

ERROR ANALYSIS FOR CONSTRAINED FIRST-ORDER SYSTEM LEAST-SQUARES FINITE-ELEMENT METHODS*

J. H. ADLER[†] AND P. S. VASSILEVSKI[‡]

Abstract. In this paper, a general error analysis is provided for finite-element discretizations of partial differential equations in a saddle-point form with divergence constraint. In particular, this extends upon the work of [J. H. Adler and P. S. Vassilevski, *Springer Proc. Math. Statist.* 45, Springer, New York, 2013, pp. 1–19], giving a general error estimate for finite-element problems augmented with a divergence constraint and showing that these estimates are obtained for problems such as diffusion and Stokes’ using the first-order system least-squares (FOSLS) finite-element method. The main result is that by enforcing the constraint on a \mathbf{H}^1 -equivalent FOSLS formulation one maintains optimal convergence of the FOSLS functional (i.e., the energy norm of the error) while guaranteeing the conservation of the divergence constraint (i.e., mass conservation in some examples). The error estimates and results depend on using finite elements for the constraint space that are inf-sup stable when paired with the spaces used for the original unknowns. This includes using discontinuous spaces on coarse meshes and pairing with standard bilinear or biquadratic elements in order to confirm the results.

Key words. FOSLS, constrained minimization, mass conservation, saddle-point problem

AMS subject classifications. 65F10, 65N20, 65N30

DOI. 10.1137/130943091

1. Introduction. Least-squares finite-element methods (LSFEMs) and, in particular, the first-order system least squares (FOSLS) approach have been used widely in various applications in physics and engineering, e.g., [5, 6, 10, 12, 18, 34, 40]. This method is a finite-element discretization, which approximates the solution of a system of linear partial differential equations (PDEs) by minimizing the L^2 norm of the residual of the PDE [16, 17, 35, 36, 19]. One advantage is that this process yields symmetric positive definite (SPD) algebraic systems, which are amenable to multi-level techniques. This is true for any PDE system, including systems like Stokes, where a mixed finite-element method would yield a saddle-point problem and an indefinite linear system [13]. While a mixed method would require satisfying an inf-sup or Ladyzenskaja–Babuska–Brezzi condition [13, 15] in order to show that the finite-element spaces chosen are stable, LSFEMs do not require such restrictions. Thus, simple H^1 spaces can be used for all unknowns in the system. (\mathbf{H}^1 denotes a product of scalar H^1 spaces.) Another advantage is that these methods yield sharp and reliable a posteriori estimates [4]. This is useful for implementing adaptive local refinement techniques, which allow the approximations to be resolved more accurately in regions of higher error [14, 23].

However, a main concern with the method is that, in some applications, a loss of conservation for certain properties is noted. For instance, the Stokes or Navier–Stokes

*Submitted to the journal’s Methods and Algorithms for Scientific Computing section October 28, 2013; accepted for publication (in revised form) March 11, 2014; published electronically June 3, 2014. This work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under contract DE-AC52-07NA27344 and was sponsored by the National Science Foundation under grant DMS-1216972.

<http://www.siam.org/journals/sisc/36-3/94309.html>

[†]Department of Mathematics, Tufts University, Medford, MA 02155 (james.adler@tufts.edu).

[‡]Center for Applied Scientific Computing, Lawrence Livermore National Laboratory, Livermore, CA 94551 (panayot@llnl.gov).

system contains an equation for the conservation of momentum and one for the conservation of mass [24, 25]. Since the least-squares principle minimizes both equations equally, both quantities are conserved only up to the error tolerance given for the simulation. Therefore, attempts to improve the conservation of mass would result in a loss of accuracy in the conservation of momentum. Despite this, in several applications, conservation of a certain quantity is considered essential to capturing the true physics of the system.

In this paper, we follow up on previous work, in which methods are considered for improving the local conservation of a divergence constraint, such as mass conservation, using the FOSLS finite-element method. In [1], an approach is discussed that simply corrects the solution approximated by the FOSLS discretization so that it conserves the divergence constraint. The goal was to keep the discretization as is, preserving all of the special properties of the least-squares minimization, while still obtaining the appropriate conservation. As a result, the a posteriori error estimates and the simple finite-element spaces can still be used. This was accomplished by simply augmenting the FOSLS system with the divergence constraint and a Lagrange multiplier. Then, a subdomain correction method was used to solve the resulting saddle-point system at relatively little extra cost. This approach is not new and was done, for example, in [20] for the Stokes problem. However, choosing stable pairs of the finite-element spaces used with the Lagrange multiplier space was not really considered. In our previous paper, it was shown that only when a stable pair of elements with the constraint is used (i.e., $\mathcal{P}_2 - \mathcal{P}_0$ or $\mathcal{P}_2 - \mathcal{P}_1$) are optimal results obtained. This results from the fact that only for the stable combinations is there enough room to minimize the FOSLS functional. All finite-element pairs yield improved conservation as this is enforced directly. However, for the unstable pairings as the constraint is enforced, the solutions space is restricted and, as a result, when the FOSLS functional is minimized, there is no longer enough room to minimize certain terms in the functional. Thus, the truly constrained solution is not found. The functional is no longer estimating the \mathbf{H}^1 error accurately and the a posteriori error estimator is lost. Therefore, the conclusion is that the constraint always needs to be chosen from a space which gives a stable finite-element pair with whatever unknowns from the FOSLS system that one wishes to conserve. This requires considering an inf-sup condition for the FOSLS unknown and Lagrange multiplier pairs, but in many applications these pairs of spaces are well known [15, 24, 25].

There have been many other attempts to improve conservation properties for least-squares problems that try to avoid the issues of considering an inf-sup condition and solving saddle-point systems. For instance, simply solving the systems with more accuracy improves mass conservation. This includes using techniques such as adaptive refinement to increase the spatial resolution of the discretization [5, 10], and higher-order elements or higher-order time-stepping methods for time-dependent problems [37]. However, this requires more computational resources as it is simply driving the error down in all terms of the functional. In three-dimensional simulations this would be even more costly. Other approaches that have resulted in improved mass conservation include using divergence-free finite-element spaces [6, 2, 21, 22] or methods in which discontinuous velocity approximations are allowed [8, 9]. In addition, methods that reformulate the first-order system are used that make the approximations more conservative [27], as well as those that use a compatible least-squares method [7], which uses ideas from mixed Galerkin methods to improve the mass conservation. Finally, an alternative approach called FOSLL* [31, 32] has been developed, in which an adjoint system is considered, and the error is minimized in the L^2 norm directly.

This has been shown to improve conservation in satisfying the divergence constraint in incompressible fluid flow and electromagnetic problems.

All these methods avoid the need to solve a saddle-point system. However, there are many robust methods for solving such systems, such as those described in [1], using local subdomain corrections with an overlapping Schwarz (Vanka-like) smoother and a coarse grid correction to solve the constrained problem [41, 42, 43]. Since this postprocessing is done on local subdomains and/or on coarse grids, only a fractional amount of computational cost is added to the solution process. Thus, the Lagrange multiplier approach can be a robust method for ensuring mass conservation as well as maintaining the important properties of a least-squares minimization process.

In this paper, we extend upon the work of [1] and perform a model error analysis of the method. In particular, we give a general error estimate for finite-element problems augmented with a divergence constraint and show that we obtain these estimates that apply to problems such as diffusion and Stokes using the FOSLS method. The main result is that by enforcing the constraint on an \mathbf{H}^1 -bounded FOSLS formulation one maintains optimal convergence of the FOSLS functional (i.e., the energy norm of the error). This convergence may depend on the constraint space that is used; however, for standard problems as shown below, this does not pollute the convergence compared to the pure FOSLS method and has better local conservation properties.

The paper is outlined as follows. In section 2, we introduce some notation and briefly describe the FOSLS discretization along with the constraint method. Section 3 describes the error analysis of the Lagrange multiplier approach, along with a bound on the energy norm of the error. Numerical results are then given in section 4. Finally, concluding remarks and a discussion of future work is given in section 5.

2. Constrained first-order system least-squares. Consider a PDE system that is first put into a differential first-order system of equations, denoted by $L\mathbf{U} = \mathbf{f}$. Here, L is a mapping from an appropriate Hilbert space, \mathcal{V} , to an L^2 product space. In many contexts, \mathcal{V} is chosen to be an \mathbf{H}^1 product space with appropriate boundary conditions.

This minimization is written as

$$(2.1) \quad \mathbf{U}_* = \arg \min_{\mathbf{U} \in \mathcal{V}} \mathcal{G}(\mathbf{U}; \mathbf{f}) := \arg \min_{\mathbf{U} \in \mathcal{V}} \|\mathbf{L}\mathbf{U} - \mathbf{f}\|_0^2,$$

where \mathbf{U}_* is the solution in an appropriate \mathbf{H}^1 space. The minimization results in the weak form of the problem: Find $\mathbf{U}_* \in \mathcal{V}$ such that

$$(2.2) \quad a(\mathbf{U}, \mathbf{v}) := (L\mathbf{U}_*, L\mathbf{v}) = (\mathbf{f}, L\mathbf{v}) \quad \text{for all } \mathbf{v} \in \mathcal{V},$$

where (\cdot, \cdot) is the usual L^2 inner product on the product space, $(L^2)^k$, for k equations in the linear system. If the following properties of the bilinear form, $a(\mathbf{U}, \mathbf{v})$ are assumed: \exists constants, c_1 and c_2 , such that

$$(2.3) \quad \text{continuity} \quad a(\mathbf{U}, \mathbf{v}) \leq c_2 \|\mathbf{U}\|_{\mathcal{V}} \|\mathbf{v}\|_{\mathcal{V}} \quad \text{for all } \mathbf{U}, \mathbf{v} \in \mathcal{V},$$

$$(2.4) \quad \text{coercivity} \quad a(\mathbf{U}, \mathbf{U}) \geq c_1 \|\mathbf{U}\|_{\mathcal{V}}^2 \quad \text{for all } \mathbf{u} \in \mathcal{V},$$

then, by the Riesz representation theorem, this bilinear form is an inner product on \mathcal{V} [30]. In addition, these properties imply the existence of a unique solution, $\mathbf{U}_* \in \mathcal{V}$, for the weak problem (2.2). Here, c_1 and c_2 depend only on the operator, L , and the domain of the problem. They are independent of \mathbf{U} and \mathbf{v} .

