u}) \equiv J(\mathbf{y}, \mathbf{u}) + \frac{\rho}{2} \|A\mathbf{y} + B\mathbf{u} - \mathbf{f}_3\|_{A_0^{-1}}^2$$

$$= J(\mathbf{y}, \mathbf{u}) + \frac{\rho}{2} (A\mathbf{y} + B\mathbf{u} - \mathbf{f}_3)^{\mathrm{T}} A_0^{-1} (A\mathbf{y} + B\mathbf{u} - \mathbf{f}_3)$$

subject to the same constraint $A\mathbf{y} + B\mathbf{u} - \mathbf{f}_3 = \mathbf{0}$. Here, matrix $A_0$ will be assumed to be a symmetric positive definite matrix of dimension $n$, spectrally equivalent to $A$, while $\rho \geqslant 0$ is a parameter. By construction, the term $\|A\mathbf{y} + B\mathbf{u} - \mathbf{f}_3\|^2_{A_0^{-1}}$ will be zero in the constraint set, so that the solution of the constrained minimization problem is unaltered. Defining an augmented Lagrangian functional $\mathscr{L}_{\text{aug}}(\mathbf{y}, \mathbf{u}, \mathbf{p})$:

$$\mathscr{L}_{\text{aug}}(\mathbf{y}, \mathbf{u}, \mathbf{p}) = J_{\text{aug}}(\mathbf{y}, \mathbf{u}) + \mathbf{p}^{\text{T}}(A\mathbf{y} + B\mathbf{u} - \mathbf{f}_3) \tag{28}$$

and seeking its saddle point will yield the following modified saddle point system:

$$\begin{bmatrix} M + \rho A^{\text{T}} A_0^{-1} A & \rho A^{\text{T}} A_0^{-1} B & A^{\text{T}} \\ \rho B^{\text{T}} A_0^{-1} A & G + \rho B^{\text{T}} A_0^{-1} B & B^{\text{T}} \\ A & B & 0 \end{bmatrix} \begin{bmatrix} \mathbf{y} \\ \mathbf{u} \\ \mathbf{p} \end{bmatrix} = \begin{bmatrix} \mathbf{f}_1 + \rho A^{\text{T}} A_0^{-1} \mathbf{f}_3 \\ \mathbf{f}_2 + \rho B^{\text{T}} A_0^{-1} \mathbf{f}_3 \\ \mathbf{f}_3 \end{bmatrix} \tag{29}$$

The above system can alternatively be obtained from (13) by multiplying the third block row of (13) by $\rho A^{\text{T}} A_0^{-1}$ and adding it to the first block row, and multiplying the third block row of (13) by $\rho B^{\text{T}} A_0^{-1}$ and adding it to the second block row.

To simplify our discussion, we shall employ the notation:

$$K \equiv \begin{bmatrix} M + \rho A^{\text{T}} A_0^{-1} A & \rho A^{\text{T}} A_0^{-1} B \\ \rho B^{\text{T}} A_0^{-1} A & G + \rho B^{\text{T}} A_0^{-1} B \end{bmatrix}, \quad N^{\text{T}} \equiv \begin{bmatrix} A^{\text{T}} \\ B^{\text{T}} \end{bmatrix}, \quad \mathbf{w} \equiv \begin{bmatrix} \mathbf{y} \\ \mathbf{u} \end{bmatrix} \tag{30}$$

Using this, the augmented saddle point system can be represented compactly as

$$\begin{bmatrix} K & N^{\text{T}} \\ N & 0 \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ \mathbf{p} \end{bmatrix} = \begin{bmatrix} \mathbf{f} \\ \mathbf{g} \end{bmatrix} \tag{31}$$

where $\mathbf{f} = ((\mathbf{f}_1 + \rho A^{\text{T}} A_0^{-1} \mathbf{f}_3)^{\text{T}}, (\mathbf{f}_2 + \rho B^{\text{T}} A_0^{-1} \mathbf{f}_3)^{\text{T}})^{\text{T}}$ and $\mathbf{g} = \mathbf{f}_3$. This coefficient matrix is *symmetric indefinite*, and we shall consider two algorithms for solving it using a block diagonal preconditioner of the form $\text{blockdiag}(K_0, T_0)$, where $K_0$ and $T_0$ are matrices spectrally equivalent to $K$ and $T = N K^{-1} N^{\text{T}}$, respectively.

Our first augmented Lagrangian method will solve (31) using the MINRES algorithm with $\text{blockdiag}(K_0, T_0)$ as a preconditioner. Our second method will transform (31) into a *symmetric positive definite* system [23], and solve it using the CG algorithm. Analysis of system (31) with preconditioner $\text{blockdiag}(K_0, T_0)$ shows that effective MINRES or CG algorithms can be formulated, provided $K_0$ and $T_0$ are spectrally equivalent to $K$ and $T = N K^{-1} N^{\text{T}}$, respectively, [14, 15, 18–21, 23]. We now consider $\text{blockdiag}(A^{\text{T}} A_0^{-1} A, G)$ as a preconditioner for $K$.

*Lemma 4.1*
Let $G$ be positive definite, and suppose the following hold.

1. Let $\mathbf{y}^{\text{T}} M \mathbf{y} \leqslant \gamma_1 (\mathbf{y}^{\text{T}} A^{\text{T}} A_0^{-1} A \mathbf{y})$ for some $\gamma_1 > 0$ independent of $h$.
2. Let $\mathbf{v}^{\text{T}} (B^{\text{T}} A_0^{-1} B) \mathbf{v} \leqslant \gamma_2 (\mathbf{v}^{\text{T}} G \mathbf{v})$ for some $\gamma_2 > 0$ independent of $h$.
3. Let $\beta_* \equiv \frac{1}{2}((2 + \gamma_2) - \sqrt{\gamma_2^2 + 4\gamma_2})$ and $\beta_{**} \equiv \frac{1}{2}((2 + \gamma_1 + \gamma_2) + \sqrt{(\gamma_1 - \gamma_2)^2 + 4\gamma_2})$.