Next, \mathbf{U}_* is approximated by restricting (2.1) to a finite-dimensional space, $\mathbf{V}_h \subseteq \mathbf{V}$, which leads to (2.2) restricted to \mathbf{V}_h . Since \mathbf{V}_h is a subspace of \mathbf{V} , the discrete problem is also well-posed. Choosing an appropriate basis, $\mathbf{V}_h = \text{span}\{\Phi_j\}$, and restricting (2.2) to this basis, yields an algebraic system of equations involving the matrix, A , with elements

$$(2.5) \quad (A)_{ij} = (L\Phi_j, L\Phi_i).$$

It has been shown that in the context of an SPD \mathbf{H}^1 -equivalent bilinear form restricted to a finite-element subspace, a multilevel technique exists that yields optimal convergence to the linear system [17].

To illustrate this further and to demonstrate some results of the constrained method, two test problems that have a divergence constraint are considered: diffusion and Stokes. In [1], a diffusion problem in two dimensions was considered, $-\nabla \cdot \nabla p = g$, rewritten as a first-order system:

$$(2.6) \quad -\nabla \cdot \mathbf{u} = g,$$

$$(2.7) \quad \nabla \times \mathbf{u} = 0,$$

$$(2.8) \quad \mathbf{u} - \nabla p = 0.$$

Here, the gradient of the solution is introduced as an extra variable and the extra curl equation is added so that the weak system is continuous and coercive and, therefore, \mathbf{H}^1 -equivalent [16, 17]. Then, the following functional is minimized:

$$\mathcal{G}(\mathbf{u}, p; \mathbf{f}) = \|\nabla \cdot \mathbf{u} + g\|_0^2 + \|\nabla \times \mathbf{u}\|_0^2 + \|\mathbf{u} - \nabla p\|_0^2.$$

The resulting discrete system is

$$A\mathbf{U} = b,$$

where $\mathbf{U} = (\mathbf{u}, p)^T$. Here, A is the matrix as defined in (2.5), where L now refers to system (2.6)–(2.8). Similarly, the right-hand-side vector, b , is defined as $b_i = (\mathbf{f}, L\Phi_i)$, where $\mathbf{f} = (g, 0, 0)^T$. When minimizing this functional, equal weight is given to each term in the system. Therefore, if better accuracy is needed on a certain term, such as the divergence constraint, accuracy is lost in the other portions. In many applications, however, exact conservation of certain terms is important for developing an accurate model of a physical system. For instance, one may want to conserve the “mass” of the system. This is defined as

$$(2.9) \quad \int_{\Omega} -\nabla \cdot \mathbf{u} \, d\Omega = \int_{\Omega} g \, d\Omega.$$

In addition, in many applications *local* mass conservation is desired instead, where the mass is conserved in all regions of the domain, including a single element.

This notion of mass has more physical meaning in the case of modeling an incompressible fluid via Stokes equations:

$$\begin{aligned} -\nabla \cdot \nabla \mathbf{u} + \nabla p &= \mathbf{g}, \\ \nabla \cdot \mathbf{u} &= 0. \end{aligned}$$

Now, (2.9) states that the amount of flow in or out of the system is equal to the flow contributed by the source (in this case 0). Thus, a velocity-vorticity-pressure first-order formulation of Stokes equations is the second system considered. The vorticity,

$\omega = \nabla \times \mathbf{u}$, is introduced and the system is augmented to make it formally \mathbf{H}^1 -equivalent:

$$(2.10) \quad \nabla \times \omega + \nabla p = \mathbf{g},$$

$$(2.11) \quad \nabla \cdot \mathbf{u} = 0,$$

$$(2.12) \quad \omega - \nabla \times \mathbf{u} = \mathbf{0},$$

$$(2.13) \quad \nabla \cdot \omega = 0,$$

It should be noted that in two dimensions the vorticity is a scalar, and (2.13) is not present. Here, preserving the incompressibility condition, (2.11), is difficult at the discrete level but can be considered as an added constraint to the system. The FOSLS function in this instance is

$$\mathcal{G}(\mathbf{u}, \omega, p; \mathbf{f}) = \|-\nabla \times \omega + \nabla p - \mathbf{g}\|_0^2 + \|\nabla \cdot \mathbf{u}\|_0^2 + \|\omega - \nabla \times \mathbf{u}\|_0^2 + \|\nabla \cdot \omega\|_0^2,$$

and we obtain a discrete linear system, $A\mathbf{U} = b$, where $\mathbf{U} = (\mathbf{u}, p, \omega)^t$, and $b_i = (\mathbf{f}, L\Phi_i)$, where $\mathbf{f} = (\mathbf{g}, 0, \mathbf{0}, 0)^t$.

In both cases, to enforce mass conservation, a Lagrange multiplier, λ , is introduced and the FOSLS system is augmented as follows:

$$(2.14) \quad \begin{pmatrix} A & B^T \\ B & 0 \end{pmatrix} \begin{pmatrix} \mathbf{U} \\ \lambda \end{pmatrix} = \begin{pmatrix} b \\ \hat{b} \end{pmatrix}.$$

Here, A and \mathbf{U} are as before for the FOSLS discretization, λ is the Lagrange multiplier, and B is a finite-element assembly of the constraint, in these two examples, $\nabla \cdot \mathbf{u} = f$.

For the rest of the paper, we consider a triangulation of a mesh in two dimensions, \mathcal{T}_h with grid spacing h . In addition, consider the polynomial spaces of order k defined on this triangulation as \mathcal{P}_k . Let $\Phi_j \in [\mathcal{P}_{k_1}]^2$ be a vector and let $w_i \in \mathcal{P}_{k_2}$ be a scalar. Thus, B is defined as

$$B_{ij} = (\nabla \cdot \Phi_j, w_i)$$

and

$$\hat{b}_i = (f, w_i).$$

In the case of the diffusion problem, (2.6)–(2.8), $f := -g$, and in the case of Stokes, (2.10)–(2.13), $f = 0$. Choosing a discontinuous space for w ensures the local conservation of mass (on the element level). For this paper, we consider \mathcal{P}_0 elements (i.e., a finite-element space of discontinuous piecewise constants) only.

3. Error estimates. The following section analyzes this system from a variational point of view and develops error estimates for the solution to such a system.

3.1. A general error estimate for problems with divergence constraint.

Consider a constrained minimization problem that leads to the following saddle-point problem: Find $\mathbf{U} = \begin{bmatrix} * \\ \mathbf{u} \end{bmatrix} \in \mathbf{H}^1$, $\mathbf{u} \in \mathbf{H}_0(\text{div})$, and $\lambda \in L_0^2$ such that

$$(3.1) \quad a(\mathbf{U}, \mathbf{v}) + (\lambda, \nabla \cdot \underline{\theta}) = (\mathbf{f}, \mathbf{v}) \quad \text{for all } \mathbf{v} = \begin{bmatrix} * \\ \underline{\theta} \end{bmatrix} \in \mathbf{H}^1, \underline{\theta} \in \mathbf{H}_0(\text{div}),$$

$$(\nabla \cdot \mathbf{u}, w) = (f, w) \quad \text{for all } w \in L_0^2.$$

Here, we define $\mathbf{H}_0(\text{div}) = \{\mathbf{u} \in \mathbf{L}_0^2 : \nabla \cdot \mathbf{u} \in L_0^2 \text{ and } \mathbf{n} \cdot \mathbf{u} = 0 \text{ on } \partial\Omega\}$, where \mathbf{n} is the unit normal to the domain and \mathbf{L}_0^2 denotes a product of scalar L_0^2 spaces. Unless noted otherwise, these spaces are defined over the domain Ω ; $\mathbf{H}^1 = \mathbf{H}^1(\Omega)$ and $L_0^2 = L_0^2(\Omega)$. A general bilinear form $a(\cdot, \cdot)$ is considered, which is symmetric and positive definite, giving rise to the energy norm

$$(3.2) \quad \|\mathbf{v}\|_a = \sqrt{a(\mathbf{v}, \mathbf{v})}.$$

In the FOSLS setting, $a(\cdot, \cdot)$ is defined as in (2.2) and $\|\mathbf{v}\|_a = \mathcal{F}(\mathbf{v}; 0)$, where $\mathcal{F}(\mathbf{v}; \mathbf{f}) = \sqrt{\mathcal{G}(\mathbf{v}; \mathbf{f})}$. It is assumed that $a(\cdot, \cdot)$ is \mathbf{H}^1 -bounded or continuous, as in (2.3), i.e.,

$$(3.3) \quad a(\mathbf{U}, \mathbf{v}) \leq C \|\mathbf{U}\|_a \|\mathbf{v}\|_1 \text{ for all } \mathbf{U}, \mathbf{v} \in \mathbf{H}^1,$$

and that it is weakly coercive, in the sense that

$$(3.4) \quad C \|\nabla \cdot \mathbf{u}\|^2 \leq a(\mathbf{U}, \mathbf{U}) \text{ for any } \mathbf{U} = \begin{bmatrix} * \\ \mathbf{u} \end{bmatrix} \in \mathbf{H}^1 \text{ and } \mathbf{u} \in \mathbf{H}_0(\text{div}).$$

Note, that in the case of a continuous and coercive FOSLS bilinear form, where the divergence constraint is part of the functional itself, (3.4) is satisfied by the coercivity condition, (2.4).