Then, for $(\mathbf{y}^T, \mathbf{u}^T)^T \neq \mathbf{0}$, the following bounds will hold:

$$\beta_* \leqslant \frac{\begin{bmatrix} \mathbf{y} \\ \mathbf{u} \end{bmatrix}^T \begin{bmatrix} M + A^T A_0^{-1} A & A^T A_0^{-1} B \\ B^T A_0^{-1} A & G + B^T A_0^{-1} B \end{bmatrix} \begin{bmatrix} \mathbf{y} \\ \mathbf{u} \end{bmatrix}}{\begin{bmatrix} \mathbf{y} \\ \mathbf{u} \end{bmatrix}^T \begin{bmatrix} A^T A_0^{-1} A & 0 \\ 0 & G \end{bmatrix} \begin{bmatrix} \mathbf{y} \\ \mathbf{u} \end{bmatrix}} \leqslant \beta_{**} \tag{32}$$

*Proof*
Expand the quadratic form associated with the block matrix

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{u} \end{bmatrix}^T \begin{bmatrix} M + A^T A_0^{-1} A & A^T A_0^{-1} B \\ B^T A_0^{-1} A & G + B^T A_0^{-1} B \end{bmatrix} \begin{bmatrix} \mathbf{y} \\ \mathbf{u} \end{bmatrix}$$

$$= (\mathbf{y}^T M \mathbf{y} + \mathbf{y}^T A^T A_0^{-1} A \mathbf{y} + \mathbf{u}^T G \mathbf{u} + \mathbf{u}^T B^T A_0^{-1} B \mathbf{u}) + 2\mathbf{y}^T A^T A_0^{-1} B \mathbf{u} \tag{33}$$

and employ Schwarz's inequality, using the identity $2ab \leqslant \theta a^2 + b^2/\theta$ for $0 < \theta < \infty$

$$2|\mathbf{y}^T A^T A_0^{-1} B \mathbf{u}| \leqslant \theta \mathbf{y}^T A^T A_0^{-1} A \mathbf{y} + \frac{1}{\theta} \mathbf{u}^T B^T A_0^{-1} B \mathbf{u} \tag{34}$$

To obtain an *upper bound*, substitute (34) into (33) and choose $\theta$

$$(\mathbf{y}^T M \mathbf{y} + \mathbf{y}^T A^T A_0^{-1} A \mathbf{y} + \mathbf{u}^T G \mathbf{u} + \mathbf{u}^T B^T A_0^{-1} B \mathbf{u}) + 2\mathbf{y}^T A^T A_0^{-1} B \mathbf{u}$$

$$\leqslant (1 + \gamma_1 + \theta) \mathbf{y}^T A^T A_0^{-1} A \mathbf{y} + \left(1 + \gamma_2 + \frac{\gamma_2}{\theta}\right) \mathbf{u}^T G \mathbf{u} \leqslant \beta_{**}(\mathbf{y}^T A^T A_0^{-1} A \mathbf{y} + \mathbf{u}^T G \mathbf{u})$$

by equating $\beta_{**} \equiv (1 + \gamma_1 + \theta) = (1 + \gamma_2 + \gamma_2/\theta)$. To obtain a *lower bound*, expand (33) as follows:

$$(\mathbf{y}^T M \mathbf{y} + \mathbf{y}^T A^T A_0^{-1} A \mathbf{y} + \mathbf{u}^T G \mathbf{u} + \mathbf{u}^T B^T A_0^{-1} B \mathbf{u}) + 2\mathbf{y}^T A^T A_0^{-1} B \mathbf{u}$$

$$\geqslant (\mathbf{y}^T A^T A_0^{-1} A \mathbf{y} + \mathbf{u}^T G \mathbf{u} + \mathbf{u}^T B^T A_0^{-1} B \mathbf{u}) - 2|\mathbf{y}^T A^T A_0^{-1} B \mathbf{u}|$$

$$\geqslant (\mathbf{y}^T A^T A_0^{-1} A \mathbf{y} + \mathbf{u}^T G \mathbf{u} + \mathbf{u}^T B^T A_0^{-1} B \mathbf{u}) - (1 - \tilde{\theta})(\mathbf{y}^T A^T A_0^{-1} A \mathbf{y})$$

$$- \frac{1}{1 - \tilde{\theta}} \mathbf{u}^T B^T A_0^{-1} B \mathbf{u}$$

$$\geqslant \tilde{\theta}(\mathbf{y}^T A^T A_0^{-1} A \mathbf{y}) + \mathbf{u}^T G \mathbf{u} - \frac{\tilde{\theta}}{1 - \tilde{\theta}} \mathbf{u}^T B^T A_0^{-1} B \mathbf{u}$$

$$\geqslant \tilde{\theta}(\mathbf{y}^{\mathrm{T}} A^{\mathrm{T}} A_0^{-1} A \mathbf{y}) + \frac{1 - (1 + \gamma_2)\tilde{\theta}}{1 - \tilde{\theta}} \mathbf{u}^{\mathrm{T}} G \mathbf{u}$$

$$\geqslant \beta_*(\mathbf{y}^{\mathrm{T}} A^{\mathrm{T}} A_0^{-1} A \mathbf{y} + \mathbf{u}^{\mathrm{T}} G \mathbf{u})$$

where we require $0 < \tilde{\theta} < 1$ and such that $\beta_* \equiv \tilde{\theta} = (1 - (1 + \gamma_2)\tilde{\theta})/(1 - \tilde{\theta})$. $\qquad \square$