Remark 3.1. The weak-coercivity assumption, (3.4), can always be ensured by testing the second equation in (3.1) with $w = \nabla \cdot \underline{\theta} \in L_0^2$ for any $\underline{\theta} \in \mathbf{H}_0(\text{div})$ and adding it to the first one. This leads to the equivalent saddle-point system

$$(3.5) \quad \begin{aligned} (\nabla \cdot \mathbf{u}, \nabla \cdot \underline{\theta}) + a(\mathbf{U}, \mathbf{v}) + (\lambda, \nabla \cdot \underline{\theta}) &= (\mathbf{f}, \mathbf{v}) + (f, \nabla \cdot \underline{\theta}) \\ \text{for all } \mathbf{v} = \begin{bmatrix} * \\ \underline{\theta} \end{bmatrix} \in \mathbf{H}^1, \underline{\theta} \in \mathbf{H}_0(\text{div}), \\ (\nabla \cdot \mathbf{u}, w) &= (f, w) \text{ for all } w \in L_0^2. \end{aligned}$$

In this way the bilinear form becomes $\tilde{a}(\mathbf{U}, \mathbf{v}) := (\nabla \cdot \mathbf{u}, \nabla \cdot \underline{\theta}) + a(\mathbf{U}, \mathbf{v})$, which is weakly coercive in the sense of (3.4) as long as the original one is SPD (which is assumed).

The discrete problem reads as follows: Find $\mathbf{U}_h = [\mathbf{u}_h^*] \in \mathcal{V}_h \subset \mathbf{H}^1$, $\mathbf{u}_h \cdot \mathbf{n} = 0$ on $\partial\Omega$, and $\lambda_H \in \mathcal{W}_H \subset L_0^2$, such that

$$\begin{aligned} a(\mathbf{U}_h, \mathbf{v}) + (\lambda_H, \nabla \cdot \underline{\theta}) &= (\mathbf{g}, \mathbf{v}) \text{ for all } \mathbf{v} = \begin{bmatrix} * \\ \underline{\theta} \end{bmatrix} \in \mathcal{V}_h, \\ (\nabla \cdot \mathbf{u}_h, w) &= (f, w) \text{ for all } w \in \mathcal{W}_H. \end{aligned}$$

Thus, the following assumptions are made. There is an H^1 -conforming scalar finite element space, \mathcal{V}_h , which is associated with a triangulation \mathcal{T}_h of the computational domain, Ω (polygon in two dimensions or polytope in three dimensions). To discretize the bilinear form $a(\cdot, \cdot)$, which is defined on $\mathbf{H}^1 \times \mathbf{H}^1$, where $\mathbf{H}^1 = (H^1)^d$ for a given $d \geq 1$, $\mathcal{V}_h = (\mathcal{V}_h)^d$ is used. To be specific, homogeneous normal boundary conditions are imposed for the component of $\mathbf{U} \in \mathbf{H}^1$ used in the divergence constraint denoted by \mathbf{u} , i.e., $\mathbf{u} \cdot \mathbf{n} = 0$ on $\partial\Omega$. This implies that the Lagrange multiplier space \mathcal{W}_H consists of functions with zero mean value, i.e., $\mathcal{W}_H \subset L_0^2(\Omega)$.

The Lagrange multiplier finite-element space, \mathcal{W}_H , is associated with another (quasi-uniform) mesh, \mathcal{T}_H , which is assumed to be sufficiently coarser with respect to \mathcal{T}_h . This condition in particular implies that the discrete problem is not over-constrained. It is assumed that \mathcal{W}_H consists of discontinuous piecewise polynomials of a certain degree. Thus, \mathcal{W}_H can be paired with a Raviart–Thomas space of a certain polynomial order, $\mathbf{R}_H \subset H_0(\text{div})$ (not needed in the discretization), such that the following inf-sup condition holds (cf. [38]):

$$(3.6) \quad \beta \|w\|_{1,H} \leq \sup_{\psi \in \mathbf{R}_H} \frac{(w, \nabla \cdot \psi)}{\|\psi\|_0} \simeq \sup_{\psi \in \mathbf{R}_H} \frac{(w, \nabla \cdot \psi)}{(\|\psi\|_0^2 + H^2 \|\nabla \cdot \psi\|_0^2)^{\frac{1}{2}}},$$

where $\|\cdot\|_{1,H}$ comes from the so-called interior penalty quadratic form, defined as

$$(3.7) \quad \|w\|_{1,H}^2 \equiv \sum_{T \in \mathcal{T}_H} \int_T |\nabla w|^2 \, dx + \sum_{E \in \mathcal{E}_H} \frac{1}{H} \int_E [w]^2 \, d\rho.$$

Here, $[\cdot] = [\cdot]_E$ stands for the jump across the interface E between two adjacent elements from \mathcal{T}_H . These (interior) interfaces form the set \mathcal{E}_H .

3.1.1. Error analysis. Introduce the errors $\mathbf{E} = \mathbf{U} - \mathbf{U}_h$ with components $\mathbf{e} = \mathbf{u} - \mathbf{u}_h$ and $\delta = \lambda - \lambda_H$. They satisfy the Galerkin relations

$$(3.8) \quad a(\mathbf{E}, \mathbf{v}) + (\delta, \nabla \cdot \underline{\theta}) = 0 \quad \text{for all } \mathbf{v} = \begin{bmatrix} * \\ \underline{\theta} \end{bmatrix} \in \mathcal{V}_h,$$

$$(\nabla \cdot \mathbf{e}, w) = 0 \quad \text{for all } w \in \mathcal{W}_H.$$

Next, some auxiliary estimates are proved.

LEMMA 3.1. *For any $s \in (\frac{1}{2}, 1]$, and any $w \in \mathcal{V}_h$ and $\psi \in \mathbf{H}_0^1$, there exists a $C > 0$ such that*

$$(3.9) \quad (w, \nabla \cdot \psi) \leq C \|w\|_{1,H} (\|\psi\|_0^2 + H^{2s} |\psi|_s^2)^{\frac{1}{2}}.$$

Proof. The result can be established as follows. Let \mathcal{E}_H be the set of (interior) edges (faces in three dimensions) of elements from the triangulation \mathcal{T}_H . For each $E \in \mathcal{E}_H$ denote by T_E^+ and T_E^- the two neighboring elements that share E . Then based on integration by parts, a trace inequality for H^s -functions, for any $s \in (\frac{1}{2}, 1]$, on the domains T_E^+ and T_E^- with diameter $\mathcal{O}(H)$, and Cauchy–Schwarz inequalities, we have

$$\begin{aligned} & (w, \nabla \cdot \psi) \\ &= \sum_{T \in \mathcal{T}_H} \int_T \nabla w \cdot \psi + \sum_{E \in \mathcal{E}_H} \int_E [w] \psi \cdot \mathbf{n}_E \\ &\leq C \left[\sum_{T \in \mathcal{T}_H} \|\nabla w\|_T \|\psi\|_T + \sum_{E \in \mathcal{E}_H} \left(H^{-1} \int_E [w]^2 \right)^{\frac{1}{2}} \left(H \int_E (\psi \cdot \mathbf{n}_E)^2 \right)^{\frac{1}{2}} \right] \\ &\leq C \left[\sum_{T \in \mathcal{T}_H} \|\nabla w\|_T \|\psi\|_T + \sum_{E \in \mathcal{E}_H} \left(H^{-1} \int_E [w]^2 \right)^{\frac{1}{2}} \left(\|\psi\|_{T_E^+ \cup T_E^-}^2 + H^{2s} |\psi|_{s, T_E^+ \cup T_E^-}^2 \right)^{\frac{1}{2}} \right] \\ &\leq C \left[\sum_{T \in \mathcal{T}_H} \|\nabla w\|_T^2 + \sum_{E \in \mathcal{E}_H} \frac{1}{H} \int_E [w]^2 \right]^{\frac{1}{2}} \left[\|\psi\|^2 + H^{2s} |\psi|_s^2 \right]^{\frac{1}{2}} \\ &= C \|w\|_{1,H} (\|\psi\|^2 + H^{2s} |\psi|_s^2)^{\frac{1}{2}}. \end{aligned}$$

The trace inequality that we use states that

$$H \int_E (\boldsymbol{\psi} \cdot \mathbf{n}_E)^2 \leq C \left(\|\boldsymbol{\psi}\|_{T_E^+ \cup T_E^-}^2 + H^{2s} |\boldsymbol{\psi}|_{s, T_E^+ \cup T_E^-}^2 \right),$$

where C is mesh-independent but depends on the transformation from the reference element to the H -dependent mesh. Since all meshes we consider are generated from uniform refinement, this constant is bounded. \square

LEMMA 3.2. *The following inf-sup estimate holds if h/H is sufficiently small:*

$$(3.10) \quad C \|w\|_{1,H} \leq \sup_{\underline{\boldsymbol{\theta}} \in \mathcal{V}_h} \frac{(w, \nabla \cdot \underline{\boldsymbol{\theta}})}{(\|\underline{\boldsymbol{\theta}}\|_0^2 + H^2 \|\nabla \underline{\boldsymbol{\theta}}\|_0^2)^{\frac{1}{2}}} \text{ for all } w \in \mathcal{W}_H.$$

Proof. Given $\mathbf{z} \in \mathbf{H}_0^1$, let $\mathbf{z}_h \in \mathcal{V}_h$ (vanishing on $\partial\Omega$) be a stable finite interpolant such that

$$(3.11) \quad \|\mathbf{z}_h\|_0 \leq C \|\mathbf{z}\|_0, \quad \|\nabla \mathbf{z}_h\|_0 \leq C \|\nabla \mathbf{z}\|_0, \quad \text{and} \quad \|\mathbf{z} - \mathbf{z}_h\|_0 \leq Ch \|\nabla \mathbf{z}\|_0.$$