*Remark 11*

The upper bound $\beta_{**}$ in (32) can be replaced by $\max\{2 + \gamma_1, 2\gamma_2 + 1\}$. This simpler upper bound can be derived by substituting (34) into (33) and choosing $\theta = 1$. Similarly, the lower bound $\beta_*$ in (32) can be replaced by $\min\{1/(2 + 2\gamma_2), (1 + \gamma_2)/(2\gamma_2 + 1)\}$. This can be derived by choosing $\tilde{\theta} = 1/(2 + 2\gamma_2)$. Replacing $\beta_{**}$ by $\max\{2 + \gamma_1, 2\gamma_2 + 1\}$ will lead to a more tractable expression for the optimal parameter $\rho_{\mathrm{opt}}$ which minimizes the condition number in (32), when a scaling parameter $\rho > 0$ is introduced in the augmented Lagrangian formulation and $A_0^{-1}$ is replaced by $\rho A_0^{-1}$ in Lemma 4.1. If $\gamma_1$, $\gamma_2$ are as defined earlier (corresponding to the choice $\rho = 1$), then the following bounds will hold when $(\mathbf{y}^{\mathrm{T}}, \mathbf{u}^{\mathrm{T}})^{\mathrm{T}} \neq \mathbf{0}$:

$$\beta_*(\rho) \leqslant \frac{\begin{bmatrix} \mathbf{y} \\ \mathbf{u} \end{bmatrix}^{\mathrm{T}} \begin{bmatrix} M + \rho A^{\mathrm{T}} A_0^{-1} A & \rho A^{\mathrm{T}} A_0^{-1} B \\ \rho B^{\mathrm{T}} A_0^{-1} A & G + \rho B^{\mathrm{T}} A_0^{-1} B \end{bmatrix} \begin{bmatrix} \mathbf{y} \\ \mathbf{u} \end{bmatrix}}{\begin{bmatrix} \mathbf{y} \\ \mathbf{u} \end{bmatrix}^{\mathrm{T}} \begin{bmatrix} \rho A^{\mathrm{T}} A_0^{-1} A & 0 \\ 0 & G \end{bmatrix} \begin{bmatrix} \mathbf{y} \\ \mathbf{u} \end{bmatrix}} \leqslant \max\left\{2 + \frac{\gamma_1}{\rho}, 2\gamma_2\rho + 1\right\} \quad (35)$$

where $\beta_*(\rho) = \frac{1}{2}((2 + \gamma_2\rho) - \sqrt{\gamma_2^2\rho^2 + 4\gamma_2\rho})$. Denote the condition number in (35) by $\kappa(\rho)$

$$\kappa(\rho) \equiv \frac{\max\{2 + (\gamma_1/\rho), 2\gamma_2\rho + 1\}}{\beta_*(\rho)}$$

The scaling parameter $\rho = \rho_{\mathrm{opt}}$ can be chosen to minimize the above condition number. It can be verified easily that $(2 + \gamma_1/\rho)/\beta_*(\rho)$ intersects $(2\gamma_2\rho + 1)/\beta_*(\rho)$ when $\rho = \rho_*$

$$\rho_* = \frac{1 + \sqrt{1 + 8\gamma_1\gamma_2}}{4\gamma_2}$$

Furthermore, it can be verified that $(2\gamma_2\rho + 1)/\beta_*(\rho)$ is monotonically increasing for $\rho \geqslant \rho_*$. Consequently, the *optimal* choice $\rho_{\mathrm{opt}}$ of parameter $\rho$ will occur for $\rho \in (0, \rho_*]$

$$\kappa(\rho_{\mathrm{opt}}) = \min_{0 < \rho \leqslant \rho_*} \kappa(\rho)$$

An explicit expression can be derived for $\rho_{\mathrm{opt}}$ (using Maple or Mathematica), however, we shall omit the resulting expression, since it is lengthy.

*Remark 12*

The limiting case $\rho = 0$ yields the original saddle point system (13). If $M$ is non-singular, then choosing a preconditioner $K_0$ for $K = \text{blockdiag}(M, G)$ will be simple. If $\Omega_0 = \Omega$, we may choose $K_0 = \text{blockdiag}(M_0, G_0)$, where $M_0$ and $G_0$ are spectrally equivalent to $M$ and $G$, respectively.

*Remark 13*

In applications to control system (29), matrix $A_0$ can be chosen as a preconditioner spectrally equivalent to $A$. Then, $A_0^{-1}$ will be spectrally equivalent to $A^{-1}$, and $A^{\mathrm{T}} A_0^{-1} A$ will be spectrally equivalent to $A$. An application of Poincare–Freidrichs inequality will yield the bound assumed in Lemma 4.1 with $\gamma_1$ independent of $h$. Furthermore, when $A_0$ is spectrally equivalent to $A$, it will also hold that $B A_0^{-1} B^{\mathrm{T}}$ is spectrally equivalent to $Q^{\mathrm{T}} S^{-1} Q$, and the arguments employed in Lemma 3.1 will yield the bound assumed in Lemma 4.1 with $\gamma_2$ independent of $h$.

We shall consider two approaches which employ a preconditioner $K_0$ for $K$ to precondition the augmented Lagrangian saddle point system. The first approach, described in Section 4.1, solves the augmented saddle point system using the MINRES algorithm with a block diagonal preconditioner. The second approach, described in Section 4.2, reformulates the augmented saddle point system as a symmetric positive definite system and solves it using a CG algorithm.

### 4.1. Minimum residual acceleration

Consider now the solution of system (31) for $\rho > 0$, with $K$, $N$, **w**, **p** defined by (30). Since system (31) is symmetric but *indefinite*, the CG algorithm cannot be employed to solve it. Instead, our first method employs the MINRES algorithm [17, 18] for symmetric indefinite systems.

Typically, the rate of convergence of the MINRES algorithm to solve a saddle point system depends on the intervals $[-d, -c]$ and $[a, b]$ containing the negative and positive eigenvalues of the preconditioned system [6, 15, 18–21]. Theoretical convergence bounds for the MINRES algorithm are generally weaker than that for the CG algorithm, however, its rate of convergence will be independent of a parameter provided the intervals containing the eigenvalues are fixed and bounded away from zero, independent of the same parameter. In particular, if a symmetric positive definite preconditioner of the form $\text{blockdiag}(K_0, T_0)$ is employed to precondition (31), and $K_0$ and $T_0$ are spectrally equivalent to $K$ and $T = N K^{-1} N^{\mathrm{T}}$, respectively, independent of a parameter, then the rate of convergence of the preconditioned MINRES algorithm will also be independent of those parameters [15, 16, 18–21, 33]. The next result considers $K_0 = \text{blockdiag}(\rho A_0, G_0)$ as a preconditioner for $K$ and a matrix $A_*$ spectrally equivalent to $T_0 = \rho^{-1} A + B G^{-1} B^{\mathrm{T}}$.

*Lemma 4.2*

Suppose the following conditions hold.

1. Let $A_0$ be spectrally equivalent to $A$, independent of $h$.
2. Let $G_0$ be spectrally equivalent to $G$, independent of $h$.
3. Let $A_*$ be spectrally equivalent to $\rho^{-1} A + B G^{-1} B^{\mathrm{T}}$, independent of $h$.

Then, the rate of convergence of the MINRES algorithm to solve (31) using the preconditioner $L_0 = \text{blockdiag}(\rho A_0, G_0, A_*)$ will be independent of $h$ (but not $\alpha_1, \alpha_2$) for $\rho > 0$.

*Proof*

When $A_0$ is spectrally equivalent to $A$, an application of Lemma 4.1 will yield $\gamma_1$ and $\gamma_2$ to be independent of $h$, and blockdiag$(\rho A^T A_0^{-1} A, G)$ to be spectrally equivalent to $K$, independent of $h$. Matrix $A^T A_0^{-1} A$ will be spectrally equivalent to $A$ and to $A_0$, thus, replacing $A^T A_0^{-1} A$ by $A_0$ and $G$ by $G_0$ will yield that $K_0 = \text{blockdiag}(\rho A_0, G_0)$ to be spectrally equivalent to $K$, independent of $h$. Spectral equivalence between $K$ and $K_0$ immediately yields spectral equivalence between $T = N K^{-1} N^T$ and $N K_0^{-1} N^T$. Substituting $K_0 = \text{blockdiag}(\rho A_0, G_0)$ into $N K_0^{-1} N^T$ yields $\rho^{-1} A^T A_0^{-1} A + B^T G_0^{-1} B$, which is spectrally equivalent to $\rho^{-1} A + B^T G^{-1} B$. Analysis of saddle point algorithms show that the rate of convergence of iterative algorithms to solve a system of form (31) using a preconditioner blockdiag$(K_0, T_0)$, will be independent of a parameter, provided $K_0$ and $T_0$ are spectrally equivalent to $K$ and $T = N K^{-1} N^T$, independent of that parameter. Thus, it will be sufficient to require $A_*$ to be spectrally equivalent to $\rho^{-1} A + B^T G^{-1} B$.

$\square$

*Remark 14*

Each application $L_0^{-1}$ of $L_0 = \text{blockdiag}(\rho A_0, G_0, A_*)$ will require the action of $A_0^{-1}$ once, $G_0^{-1}$ once and $A_*^{-1}$ once. Each multiplication by $K$ can be computed using

$$K \begin{bmatrix} \mathbf{y} \\ \mathbf{u} \end{bmatrix} = \begin{bmatrix} M\mathbf{y} \\ G\mathbf{u} \end{bmatrix} + \begin{bmatrix} A^T \\ B^T \end{bmatrix} A_0^{-1}(A^T\mathbf{y} + B\mathbf{u}) \tag{36}$$

This will require the action of $A_0^{-1}$ *once*.

*Remark 15*

Matrix $A_*$ should be spectrally equivalent to $\rho^{-1} A + B^T G^{-1} B$. When $G = \alpha_1 Q$, this requires matrix $A_*$ to be spectrally equivalent to

$$A_* \asymp \rho^{-1} \begin{bmatrix} A_{II} & A_{IB} \\ A_{IB}^T & \dfrac{\rho}{\alpha_1} Q + A_{BB} \end{bmatrix}$$

This will correspond to a discretization of a scaled Laplacian with Robin boundary conditions on $\partial\Omega$. In this case, any suitable Robin preconditioner $A_*$ (using domain decomposition, for instance) can be employed. When $G = \alpha_2 Q^T S^{-1} Q$, matrix $A_*$ will be required to satisfy

$$A_* \asymp \rho^{-1} \begin{bmatrix} A_{II} & A_{IB} \\ A_{IB}^T & A_{BB} \end{bmatrix} + \alpha_2^{-1} \begin{bmatrix} 0 & 0 \\ 0 & S \end{bmatrix} = \rho^{-1} \begin{bmatrix} A_{II} & A_{IB} \\ A_{IB}^T & \dfrac{\rho}{\alpha_2} S + A_{BB} \end{bmatrix}$$

where $S = (A_{BB} - A_{IB}^T A_{II}^{-1} A_{IB})$. Since $A^T = A > 0$, it will hold that $S^T = S > 0$. The following algebraic property can also be shown to hold (in the sense of quadratic forms):

$$\begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \leqslant \begin{bmatrix} 0 & 0 \\ 0 & S \end{bmatrix} \leqslant \begin{bmatrix} A_{II} & A_{IB} \\ A_{IB}^T & A_{BB} \end{bmatrix}$$

As a result, it will hold that

$$\rho^{-1} \begin{bmatrix} A_{II} & A_{IB} \\ A_{IB}^{\mathrm{T}} & A_{BB} \end{bmatrix} \leqslant \rho^{-1} \begin{bmatrix} A_{II} & A_{IB} \\ A_{IB}^{\mathrm{T}} & A_{BB} + \dfrac{\rho}{\alpha_2} S \end{bmatrix} \leqslant (\rho^{-1} + \alpha_2^{-1}) \begin{bmatrix} A_{II} & A_{IB} \\ A_{IB}^{\mathrm{T}} & A_{BB} \end{bmatrix}$$

Thus, it is sufficient that $A_*$ be spectrally equivalent to the Neumann matrix $A$.