Now use estimate (3.9) from Lemma 3.1 for $\boldsymbol{\psi} = \mathbf{z} - \mathbf{z}_h$, which gives, for any parameter $\tau > 0$,

$$\begin{aligned} (w, \nabla \cdot (\mathbf{z} - \mathbf{z}_h)) &\leq C \|w\|_{1,H} H \left(\frac{h}{H} + \left(\frac{h}{H} \right)^{1-s} \right) \|\nabla \mathbf{z}\| \\ &\leq C \frac{H}{\sqrt{\tau}} \left(\frac{h}{H} + \left(\frac{h}{H} \right)^{1-s} \right) \|w\|_{1,H} (\|\mathbf{z}\|_0^2 + \tau \|\nabla \mathbf{z}\|_0^2)^{\frac{1}{2}}. \end{aligned}$$

Based on a regular decomposition result, found, for example, in [29], for each $\mathbf{u} \in \mathbf{H}_0(\Omega, \text{div})$ and $\tau > 0$, there is a $\mathbf{z} \in \mathbf{H}_0^1(\Omega)$ such that $\nabla \cdot \mathbf{z} = \nabla \cdot \mathbf{u}$ and

$$\|\mathbf{z}\|^2 + \tau \|\nabla \mathbf{z}\|^2 \leq c_0 (\|\mathbf{u}\|_0^2 + \tau \|\nabla \cdot \mathbf{u}\|_0^2)$$

for a uniform constant $c_0 > 0$ (independent of τ). This implies the estimate

$$(3.12) \quad \sup_{\mathbf{u} \in \mathbf{H}_0(\Omega, \text{div})} \frac{(\bar{\lambda}, \nabla \cdot \mathbf{u})}{(\|\mathbf{u}\|_0^2 + \tau \|\nabla \cdot \mathbf{u}\|_0^2)^{\frac{1}{2}}} \leq \sqrt{c_0} \sup_{\mathbf{z} \in \mathbf{H}_0^1(\Omega)} \frac{(\bar{\lambda}, \nabla \cdot \mathbf{z})}{(\|\mathbf{z}\|_0^2 + \tau \|\nabla \mathbf{z}\|_0^2)^{\frac{1}{2}}}.$$

Now, using the inf-sup estimate, (3.6), characterizing the $\|\cdot\|_{1,H}$ -norm (defined in (3.7)) and the above estimate, (3.12), we then obtain (since $\mathbf{R}_H \subset \mathbf{H}_0(\text{div})$)

$$\beta \|w\|_{1,H} \leq \sup_{\mathbf{u} \in \mathbf{H}_0(\text{div})} \frac{(w, \nabla \cdot \mathbf{u})}{(\|\mathbf{u}\|^2 + H^2 \|\nabla \cdot \mathbf{u}\|^2)^{\frac{1}{2}}} \leq \sqrt{c_0} \sup_{\mathbf{z} \in \mathbf{H}_0^1(\Omega)} \frac{(w, \nabla \cdot \mathbf{z})}{(\|\mathbf{z}\|^2 + H^2 \|\nabla \mathbf{z}\|^2)^{\frac{1}{2}}}.$$

Therefore,

$$\beta \|w\|_{1, H} \leq \sqrt{c_0} \sup_{\mathbf{z} \in \mathbf{H}_0^1(\Omega)} \frac{(w, \nabla \cdot (\mathbf{z} - \mathbf{z}_h))}{(\|\mathbf{z}\|^2 + H^2 \|\nabla \mathbf{z}\|^2)^{\frac{1}{2}}} + C \sup_{\mathbf{z}_h \in \mathcal{V}_h} \frac{(w, \nabla \cdot \mathbf{z}_h)}{(\|\mathbf{z}_h\|^2 + H^2 \|\nabla \mathbf{z}_h\|^2)^{\frac{1}{2}}}.$$

Using properties (3.11), the last three estimates imply, assuming that $\tau \leq CH^2$,

$$\begin{aligned} \beta \|w\|_{1, H} &\leq C \frac{H}{\sqrt{\tau}} \left(\frac{h}{H} + \left(\frac{h}{H} \right)^{1-s} \right) \|w\|_{1, H} \sup_{\mathbf{z} \in \mathbf{H}_0^1(\Omega)} \frac{(\|\mathbf{z}\|_0^2 + \tau \|\nabla \mathbf{z}\|^2)^{\frac{1}{2}}}{(\|\mathbf{z}\|_0^2 + H^2 \|\nabla \mathbf{z}\|^2)^{\frac{1}{2}}} \\ &\quad + C \sup_{\mathbf{z}_h \in \mathcal{V}_h} \frac{(w, \nabla \cdot \mathbf{z}_h)}{(\|\mathbf{z}_h\|^2 + H^2 \|\nabla \mathbf{z}_h\|^2)^{\frac{1}{2}}} \\ &\leq C \frac{H}{\sqrt{\tau}} \left(\frac{h}{H} + \left(\frac{h}{H} \right)^{1-s} \right) \|w\|_{1, H} + C \sup_{\mathbf{z}_h \in \mathcal{V}_h} \frac{(w, \nabla \cdot \mathbf{z}_h)}{(\|\mathbf{z}_h\|^2 + H^2 \|\nabla \mathbf{z}_h\|^2)^{\frac{1}{2}}}. \end{aligned}$$

Note that we started with a vector-function, \mathbf{z}_h , having all its components vanish on $\partial\Omega$. In the last estimate, the space is enlarged to \mathcal{V}_h , which requires that only $\mathbf{z}_h \cdot \mathbf{n} = 0$ on $\partial\Omega$.

Next, choose $\tau \simeq H^2$ and fix $s < 1$ so that for h/H sufficiently small the term

$$C \frac{H}{\sqrt{\tau}} \left(\frac{h}{H} + \left(\frac{h}{H} \right)^{1-s} \right) \|w\|_{1, H}$$

can be absorbed by the similar term on the left-hand side of the inequality. This completes the proof. \square

With these estimates, we now prove the main result.

THEOREM 3.3. *Under the assumption that the bilinear form $a(\cdot, \cdot)$ is \mathbf{H}^1 -bounded (as in (3.3)) and weakly coercive as in (3.4) (see also Remark 3.1), the following optimal energy error estimate holds for h/H sufficiently small:*

$$\|\mathbf{U} - \mathbf{U}_h\|_a \leq C \left\{ \inf_{\mathbf{v} \in \mathcal{V}_h} [\|\mathbf{U} - \mathbf{v}\|_a + H^{-1} \|\mathbf{u} - \underline{\theta}\| + \|\nabla(\mathbf{u} - \underline{\theta})\|] + \inf_{w \in \mathcal{W}_H} \|\lambda - w\| \right\}.$$

Above, $\mathbf{U} = [\mathbf{u}^*]$ is the exact solution and $\lambda \in L_0^2$ is the exact Lagrange multiplier, both of which are assumed to exist and to belong to the respective functional spaces. Also, \mathbf{U}_h is the finite-element solution, $\mathbf{v} = [\underline{\theta}^*] \in \mathcal{V}_h$, and $w \in \mathcal{W}_H$ are any finite-element functions.

Proof. For any $\mathbf{v} \in \mathcal{V}_h$ and $w \in \mathcal{W}_H$, the following identities hold:

$$\begin{aligned} (3.13) \quad a(\mathbf{E}, \mathbf{E}) &= a(\mathbf{E}, \mathbf{v}) + a(\mathbf{E}, \mathbf{E} - \mathbf{v}) + (\delta, \nabla \cdot \underline{\theta}) - (\delta, \nabla \cdot \underline{\theta}) \\ &= a(\mathbf{E}, \mathbf{E} - \mathbf{v}) - (\delta, \nabla \cdot \underline{\theta}) \\ &= a(\mathbf{E}, \mathbf{E} - \mathbf{v}) - (\lambda - w, \nabla \cdot \underline{\theta}) + (\lambda_H - w, \nabla \cdot \underline{\theta}) \\ &= a(\mathbf{E}, \mathbf{E} - \mathbf{v}) + (\lambda - w, \nabla \cdot (\mathbf{e} - \underline{\theta})) - (\lambda - w, \nabla \cdot \mathbf{e}) \\ &\quad + (\lambda_H - w, \nabla \cdot (\underline{\theta} - \mathbf{e})). \end{aligned}$$

From the last inf-sup condition, Lemma 3.2, and the identity (3.8), for the errors $\delta = \lambda - \lambda_H$ and $\mathbf{E} = \mathbf{U} - \mathbf{U}_h$,

$$(-\delta, \nabla \cdot \underline{\theta}) = a(\mathbf{E}, \mathbf{v}) \text{ for any } \mathbf{v} = \begin{bmatrix} * \\ \underline{\theta} \end{bmatrix} \in \mathcal{V}_h.$$

Thus, since $\lambda_H = \lambda - \delta$, $(\lambda_H - w, \nabla \cdot \underline{\theta}) = (\lambda - w, \nabla \cdot \underline{\theta}) + a(\mathbf{E}, \mathbf{v})$,

$$\begin{aligned} C\|\lambda_H - w\|_{1,H} &\leq \sup_{\mathbf{v}=[0,\underline{\theta}]^t \in \mathcal{V}_h} \frac{(\lambda_H - w, \nabla \cdot \underline{\theta})}{\|\underline{\theta}\| + H\|\nabla \underline{\theta}\|} = \sup_{\mathbf{v}=[0,\underline{\theta}]^t \in \mathcal{V}_h} \frac{a(\mathbf{E}, \mathbf{v}) + (\lambda - w, \nabla \cdot \underline{\theta})}{\|\mathbf{v}\| + H\|\nabla \mathbf{v}\|} \\ &\leq \sup_{\mathbf{v} \in \mathcal{V}_h} \frac{a(\mathbf{E}, \mathbf{v}) + (\lambda - w, \nabla \cdot \underline{\theta})}{\|\mathbf{v}\| + H\|\nabla \mathbf{v}\|} \leq CH^{-1} (\|\mathbf{E}\|_a + \|\lambda - w\|). \end{aligned}$$