## 4.2. Conjugate gradient acceleration

The CG method cannot be applied directly to solve system (31), since it is symmetric *indefinite*. However, it is shown in [23] that a general saddle point system of the form (31) can be transformed into an equivalent *symmetric positive definite* system. This resulting system may be solved by the CG method. We shall describe the transformation below. Let $K_0$ denote a symmetric positive definite preconditioner for $K$ satisfying

$$\varepsilon_1 K \leqslant K_0 \leqslant \varepsilon_2 K, \quad \text{for } 0 < \varepsilon_1 \leqslant \varepsilon_2 < 1$$

independent of $h$. Then, a symmetric positive definite system equivalent to (31) is

$$\begin{bmatrix} K^{\mathrm{T}} K_0^{-1} K - K & (K^{\mathrm{T}} K_0^{-1} - I) N^{\mathrm{T}} \\ N(K_0^{-1} K - I) & N K_0^{-1} N^{\mathrm{T}} \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ \mathbf{p} \end{bmatrix} = \begin{bmatrix} (K^{\mathrm{T}} K_0^{-1} - I)\mathbf{f} \\ N K_0^{-1} \mathbf{f} - \mathbf{g} \end{bmatrix} \tag{37}$$

The coefficient matrix $L$ in (37) can be shown to be spectrally equivalent to $L_0$ below

$$L = \begin{bmatrix} K^{\mathrm{T}} K_0^{-1} K - K & (K^{\mathrm{T}} K_0^{-1} - I) N^{\mathrm{T}} \\ N(K_0^{-1} K - I) & N K_0^{-1} N^{\mathrm{T}} \end{bmatrix} \quad \text{and} \quad L_0 = \begin{bmatrix} K_0 & 0 \\ 0 & T_0 \end{bmatrix} \tag{38}$$

where $T_0$ is any matrix spectrally equivalent to $T = N K^{-1} N^{\mathrm{T}}$, see [15, 21, 23]. We may thus obtain the solution to (31) by solving (37) employing the CG method, with $L_0$ as a preconditioner. The resulting rate of convergence will be independent of $h$.

As before, Lemma 4.1 suggests how to construct a symmetric positive definite preconditioner $K_0$ for $K$, satisfying $K_0 \leqslant \varepsilon_2 K$. Suppose $A_0$ is spectrally equivalent to $A$, additionally satisfying $A_0 \leqslant A$. Then, $A_0$ and $A^{\mathrm{T}} A_0^{-1} A$ will also be spectrally equivalent, with $A_0 \leqslant A \leqslant A^{\mathrm{T}} A_0^{-1} A$. If $G_0$ is spectrally equivalent to $G$, satisfying $G_0 \leqslant G$, then Lemma 4.1 will yield

$$K_0 \equiv \varepsilon_2 \theta_*(\rho) \begin{bmatrix} \rho A_0 & 0 \\ 0 & G_0 \end{bmatrix} \leqslant \varepsilon_2 K = \varepsilon_2 \begin{bmatrix} M + \rho A^{\mathrm{T}} A_0^{-1} A & \rho A^{\mathrm{T}} A_0^{-1} B \\ \rho B^{\mathrm{T}} A_0^{-1} A & G + \rho B^{\mathrm{T}} A_0^{-1} B \end{bmatrix} \tag{39}$$

Thus, once spectrally equivalent matrices $A_0 \leqslant A$ and $G_0 \leqslant G$ have been chosen, and parameter $\theta_*(\rho)$ has been estimated, the preconditioner $K_0$ defined by (39) can be employed to transform indefinite system (31) into the symmetric positive definite system (37). The CG method can be employed to solve system (37) with spectrally equivalent preconditioner $L_0$ defined by (38).

*Remark 16*
In practical implementations of the CG method to solve (37), the matrix–vector product with $L$ can be computed as follows, when $\mathbf{w} = (\mathbf{y}^{\mathrm{T}}, \mathbf{u}^{\mathrm{T}})^{\mathrm{T}}$:

$$L \begin{bmatrix} \mathbf{w} \\ \mathbf{p} \end{bmatrix} = \begin{bmatrix} (K^{\mathrm{T}} K_0^{-1} - I)(K\mathbf{w} + N^{\mathrm{T}}\mathbf{p}) \\ N[K_0^{-1}(K\mathbf{w} + N^{\mathrm{T}}\mathbf{p}) - \mathbf{w}] \end{bmatrix}$$

Each matrix–vector product with $K$ can be computed as in (36), requiring the action of $A_0^{-1}$ *once*. Each matrix–vector multiplication with $L$ will require two matrix–vector multiplications with $K$ and one matrix–vector multiplications with $K_0^{-1}$ (which together require the action of $A_0^{-1}$ three times). Also note that the action of inverse $L_0^{-1}$ of a block diagonal preconditioner requires the action of $A_0^{-1}$ and $A_*^{-1}$. Thus, each iteration of this preconditioned CG algorithm will require the action of $A_0^{-1}$ four times and that of $A_*^{-1}$ once, per iteration. For alternative efficient implementations, see [33].

The following bounds can be proved for the resulting CG algorithm.

*Lemma 4.3*
Suppose the following conditions hold.

1. Let $A_0$ be spectrally equivalent to $A$ satisfying $A_0 \leqslant A$.
2. Let $G_0 \leqslant G$ be spectrally equivalent to $G$.
3. Let $A_*$ be spectrally equivalent to $(1/\varepsilon_2 \theta_*(\rho))(\rho^{-1} A + B G^{-1} B^{\mathrm{T}})$.

Then, for the choice $K_0 = \varepsilon_2\, \theta_*(\rho)\mathrm{blockdiag}(\rho A_0, G_0)$, the matrix:

$$L_0 = \mathrm{blockdiag}(K_0, A_*) = \begin{bmatrix} K_0 & 0 \\ 0 & A_* \end{bmatrix}$$

will be spectrally equivalent to $L$ in (37), independent of $h$.

*Proof*
By construction, matrix $K_0$ will satisfy $K_0 \leqslant \varepsilon_2 K$ and be spectrally equivalent to $K$. As a result, by [15, 21, 23], the coefficient matrix $L$ in (37) will be symmetric positive definite, and spectrally equivalent to $L_0 = \mathrm{blockdiag}(K_0, N K_0^{-1} N^{\mathrm{T}})$. For the special choice $K_0 = \varepsilon_2\, \theta_*(\rho)\mathrm{blockdiag}(\rho A_0, G_0) \leqslant \varepsilon_2 K$, we obtain

$$(N K_0^{-1} N^{\mathrm{T}}) = \frac{1}{\varepsilon_2 \theta_*(\rho)}(\rho^{-1} A A_0^{-1} A^{\mathrm{T}} + B G_0^{-1} B^{\mathrm{T}})$$

where $A A_0^{-1} A^{\mathrm{T}}$ is spectrally equivalent to $A_0$, and $B G_0^{-1} B^{\mathrm{T}}$ has block structure

$$B G_0^{-1} B^{\mathrm{T}} = \begin{bmatrix} 0 \\ Q \end{bmatrix} G_0^{-1} \begin{bmatrix} 0 \\ Q \end{bmatrix}^{\mathrm{T}} = \begin{bmatrix} 0 & 0 \\ 0 & Q G_0^{-1} Q^{\mathrm{T}} \end{bmatrix}$$

Thus, $\rho^{-1}AA_0^{-1}A^{\mathrm{T}} + BG_0^{-1}B^{\mathrm{T}}$ will be spectrally equivalent to

$$\rho^{-1}A + BG^{-1}B^{\mathrm{T}} = \rho^{-1}\begin{bmatrix} A_{II} & A_{IB} \\ A_{IB}^{\mathrm{T}} & A_{BB} + \rho QG_0^{-1}Q^{\mathrm{T}} \end{bmatrix}$$

The desired result follows when $A_*$ is spectrally equivalent to $(1/\varepsilon_2\theta_*(\rho))(\rho^{-1}A + BG^{-1}B^{\mathrm{T}})$.
□

# 5. ALTERNATIVE APPROACHES

In this section, we shall describe two alternative *heuristic* approaches to solving (13) when matrix $M$ is *singular*. One approach describes a projected gradient method, as in [24], without the use of the augmented Lagrangian formulation or the reduced Schur complement system for **u**. In another approach, a non-symmetric block matrix preconditioner is proposed for (13) and accelerated by GMRES.

## 5.1. Projected gradient method

Suppose that $M$ is a matrix of rank $(n - k)$, where $\dim(\mathrm{Kernel}(M)) = k$. Let $H$ denote a matrix of dimension $n \times k$, whose columns form a basis for the *null space* of $M$

$$\mathrm{Range}(H) = \mathrm{Kernel}(M) \subset R^n$$

When $M$ is singular, the first block row of (13) will be solvable only when

$$H^{\mathrm{T}}(\mathbf{f}_1 - A^{\mathrm{T}}\mathbf{p}) = \mathbf{0} \Longrightarrow \mathbf{y} = M^{\dagger}(\mathbf{f}_1 - A^{\mathrm{T}}\mathbf{p}) + H\boldsymbol{\alpha}$$

where $\boldsymbol{\alpha} \in R^k$. Here, $M^{\dagger}$ denotes the Moore–Penrose pseudoinverse of $M$. Formally solving the second block row for **u**, yields

$$\mathbf{u} = G^{-1}(\mathbf{f}_2 - B^{\mathrm{T}}\mathbf{p})$$

Formally substituting the preceding two expressions for **y** and **u** into the third block row yields the following reduced system for **p**, together with $H^{\mathrm{T}}(\mathbf{f}_1 - A^{\mathrm{T}}\mathbf{p}) = \mathbf{0}$, the requirement for consistency of the first block row

$$AH\boldsymbol{\alpha} - (AM^{\dagger}A^{\mathrm{T}} + BG^{-1}B^{\mathrm{T}})\mathbf{p} = \mathbf{f}_3 - AM^{\dagger}\mathbf{f}_1 - BG^{-1}\mathbf{f}_2$$

$$H^{\mathrm{T}}A^{\mathrm{T}}\mathbf{p} = H^{\mathrm{T}}\mathbf{f}_1$$

Define the following Euclidean orthogonal *projection* $P_0 = AH(H^{\mathrm{T}}A^{\mathrm{T}}AH)^{-1}H^{\mathrm{T}}A^{\mathrm{T}}$ onto Range $(AH)$ (where Range$(AH)$ has dimension $k$). Applying $(I - P_0)$ to the preceding system, and noting that $(I - P_0)AH\boldsymbol{\alpha} = \mathbf{0}$ yields:

$$(I - P_0)(AM^{\dagger}A^{\mathrm{T}} + BG^{-1}B^{\mathrm{T}})\mathbf{p} = -(I - P_0)(\mathbf{f}_3 - AM^{\dagger}\mathbf{f}_1 - BG^{-1}\mathbf{f}_2)$$

together with the constraint $H^{\mathrm{T}}A^{\mathrm{T}}\mathbf{p} = H^{\mathrm{T}}\mathbf{f}_1$. We may decompose