As a result, from the representation (3.13) and inequality (3.9), used for $w := \lambda_H - w$, $\psi = \underline{\theta} - \mathbf{e}$, and $s = 1$,

$$\begin{aligned} \|\mathbf{E}\|_a^2 &\leq \|\mathbf{E}\|_a \|\mathbf{E} - \mathbf{v}\|_a + \|\lambda - w\| (\|\nabla \cdot (\mathbf{e} - \underline{\theta})\| + \|\nabla \cdot \mathbf{e}\|) + (\lambda_H - w, \nabla \cdot (\underline{\theta} - \mathbf{e})) \\ &\leq \|\mathbf{E}\|_a \|\mathbf{E} - \mathbf{v}\|_a + \|\lambda - w\| (\|\nabla \cdot (\mathbf{e} - \underline{\theta})\| + \|\nabla \cdot \mathbf{e}\|) \\ &\quad + C\|\lambda_H - w\|_{1,H} (\|\mathbf{e} - \underline{\theta}\|^2 + H^2 \|\nabla(\mathbf{e} - \underline{\theta})\|^2)^{\frac{1}{2}} \\ &\leq \|\mathbf{E}\|_a \|\mathbf{E} - \mathbf{v}\|_a + \|\lambda - w\| (\|\nabla \cdot (\mathbf{e} - \underline{\theta})\| + \|\nabla \cdot \mathbf{e}\|) \\ &\quad + CH^{-1} (\|\mathbf{E}\|_a + \|\lambda - w\|) (\|\mathbf{e} - \underline{\theta}\|^2 + H^2 \|\nabla(\mathbf{e} - \underline{\theta})\|^2)^{\frac{1}{2}}. \end{aligned}$$

Letting $\mathbf{v} := \mathbf{v} - \mathbf{U}_h$, hence $\underline{\theta} := \underline{\theta} - \mathbf{u}_h$, gives $\mathbf{E} - \mathbf{v} := \mathbf{U} - \mathbf{v}$ and $\mathbf{e} - \underline{\theta} := \mathbf{u} - \underline{\theta}$. Therefore,

$$\begin{aligned} \|\mathbf{E}\|_a^2 &\leq \|\mathbf{E}\|_a \|\mathbf{U} - \mathbf{v}\|_a + \|\lambda - w\| (\|\nabla \cdot (\mathbf{u} - \underline{\theta})\| + \|\nabla \cdot \mathbf{e}\|) \\ &\quad + C (\|\mathbf{E}\|_a + \|\lambda - w\|) (H^{-2} \|\mathbf{u} - \underline{\theta}\|^2 + \|\nabla(\mathbf{u} - \underline{\theta})\|^2)^{\frac{1}{2}}. \end{aligned}$$

Using a Cauchy–Schwarz inequality, for a small $\epsilon > 0$, then

$$\begin{aligned} \|\mathbf{E}\|_a \|\mathbf{U} - \mathbf{v}\|_a &\leq \frac{1}{2} \|\mathbf{E}\|_a^2 + \frac{1}{2} \|\mathbf{U} - \mathbf{v}\|_a^2, \\ C\|\mathbf{E}\|_a (H^{-2} \|\mathbf{u} - \underline{\theta}\|^2 + \|\nabla(\mathbf{u} - \underline{\theta})\|^2)^{\frac{1}{2}} &\leq \epsilon \|\mathbf{E}\|_a^2 + \frac{C^2}{4\epsilon} (H^{-2} \|\mathbf{u} - \underline{\theta}\|^2 + \|\nabla(\mathbf{u} - \underline{\theta})\|^2), \\ \|\lambda - w\| \|\nabla \cdot \mathbf{e}\| &\leq \epsilon \|\nabla \cdot \mathbf{e}\|^2 + \frac{1}{4\epsilon} \|\lambda - w\|^2. \end{aligned}$$

Finally, the desired optimal energy error estimate for the vector unknown, \mathbf{U} , is straightforward by noticing that the terms $(\frac{1}{2} + \epsilon) \|\mathbf{E}\|_a^2$ and $\epsilon \|\nabla \cdot \mathbf{e}\|^2$ can be absorbed into the term $\|\mathbf{E}\|_a^2 = a(\mathbf{E}, \mathbf{E})$ on the left-hand side, based on the weak-coercivity assumption, (3.4). \square

Remark 3.2. In the above estimates, it is not assumed that the FOSLS bilinear form, $a(\cdot, \cdot)$, is \mathbf{H}^1 -coercive. If such stronger coercivity is assumed (which does hold in certain cases), then the constraint problem is essentially like Stokes', and any error analysis available for discretized Stokes problems can be adopted.

Finally, note that an error estimate for the Lagrange multiplier λ is not given, since it is not needed by the method. Next, we describe some numerical experiments that are used to illustrate this result.

4. Numerical results. In the following section, the mass conservation and errors are measured for two types of problems. In both cases, the FOSLS formulation, $A\mathbf{U} = b$, is solved using the preconditioned conjugate gradient (PCG) method with a single $V(1, 1)$ algebraic multigrid (AMG) cycle used as a preconditioner. BoomerAMG from the HYPRE package [28] developed at the Lawrence Livermore National Laboratory is used, with symmetric hybrid Gauss–Seidel (Gauss–Seidel on nodes within the processor and block Jacobi across processors) as the smoother. The number of PCG iterations, IT , is reported in Tables 1 and 2 in order to reduce the linear residual to machine precision. To solve the constrained saddle-point problem, (2.14), many techniques can be used, including the subdomain correction ideas presented in [1]. However, as the main point of this paper is to show the error estimates, we use a simple preconditioned MINRES algorithm [39] and solve the system to machine precision. Such low tolerances, for both the constrained and the pure FOSLS problem, are used to ensure that the asymptotic regimes have been reached. In practice, many fewer iterations would be needed for both systems. It should be noted that in the following set of results, the number of PCG iterations required to solve the pure FOSLS system is almost half that of the number of MINRES iterations needed to solve the constrained problem. However, again, this many iterations would not be needed in practice and when coupled with the subdomain correction ideas presented in [1], the constrained approach would only require a minimal amount of extra work compared to the FOSLS case. Additionally, adding a weight to the divergence equations in the FOSLS system typically increases the number of iterations required on the same order of magnitude as the weight. All matrices and vectors for the tests were constructed using the modular finite-element library MFEM [33].

Let A be the discretized FOSLS system and M be the mass matrix for \mathcal{P}_0 elements on the constraint space. Since A is SPD, the preconditioner

$$P = \begin{pmatrix} A & 0 \\ 0 & M \end{pmatrix}$$

is also SPD and, in the case of $a(\cdot, \cdot)$ being \mathbf{H}^1 -equivalent, P is known to be a uniform preconditioner to system (2.14). It should be noted that the Lagrange multiplier space can also be constructed on a coarse mesh to ensure stability and satisfy the conditions of the above theorems. Therefore, M will be the appropriate mass matrix on this coarsened mesh. Since \mathcal{P}_0 elements are used, M is diagonal and is inverted exactly. In contrast, A is not inverted exactly; rather it is replaced with an AMG preconditioner, A_{AMG} , namely, one BoomerAMG $V(1, 1)$ cycle is performed, which defines A_{AMG}^{-1} .

In the tables, several quantities are reported. These include the local mass conservation measured by mass loss, $ml_L = \sum_T |\int_T (\nabla \cdot \mathbf{u} - f) dT|$, which is a mesh-independent measure of the local mass conservation. Note that

$$\begin{aligned} ml_L &\leq \sum_T \int_T |\nabla \cdot \mathbf{u} - f| dT \leq \sum_T |T|^{1/2} \|\nabla \cdot \mathbf{u} - f\|_{0,T} \\ &\leq \left(\sum_T |T| \right)^{1/2} \left(\sum_T \|\nabla \cdot \mathbf{u} - f\|_{0,T} \right)^{1/2} = |\Omega|^{1/2} \|\nabla \cdot \mathbf{u} - f\|_{0,\Omega}. \end{aligned}$$

In addition, we measure the global mass loss, $ml_G = |\int_\Omega (\nabla \cdot \mathbf{u} - f) d\Omega|$, the FOSLS functional, \mathcal{F} , which is equivalent to the energy norm of the error in the solution, \mathbf{U} , the L^2 norm of the error in the solution, \mathbf{U}_{err} , if possible, and the number of

iterations needed to solve the system, IT. Simulations are performed for the pure FOSLS discretization using both bilinear, \mathcal{P}_1 , and biquadratic, \mathcal{P}_2 , elements on a uniform triangular mesh with spacing h . The constrained FOSLS minimization is also computed using these spaces for the unknowns plus a \mathcal{P}_0 space for the Lagrange multiplier, λ , on a mesh of size H . We compare results for coarsening ratios of $h/H = 1, 1/2, 1/4$, and $1/8$.

4.1. Diffusion problem. Here, we solve a simple diffusion problem using the FOSLS formulation, (2.6)–(2.8). The test problem is created by assuming the true solution, $p = \sin(\pi x)\sin(\pi y)$ on the unit square, $\Omega = [0, 1] \times [0, 1]$. Then, homogeneous Dirichlet boundary conditions are used for p on the boundaries, and the

TABLE 1
Diffusion problem (2.6)–(2.8) results using preconditioned MINRES on the saddle-point system, (2.14), reducing the residual to machine precision and AMG preconditioned PCG on the FOSLS system.