$$\mathbf{p} = \mathbf{p}_* + \tilde{\mathbf{p}} \quad \text{where } H^{\mathrm{T}}A^{\mathrm{T}}\mathbf{p}_* = H^{\mathrm{T}}\mathbf{f}_1$$

so that $H^T A^T \tilde{\mathbf{p}} = \mathbf{0}$. The term $\mathbf{p}_*$ can be sought as $\mathbf{p}_* = A H \gamma_*$ for some $\gamma_* \in R^k$. This will yield $\mathbf{p}_* = A H (H^T A^T A H)^{-1} H^T \mathbf{f}_1$, and the following system for $\tilde{\mathbf{p}}$:

$$(I - P_0)(A M^\dagger A^T + B G^{-1} B^T)\tilde{\mathbf{p}} = \tilde{\mathbf{g}}$$

where $\tilde{\mathbf{g}} \equiv (I - P_0)(A M^\dagger \mathbf{f}_1 + B G^{-1} \mathbf{f}_2 - \mathbf{f}_3 - (A M^\dagger A^T + B G^{-1} B^T)\mathbf{p}_*)$. Since $H^T A^T \tilde{\mathbf{p}} = \mathbf{0}$, it will formally hold that $(I - P_0)\tilde{\mathbf{p}} = \tilde{\mathbf{p}}$, so that we may solve the system for $\tilde{\mathbf{p}}$ using a CG algorithm. A preconditioner $T$ may be employed, such that the action of its inverse is given by

$$T^{-1} \equiv (I - P_0)(A_0 M^\dagger A_0^T + B G^{-1} B^T)^{-1}(I - P_0)$$

The term $(A M^\dagger A^T + B G^{-1} B^T)^{-1}$ may be replaced by $A_0^{-T} M A_0^{-1}$

## 5.2. Block preconditioner

Another approach to solve (13) is to precondition this system by a non-symmetric block matrix preconditioner $L_0$, as described below, and to use GMRES acceleration [17]:

$$L_0 = \begin{bmatrix} M & 0 & A^T \\ 0 & G & B^T \\ A & 0 & 0 \end{bmatrix}$$

It is easily verified that block matrix $L_0$ is easily inverted. A heuristic analysis of the eigenvalues of $L_0^{-1} L$ (where $L$ denotes the original symmetric, indefinite saddle point coefficient matrix), can be obtained by analysing:

$$\begin{bmatrix} M & 0 & A^T \\ 0 & G & B^T \\ A & B & 0 \end{bmatrix} \begin{bmatrix} \mathbf{y} \\ \mathbf{u} \\ \mathbf{p} \end{bmatrix} = \lambda \begin{bmatrix} M & 0 & A^T \\ 0 & G & B^T \\ A & 0 & 0 \end{bmatrix} \begin{bmatrix} \mathbf{y} \\ \mathbf{u} \\ \mathbf{p} \end{bmatrix} \qquad (40)$$

The eigenvalues may be estimated as follows. From (40) we obtain

$$(1 - \lambda)M\mathbf{y} + (1 - \lambda)A^T\mathbf{p} = \mathbf{0}$$
$$(1 - \lambda)G\mathbf{u} + (1 - \lambda)B^T\mathbf{p} = \mathbf{0} \qquad (41)$$
$$(1 - \lambda)A\mathbf{y} + B\mathbf{u} = \mathbf{0}$$

If $\lambda = 1$, then $B\mathbf{u} = \mathbf{0}$, with $\mathbf{y}$ and $\mathbf{p}$ arbitrary, yielding $\mathbf{u} = \mathbf{0}$. The eigenspace will have dimension $2n$. If $\lambda \neq 1$, then:

$$(1 - \lambda)G\mathbf{u} = (B^T A^{-T} M A^{-1} B)\mathbf{u} \qquad (42)$$

If $\mathrm{Range}(A^{-1} B) \cap \mathrm{Kernel}(M) = \{0\}$, we will obtain $c_1 h^r \leqslant (1 - \lambda) \leqslant c_2 < 1$ for some $c_1, c_2 > 0$ independent of $h$. This will yield $m$ eigenvectors. Thus, all eigenvalues $\lambda$ will lie in an interval $[1 - c_2, 1]$ away from the origin.

## 6. CONCLUDING REMARKS

In this paper we have mainly described two approaches for iteratively solving the saddle point system (13). Both approaches avoid the use of GMRES acceleration and can be applied to two alternate choices of regularization terms. The first method is based on the CG solution of a Schur complement system, and requires double iteration, while the method, based on the augmented Lagrangian formulation, avoids double iteration. In both the cases, the preconditioners described yield rates of convergence independent of $h$, however, the rate of convergence may depend on the magnitude of the regularization parameters $\alpha_1 > 0$ and $\alpha_2 > 0$ (except for the FFT-based preconditioner applicable when $\Omega \subset R^2$).

Throughout the paper, we have assumed that $\sigma > 0$, so that matrix $A$ is symmetric positive definite. However, if $\sigma = 0$ in an application, then matrix $A$ will be *singular* with $\mathbf{1} = (1, \ldots, 1)^{\mathrm{T}}$ spanning the null space of $A$. In this case, all the preceding algorithms must be appropriately modified, by replacing $A^{-1}$ by $A^{\dagger}$. The action of $A^{\dagger}$ on a vector can be computed numerically by filtering out the components of this vector in the direction of $\mathbf{1}$ using a projection $(I - P_0)$ where $P_0$ denotes the Euclidean orthogonal projection onto $\mathbf{1}$.

### REFERENCES

1. Heinkenschloss M, Nguyen H. Balancing Neumann–Neumann methods for elliptic optimal control problems. In *Proceedings of the 15th International Conference on Domain Decomposition*, Kornhuber R, Hoppe RW, Periaux J, Pironneau O, Widlund OB, Xu J (eds), Lecture Notes in Computational Science and Engineering, vol. 40. Springer: Berlin, 2004; 589–596.
2. Heinkenschloss M, Nguyen H. Neumann–Neumann domain decomposition preconditioners for linear-quadratic elliptic optimal control problems. *SIAM Journal on Scientific Computing* 2006; **28**(3):1001–1028.
3. Lions JL. *Some Methods in the Mathematical Analysis of Systems and their Control*. Taylor & Francis: London, 1981.
4. Nguyen H. Domain decomposition methods for linear-quadratic elliptic optimal control problems. *CAAM Technical Report TR04-16*, *Ph.D. Thesis*, Rice University, 2004.
5. Pironneau O. *Optimal Shape Design for Elliptic Systems*. Springer: Berlin, 1983.
6. Benzi M, Golub GH, Liesen J. Numerical solution of saddle point problems. *Acta Numerica* 2005; 1–137.
7. Battermann A, Sachs EW. Block preconditioner for KKT systems in PDE-governed optimal control problems. In *Workshop on Fast Solutions of Discretized Optimization Problems*, Hoppe RHW, Hoffmann K-H, Schulz V (eds). Birkhäuser: Basel, 2001; 1–18.
8. Battermann A, Heinkenschloss M. Preconditioners for Karush-Kuhn-Tucker matrices arising in the optimal control of distributed systems. In *Optimal Control of Partial Differential Equations*, *Vorau 1997*, Desch W, Kappel F, Kunisch K (eds). Birkhäuser: Basel, Boston, Berlin, 1998; 15–32.
9. Biros G, Ghattas O. Parallel Lagrange–Newton–Krylov–Schur methods for PDE constrained optimization. Part I: Krylov–Schur solver. *SIAM Journal on Scientific Computing* 2005; **27**(2):687–713.
10. Biros G, Ghattas O. Parallel Lagrange–Newton–Krylov–Schur methods for PDE constrained optimization. Part II: The Lagrange–Newton solver, and its application to optimal control of steady viscous flows. *SIAM Journal on Scientific Computing* 2005; **27**(2):714–739.
11. Haber E, Ascher UM. Preconditioned all-at-once methods for large, sparse parameter estimation problems. *Inverse Problems* 2001; **17**:1847–1864.
12. Prudencio E, Byrd R, Cai X-C. Parallel full space SQP Lagrange–Newton–Krylov–Schwarz algorithms for PDE-constrained problems. *SIAM Journal on Scientific Computing* 2006; **27**:1305–1328.

13. Axelsson O. *Iterative Solution Methods*. Cambridge University Press: Cambridge, MA, 1996.
14. Elman H, Golub GH. Inexact and preconditioned Uzawa algorithms for saddle point problems. *SIAM Journal on Numerical Analysis* 1994; **31**(6):1645–1661.
15. Klawonn A. Block triangular preconditioners for saddle point problems with a penalty term. *SIAM Journal on Scientific Computing* 1998; **19**(1):172–184.
16. Murphy MF, Golub GH, Wathen AJ. A note on preconditioning for indefinite linear systems. *SIAM Journal on Scientific Computing* 2000; **21**(6):196–197.
17. Saad Y. *Iterative Methods for Sparse Linear Systems*. PWS Publishing Company: Massachusetts, 1996.
18. Rusten T, Winther R. A preconditioned iterative method for saddle point problems. *SIAM Journal on Mathematical Analysis* 1992; **13**(3):887–904.
19. Elman HC. Perturbation of eigenvalues of preconditioned Navier–Stokes operators. *SIAM Journal on Matrix Analysis and Applications* 1997; **18**(3):733–751.
20. Klawonn A. An optimal preconditioner for a class of saddle point problems with a penalty term. *SIAM Journal on Scientific Computing* 1998; **19**(2):540–552.
21. Zulehner W. Analysis of iterative methods for saddle point problems: a unified approach. *Mathematics of Computation* 2002; **71**:479–505.
22. Glowinski R, Le Tallec P. *Augmented Lagrangian and Operator Splitting Methods in Nonlinear Mechanics*. SIAM: Philadelphia, PA, 1989.
23. Bramble JH, Pasciak JE. A preconditioning technique for indefinite systems resulting from mixed approximations of elliptic problems. *Mathematics of Computation* 1988; **50**:1–18.
24. Farhat C, Roux F-X. A method of finite element tearing and interconnecting and its parallel solution algorithm. *International Journal for Numerical Methods in Engineering* 1991; **32**(6):1165–1370.
25. Lions JL, Magenes E. *Nonhomogeneous Boundary Value Problems and Applications*, vol. I. Springer: Berlin, 1972.
26. Brezzi F, Fortin M. *Mixed and Hybrid Finite Element Methods*. Springer: Berlin, 1991.
27. Raviart P-A, Girault V. *Finite Element Methods for Navier–Stokes Equations*. Springer: Berlin, 1986.
28. Axelsson O, Barker VA. *Finite Element Solution of Boundary Value Problems*: *Theory and Computation*. Academic Press: New York, 1984.
29. Braess D. *Finite Elements*: *Theory*, *Fast Solvers and Applications to Solid Mechanics*. Cambridge University Press: Cambridge, MA, 1997.
30. Brenner SC, Scott R. *Mathematical Theory of Finite Element Methods*. Springer: Berlin, 1994.
31. Tocelli A, Widlund OB. *Domain Decomposition Methods*: *Algorithms and Theory*. Spinger: Berlin, 2004.
32. Peisker P. On the numerical solution of the first biharmonic equation. *RAIRO—Mathematical Modelling and Numerical Analysis* 1998; **22**:655–676.
33. Dorhmann C, Lehoucq RB. A primal based penalty preconditioner for elliptic saddle point systems. *SIAM Journal on Numerical Analysis* 2006; **44**(1):270–282.