$\mathcal{P}_1 - \mathcal{P}_0$						$\mathcal{P}_2 - \mathcal{P}_0$				
h	ml_L	ml_G	\mathcal{F}	\mathbf{U}_{err}	IT	ml_L	ml_G	\mathcal{F}	\mathbf{U}_{err}	IT
FOSLS										
1/4	1.0e0	8.4e-1	3.8	3.7e-1	11	3.2e-2	1.5e-4	5.8e-1	1.7e-2	13
1/8	4.4e-1	2.2e-1	2.0	1.0e-1	14	3.2e-3	2.4e-5	1.5e-1	2.2e-3	15
1/16	2.1e-1	5.5e-2	9.9e-1	2.6e-2	16	3.9e-4	2.0e-6	3.8e-2	2.7e-4	16
1/32	1.0e-1	1.4e-2	5.0e-1	6.6e-3	17	4.8e-5	1.3e-7	9.6e-3	3.4e-5	18
1/64	5.1e-2	3.5e-3	2.5e-1	1.6e-3	17	6.0e-6	8.4e-9	2.4e-3	4.2e-6	20
1/128	2.6e-2	8.6e-4	1.2e-1	4.1e-4	18	7.5e-7	5.1e-10	6.0e-4	5.3e-7	22
1/256	1.3e-2	2.2e-4	6.2e-2	1.0e-4	19	9.3e-8	1.3e-12	1.5e-4	6.6e-8	28
Constrained FOSLS $h/H = 1$										
1/4	-	-	-	-	-	4.9e-15	2.2e-16	5.9e-1	1.7e-2	30
1/8	-	-	-	-	-	7.9e-15	4.0e-16	1.5e-1	2.1e-3	34
1/16	-	-	-	-	-	1.5e-14	5.5e-16	3.8e-2	2.7e-4	35
1/32	-	-	-	-	-	3.0e-14	1.3e-15	9.6e-3	3.4e-5	36
1/64	-	-	-	-	-	6.4e-14	1.5e-16	2.4e-3	4.2e-6	37
1/128	-	-	-	-	-	1.3e-13	1.8e-15	6.0e-4	5.3e-7	38
1/256	-	-	-	-	-	2.7e-13	1.1e-16	1.5e-4	6.6e-8	45
Constrained FOSLS $h/H = 1/2$										
1/4	1.2e-15	2.2e-16	3.9	2.8e-1	22	2.0e-15	6.6e-16	5.8e-1	1.7e-2	21
1/8	5.9e-15	3.6e-16	2.0	6.8e-2	36	2.4e-15	5.1e-16	1.5e-1	2.1e-3	25
1/16	8.2e-15	2.2e-16	1.0	1.7e-2	43	5.4e-15	8.0e-16	3.8e-2	2.7e-4	25
1/32	1.3e-14	1.5e-15	5.0e-1	4.2e-3	43	1.0e-14	7.9e-16	9.6e-3	3.4e-5	27
1/64	2.8e-14	7.5e-16	2.5e-1	1.0e-3	43	2.2e-14	1.2e-15	2.4e-3	4.2e-6	29
1/128	5.4e-14	1.1e-15	1.3e-1	2.6e-4	42	4.5e-14	9.0e-16	6.0e-4	5.3e-7	30
1/256	1.1e-13	5.9e-16	6.3e-2	6.5e-5	43	9.9e-14	1.0e-15	1.5e-4	6.6e-8	39
Constrained FOSLS $h/H = 1/4$										
1/8	1.7e-14	7.2e-16	2.0	7.6e-2	23	1.6e-15	2.2e-16	1.5e-1	2.2e-3	20
1/16	4.5e-15	8.7e-16	9.9e-1	1.7e-2	29	2.7e-15	4.5e-17	3.8e-2	2.7e-4	22
1/32	5.1e-15	2.8e-16	5.0e-1	4.2e-3	31	4.8e-15	6.5e-16	9.6e-3	3.4e-5	24
1/64	1.0e-14	3.9e-16	2.5e-1	1.0e-3	31	9.2e-15	3.4e-17	2.4e-3	4.2e-6	29
1/128	2.0e-14	1.3e-15	1.2e-1	2.6e-4	32	1.8e-14	7.1e-16	6.0e-4	5.3e-7	30
1/256	3.9e-14	7.2e-16	6.2e-2	6.5e-5	33	3.9e-14	5.1e-16	1.5e-4	6.6e-8	39
Constrained FOSLS $h/H = 1/8$										
1/16	3.3e-15	2.2e-16	1.0	2.0e-2	22	1.6e-15	1.6e-15	3.8e-2	2.7e-4	20
1/32	3.5e-15	1.1e-15	5.0e-1	4.4e-3	25	3.3e-15	2.4e-15	9.6e-3	3.4e-5	24
1/64	4.1e-15	5.6e-16	2.5e-1	1.1e-3	26	4.9e-15	3.1e-16	2.4e-3	4.2e-6	27
1/128	7.5e-15	2.9e-16	1.2e-1	2.6e-4	27	8.4e-15	3.2e-16	6.0e-4	5.3e-7	30
1/256	1.5e-14	1.7e-15	6.2e-2	6.5e-5	28	1.7e-14	8.3e-16	1.5e-4	6.6e-8	39

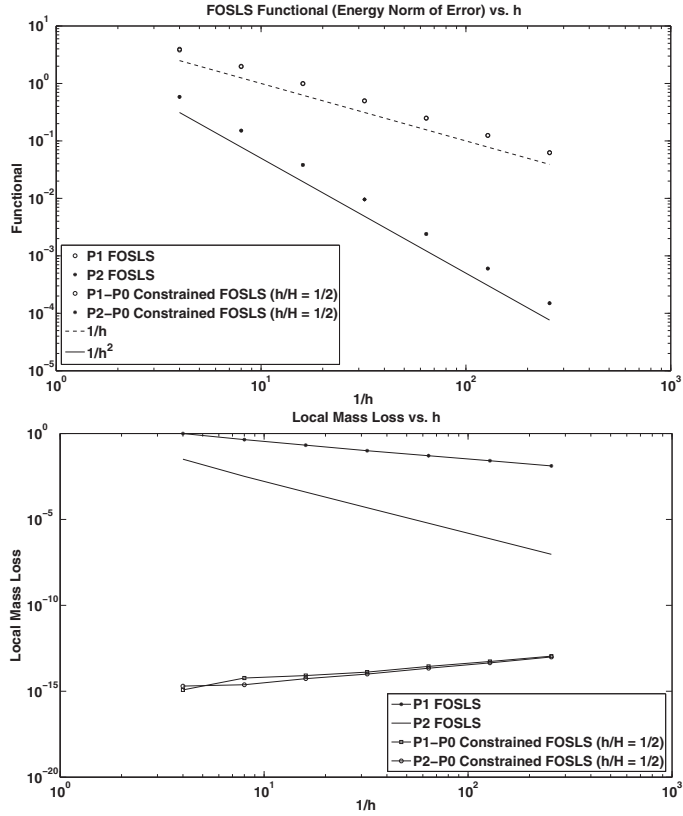


FIG. 1. Diffusion problem (2.6)–(2.8). Top: FOSLS functional (energy norm of error), \mathcal{F} , versus grid size. Bottom: Local mass loss, ml_L , versus grid size. Results for pure FOSLS solution and constrained FOSLS solution with $H = 2h$.

tangential components of $\mathbf{u} = \nabla p$ are zero as well. The right-hand side becomes $f = 2\pi^2 \sin(\pi x) \sin(\pi y)$.

The results in Table 1 show that FOSLS alone does not satisfy the divergence constraint especially accurately. As the resolution is increased, either by h-refinement or moving to biquadratics, the results do improve, but only because the overall accuracy is getting better. Adding the constraint to the system ensures that the divergence constraint is satisfied. The slight uptick in the local mass conservation with grid size is a result of numerical roundoff. Additionally, as one can see in the table and in Figure 1, the FOSLS functional (i.e., the energy norm of the error) is reduced at the appropriate rate. For biquadratics, the inf-sup condition is always satisfied and we see a quadratic decay in the energy norm error with mesh refinement as expected. For bilinears, the inf-sup is satisfied only when the constraint space is coarser (at least $h/H = 1/2$) and, thus, we see a linear decay in the energy norm error when this is satisfied.

4.2. Stokes problem. Next we consider a steady state Stokes flow in a unit cube, $\Omega = [0, 1] \times [0, 1]$. The velocity-vorticity-pressure formulation (2.10)–(2.12) is used. We assume that the normal components of the velocity are zero around the boundary, $\mathbf{n} \cdot \mathbf{u} = 0$, and the vorticity is zero everywhere except on the $x = 1$ boundary, where we define $p = 0$ in order to nail down the constant in the pressure.

Similar to the results for the diffusion equations, Table 2 and Figure 2 show that mass conservation is regained for inf-sup stable finite-element pairs between the FOSLS unknowns and the Lagrange multiplier, while the FOSLS functional still measures the energy norm of the error appropriately. Due to the boundary conditions, $\mathbf{n} \cdot \mathbf{u} = 0$, global mass conservation is guaranteed in both the pure FOSLS case and the constrained FOSLS case. However, only the constrained case gives local conservation immediately even for low resolution results.

Finally, to test the method on a more interesting problem, we consider steady-state Stokes flow down a long tube, $\Omega = [0, 8] \times [0, 1]$. In past work [6, 26, 27], it has been demonstrated that for problems where the inflow of the tube is prescribed, $u_1 = \sin(\pi x)$, along with zero boundary conditions for \mathbf{u} on the top and bottom, and a zero pressure boundary condition on the outflow, basic FOSLS formulations show a significant mass loss down the tube; see Figure 3. Applying the constrained formulation to system (2.10)–(2.12), however, shows that the mass is regained down the tube, as shown in Figure 3.

TABLE 2

Stokes problem, (2.10)–(2.12), results using preconditioned MINRES on the saddle-point system, (2.14), reducing the residual to machine precision and AMG preconditioned PCG on the FOSLS system.

$\mathcal{P}_1 - \mathcal{P}_0$					$\mathcal{P}_2 - \mathcal{P}_0$			
h	ml_L	ml_G	\mathcal{F}	it	ml_L	ml_G	\mathcal{F}	it
FOSLS								
1/4	5.6e-1	0.0	5.3	15	4.8e-3	2.8e-17	8.2e-1	15
1/8	3.6e-1	1.7e-17	2.8	16	7.6e-4	2.2e-18	2.1e-1	16
1/16	1.9e-1	1.4e-18	1.4	17	1.1e-4	6.5e-17	5.4e-2	17
1/32	9.8e-2	1.2e-17	7.0e-1	18	1.5e-5	2.4e-17	1.4e-2	20
1/64	4.9e-2	1.8e-17	3.5e-1	18	1.9e-6	2.3e-17	3.4e-3	22
1/128	2.5e-2	1.2e-17	1.8e-1	19	2.4e-7	7.8e-18	8.5e-4	25
Constrained FOSLS $h/H = 1$								
1/4	-	-	-	-	2.9e-15	8.3e-17	8.2e-1	27
1/8	-	-	-	-	4.1e-15	1.3e-17	2.1e-1	31
1/16	-	-	-	-	5.5e-15	6.8e-17	5.4e-2	32
1/32	-	-	-	-	9.7e-15	2.3e-17	1.4e-2	35
1/64	-	-	-	-	1.9e-14	1.9e-17	3.4e-3	37
1/128	-	-	-	-	3.9e-14	7.4e-18	8.5e-4	41
Constrained FOSLS $h/H = 1/2$								
1/4	4.0e-16	1.4e-17	5.4	22	8.9e-15	1.4e-17	8.2e-1	18
1/8	1.9e-15	9.1e-18	2.8	37	4.8e-15	1.3e-17	2.1e-1	22
1/16	3.3e-15	1.0e-19	1.4	40	2.6e-15	4.0e-17	5.4e-2	25
1/32	5.8e-15	3.8e-17	7.0e-1	40	3.2e-15	6.3e-17	1.4e-2	27
1/64	8.4e-15	7.9e-18	3.5e-1	40	8.1e-15	8.1e-18	3.4e-3	30
1/128	1.6e-14	2.1e-17	1.8e-1	41	1.4e-14	5.0e-17	8.5e-4	34
Constrained FOSLS $h/H = 1/4$								
1/8	2.9e-15	2.8e-17	2.8	22	9.8e-15	1.3e-17	2.1e-1	19
1/16	8.5e-15	1.5e-17	1.4	26	5.2e-15	1.1e-17	5.4e-2	22
1/32	2.5e-15	5.9e-18	7.0e-1	28	4.3e-15	5.6e-17	1.4e-2	25
1/64	3.7e-15	4.5e-17	3.5e-1	29	5.0e-15	2.8e-18	3.4e-3	29
1/128	6.2e-15	2.2e-17	1.8e-1	30	5.7e-15	1.1e-16	8.5e-4	34
Constrained FOSLS $h/H = 1/8$								
1/16	7.3e-16	2.8e-17	1.4	21	2.0e-15	3.9e-17	5.4e-2	21
1/32	5.6e-15	2.3e-17	7.0e-1	24	8.0e-16	8.7e-17	1.4e-2	25
1/64	4.9e-15	4.9e-18	3.5e-1	25	1.7e-15	1.1e-16	3.4e-3	29
1/128	2.1e-15	1.1e-17	1.8e-1	27	5.7e-15	1.5e-17	8.5e-4	32

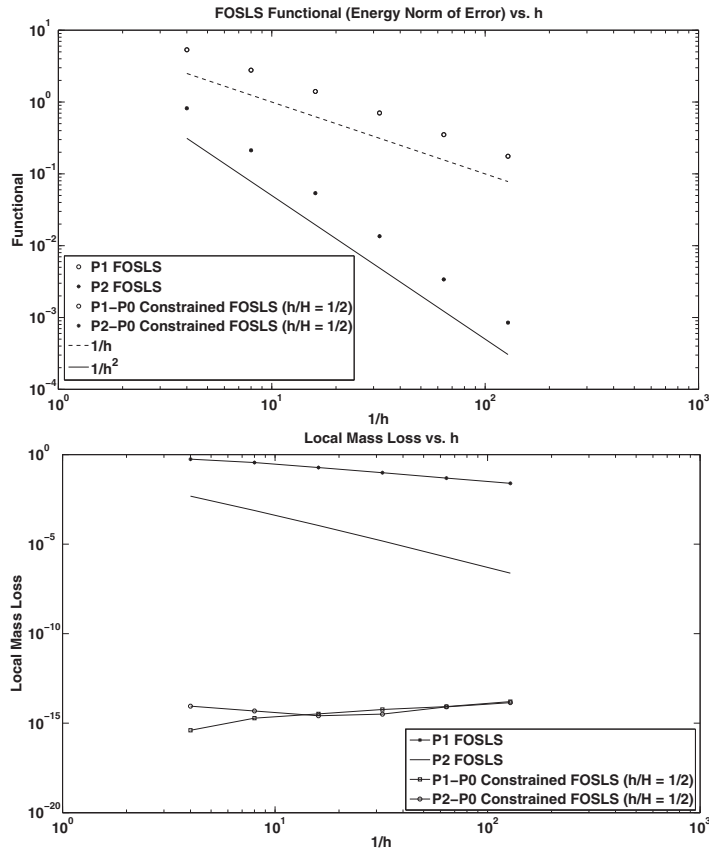


FIG. 2. Stokes problem (2.10)–(2.12). Top: FOSLS functional (energy norm of error), \mathcal{F} , versus grid size. Bottom: Local mass loss, m_L , versus grid size. Results for pure FOSLS solution and constrained FOSLS solution with $H = 2h$.

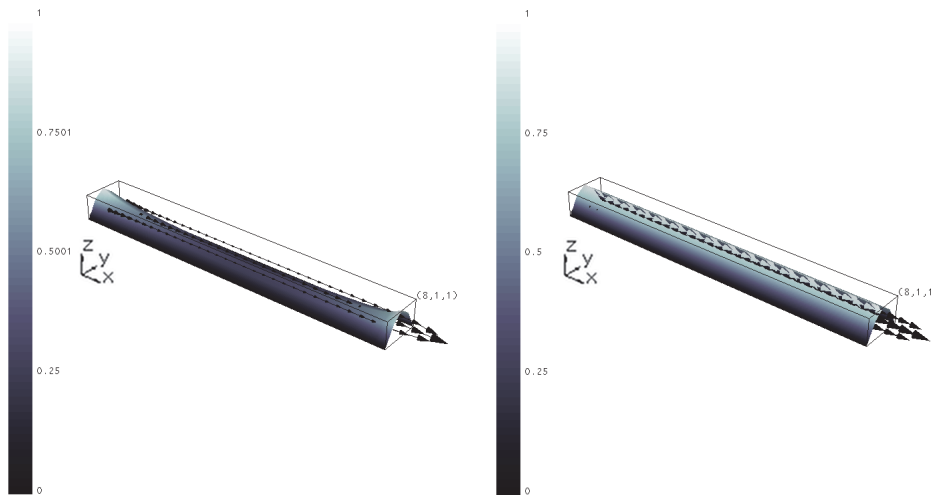


FIG. 3. Velocity field, \mathbf{u} , of Stokes; flow down a tube. Left: Pure FOSLS simulation with significant mass loss in the center of the tube. Right: Constrained FOSLS simulation, where the flow (mass) has been regained. \mathcal{P}_1 elements are used for the FOSLS unknowns and \mathcal{P}_0 for the Lagrange multiplier on a coarse mesh such that $H = 2h$.

5. Discussion. Along with [1], the above results show that conservation of a divergence constraint is readily obtained by solving the constrained saddle-point problem as described in (3.1) and (2.14). This paper expands upon the results by confirming an error estimate of the solution in the energy norm. In the context of a FOSLS finite-element formulation this guarantees that the FOSLS functional, \mathcal{F} , still maintains the property that it is an a posteriori error estimate in the energy norm when the constrained problem is solved. This assumes that appropriate spaces are chosen for the Lagrange multiplier space so that the inf-sup conditions are satisfied. In addition, this results in the need for solving a saddle-point problem. However, many advanced techniques have been developed to solve such problems and can be used with relatively little extra computational cost (especially when the constraint space is defined on a coarser mesh).

It should be noted that in all the test examples described above, a \mathcal{P}_0 space was used for the constraint. Looking at the bound in Theorem 3.3, one notices that the term, $\inf_{w \in \mathcal{W}_H} \|\lambda - w\|$ may imply a loss of order in convergence for the case of \mathcal{P}_2 elements for the FOSLS unknowns. However, this is only an upper bound and for the simple domains used (hence \mathbf{H}^1 -coercivity does hold, see Remark 3.2), the standard quadratic convergence of the functional is maintained. In general, using discontinuous \mathcal{P}_1 (or higher-order) elements for the constraint should bring the order of convergence back up.

Future work involves testing with higher-order spaces, as well as implementations for other test problems in which satisfying a divergence constraint is important. This includes applications such as incompressible magnetohydrodynamics, where a solenoidal constraint on the magnetic field, $\nabla \cdot \mathbf{B} = 0$, must be satisfied in addition to the fluid incompressibility [3, 11]. Also, this work can be extended to other formulations and other types of constraints. This includes considering more general Petrov–Galerkin type problems, where many constraints can be added to the system to ensure some kind of conservation properties that are of interest.

Acknowledgments. The authors thank Dr. Junping Wang for useful comments and suggestions as well the referees for their helpful remarks to improve the results.

REFERENCES

- [1] J. H. ADLER AND P. S. VASSILEVSKI, *Improving conservation for first-order system least squares finite-element methods*, in Numerical Solution of Partial Differential Equations: Theory, Algorithms, and Their Applications, O. P. Iliev, S. D. Margenov, P. D. Minev, P. S. Vassilevski, and L. T. Zikatanov, eds., Springer Proc. Math. Statist. 45, Springer, New York, 2013, pp. 1–19.
- [2] G. A. BAKER, W. N. JUREIDINI, AND O. A. KARAKASHIAN, *Piecewise solenoidal vector-fields and the Stokes problem*, SIAM J. Numer. Anal., 27 (1990), pp. 1466–1485.
- [3] T. BARTH, *On the role of involutions in the discontinuous Galerkin discretization of Maxwell and magnetohydrodynamic systems*, in Compatible Spatial Discretizations, Springer, New York, 2006, pp. 69–88.
- [4] M. BERNDT, T. A. MANTEUFFEL, AND S. F. MCCORMICK, *Local error estimates and adaptive refinement for first-order system least squares (FOSLS)*, Electron. Trans. Numer. Anal., 6 (1997), pp. 35–43.
- [5] P. B. BOCHEV, Z. CAI, T. A. MANTEUFFEL, AND S. F. MCCORMICK, *Analysis of velocity-flux first-order system least-squares principles for the Navier-Stokes equations: Part I*, SIAM J. Numer. Anal., 35 (1998), pp. 990–1009.
- [6] P. B. BOCHEV AND M. D. GUNZBURGER, *Analysis of least-squares finite-element methods for the Stokes equations*, Math. Comput., 63 (1994), pp. 479–506.
- [7] P. B. BOCHEV AND M. D. GUNZBURGER, *A locally conservative least-squares method for Darcy flows*, Commun. Numer. Methods Engrg. Biomed. Appl., 24 (2008), pp. 97–110.

- [8] P. B. BOCHEV, J. LAI, AND L. OLSON, *A locally conservative, discontinuous least-squares finite element method for the Stokes equations*, *Internat. J. Numer. Methods Fluids*, 68 (2011), pp. 782–804.
- [9] P. B. BOCHEV, J. LAI, AND L. OLSON, *A non-conforming least-squares finite element method for incompressible fluid flow problems*, *Internat. J. Numer. Methods Fluids*, 72 (2012), pp. 375–402.
- [10] P. B. BOCHEV, T. A. MANTEUFFEL, AND S. F. MCCORMICK, *Analysis of velocity-flux least-squares principles for the Navier-Stokes equations: Part II*, *SIAM J. Numer. Anal.*, 36 (1999), pp. 1125–1144.
- [11] J. U. BRACKBILL AND D. C. BARNES, *The effect of nonzero $\nabla \cdot \mathbf{B}$ on the numerical solution of the magnetohydrodynamic equations*, *J. Comput. Phys.*, 35 (1980), pp. 426–430.
- [12] J. H. BRAMBLE, T. V. KOLEV, AND J. E. PASCIAK, *A least-squares approximation method for the time-harmonic Maxwell equations*, *J. Numer. Math.*, 13 (2005), pp. 237–263.
- [13] S. C. BRENNER AND L. R. SCOTT, *Mathematical Theory of Finite Element Methods*, 2nd ed., Springer, New York, 2002.
- [14] M. BREZINA, J. GARCIA, T. MANTEUFFEL, S. MCCORMICK, J. RUGE, AND L. TANG, *Parallel adaptive mesh refinement for first-order system least squares*, *Numer. Linear Algebra Appl.*, 19 (2012), pp. 343–366.
- [15] F. BREZZI AND M. FORTIN, *Mixed and Hybrid Finite Elements Methods*, Springer Ser. Comput. Math., Springer, New York, 1991.
- [16] Z. CAI, R. LAZAROV, T. A. MANTEUFFEL, AND S. F. MCCORMICK, *First-order system least squares for second-order partial differential equations: Part I*, *SIAM J. Numer. Anal.*, (1994), pp. 1785–1799.
- [17] Z. CAI, T. A. MANTEUFFEL, AND S. F. MCCORMICK, *First-order system least squares for second-order partial differential equations: Part II*, *SIAM J. Numer. Anal.*, 34 (1997), pp. 425–454.
- [18] Z. CAI AND G. STARKE, *Least-squares methods for linear elasticity*, *SIAM J. Numer. Anal.*, 42 (2004), pp. 826–842.
- [19] G. F. CAREY, A. I. PEHLIVANOV, AND P. S. VASSILEVSKI, *Least-squares mixed finite element methods for non-selfadjoint elliptic problems II: Performance of block-ILU factorization methods*, *SIAM J. Sci. Comput.*, 16 (1995), pp. 1126–1136.
- [20] C. L. CHANG AND J. J. NELSON, *Least-squares finite element method for the Stokes problem with zero residual of mass conservation*, *SIAM J. Numer. Anal.*, 34 (1997), pp. 480–489.
- [21] B. COCKBURN, G. KANSCHAT, AND D. SCHOTZAU, *A locally conservative LDG method for the incompressible Navier-Stokes equations*, *Math. Comp.*, 74 (2005), pp. 1067–1095.
- [22] B. COCKBURN, G. KANSCHAT, AND D. SCHÖTZAU, *A note on discontinuous Galerkin divergence-free solutions of the Navier-Stokes equations*, *J. Sci. Comput.*, 31 (2007), pp. 61–73.
- [23] H. DE STERCK, T. MANTEUFFEL, S. MCCORMICK, J. NOLTING, J. RUGE, AND L. TANG, *Efficiency-based h- and hp-refinement strategies for finite element methods*, in *Numerical Linear Algebra with Applications*, University of Waterloo, Waterloo, ON, 2008, pp. 89–114.
- [24] V. GIRAULT AND P. A. RAVIART, *Finite Element Approximation of the Navier-Stokes Equations*, Springer, New York, 1979.
- [25] V. GIRAULT AND P. A. RAVIART, *Finite Element Methods for Navier-Stokes Equations: Theory and Algorithms*, Springer Ser. Comput. Math., Springer, New York, 1986.
- [26] J. J. HEYS, E. LEE, T. A. MANTEUFFEL, AND S. F. MCCORMICK, *On mass-conserving least-squares methods*, *SIAM J. Sci. Comput.*, 28 (2006), pp. 1675–1693.
- [27] J. J. HEYS, E. LEE, T. A. MANTEUFFEL, AND S. F. MCCORMICK, *An alternative least-squares formulation of the Navier-Stokes equations with improved mass conservation*, *J. Comput. Phys.*, 226 (2007), pp. 994–1006.
- [28] *High Performance Preconditioners*, <http://www.llnl.gov/CASC/hypre/>.
- [29] T. V. KOLEV AND P. S. VASSILEVSKI, *Regular decompositions for $h(\text{div})$ spaces.*, *Comput. Methods Appl. Math.*, 12 (2012), pp. 437–447.
- [30] P. D. LAX, *Functional Analysis*, Wiley-Interscience, New York, 2002.
- [31] E. LEE AND T. A. MANTEUFFEL, *FOSLL* method for the eddy current problem with three-dimensional edge singularities*, *SIAM J. Numer. Anal.*, 45 (2007), pp. 787–809.
- [32] T. A. MANTEUFFEL, S. F. MCCORMICK, J. RUGE, AND J. G. SCHMIDT, *First-order system LL* (FOSLL*) for general scalar elliptic problems in the plane*, *SIAM J. Numer. Anal.*, 43 (2006), pp. 2098–2120.
- [33] MFEM, <http://mfem.googlecode.com>.
- [34] S. MÜNZENMAIER AND G. STARKE, *First-order system least squares for coupled Stokes-Darcy flow*, *SIAM J. Numer. Anal.*, 49 (2011), pp. 387–404.

- [35] A. I. PEHLIVANOV AND G. F. CAREY, *Error-estimates for least-squares mixed finite-elements*, RAIRO Math. Model. Numer. Anal., 28 (1994), pp. 499–516.
- [36] A. I. PEHLIVANOV, G. F. CAREY, AND P. S. VASSILEVSKI, *Least-squares mixed finite element methods for non-selfadjoint elliptic problems I: Error analysis*, Numer. Math., 72 (1996), pp. 501–522.
- [37] J. P. PONTAZA AND J. N. REDDY, *Space-time coupled spectral/hp least-squares finite element formulation for the incompressible Navier-Stokes equations*, J. Comput. Phys., 197 (2004), pp. 418–459.
- [38] T. RUSTEN, P. S. VASSILEVSKI, AND R. WINTHER, *Interior penalty preconditioners for mixed finite element approximations of elliptic problems*, Math. Comp., 65 (1996), pp. 447–466.
- [39] Y. SAAD, *Iterative Methods for Sparse Linear Systems*, SIAM, Philadelphia, 2003.
- [40] G. STARKE, *A first-order system least squares finite element method for the shallow water equations*, SIAM J. Numer. Anal., 42 (2005), pp. 2387–2407.
- [41] A. TOSELLI AND O. B. WIDLUND, *Domain Decomposition Methods—Algorithms and Theory*, Springer, New York, 2005.
- [42] S. P. VANKA, *Block-implicit multigrid solution of Navier-Stokes equations in primitive variables*, J. Comput. Phys., 65 (1986), pp. 138–158.
- [43] P. S. VASSILEVSKI, *Multilevel block factorization preconditioners*, in Matrix-Based Analysis and Algorithms for Solving Finite Element Equations, Springer, New York, 2008.