

DECOUPLING THREE-DIMENSIONAL MIXED PROBLEMS USING DIVERGENCE-FREE FINITE ELEMENTS*

ROBERT SCHEICHL†

Abstract. In this paper we describe an iterative method for indefinite saddle-point systems arising from mixed finite element discretizations of second-order elliptic boundary value problems subject to mixed boundary conditions and posed over polyhedral three-dimensional domains. The method is based on a decoupling of the vector of velocities in the saddle-point system from the vector of pressures, resulting in a symmetric positive definite velocity system and a triangular pressure system.

The crucial step in this approach is the construction of the divergence-free Raviart–Thomas–Nédélec elements from the curls of Nédélec’s edge elements. Because of the large kernel of the curl-operator, this representation is not unique. To find a basis we consider the graph made up of the nodes and edges of the mesh and eliminate the edge elements associated with a spanning tree in this graph. To prove that this technique works in the general case considered here, we employ fundamental results from algebraic topology and graph theory.

We also include some numerical experiments, where we solve the (decoupled) velocity system by ILU-preconditioned conjugate gradients and the pressure system by simple back substitutions. We compare our method with a standard ILU-based block preconditioner for the original saddle-point system, and we find that our method is faster by a factor of at least 4.5 in all cases, with the greatest improvement occurring in the nonuniform mesh case.

Key words. mixed finite elements, second-order elliptic problems, mixed boundary conditions, divergence-free space, spanning trees, decoupled iterative method

AMS subject classifications. 65N22, 65N30, 65F10, 05C05, 55U10

PII. S1064827500375886

1. Introduction. The problem we are going to consider in this paper is the following second-order elliptic problem in velocity-pressure formulation

$$(1.1) \quad \vec{u} + K \vec{\nabla} p = \vec{g},$$

$$(1.2) \quad \vec{\nabla} \cdot \vec{u} = 0,$$

subject to mixed boundary conditions over a polyhedral three-dimensional domain Ω . Such a problem arises, for example, in groundwater flow or oil recovery simulations, where \vec{u} corresponds to the velocity, p corresponds to the pressure, and K is permeability divided by dynamic viscosity.

The numerical treatment of (1.1), (1.2) involves the solution of usually very large indefinite linear equation systems. In this paper we describe a very efficient and practicable iterative method to solve these systems by decoupling the vector of velocities from the vector of pressures, resulting in a symmetric positive definite velocity system and a triangular pressure system. The crucial step in this approach is the construction of a basis for the divergence-free Raviart–Thomas–Nédélec elements. The proof that our algorithm for this construction works uses results from algebraic topology and graph theory and will also be presented in this paper.

*Received by the editors July 24, 2000; accepted for publication (in revised form) February 15, 2001; published electronically January 30, 2002. This work was supported by EPSRC CASE Award 97D00023 in collaboration with AEA Technology.

<http://www.siam.org/journals/sisc/23-5/37588.html>

†Institut Français du Pétrole, DISMA, 1 et 4 ave. de Bois-Préau, 92852 Rueil-Malmaison, France (robert.scheichl@ifp.fr).

Because the variable of prime interest in (1.1), (1.2) (especially in the applications we have in mind) is the velocity \vec{u} , the discretization schemes of most interest are those which preserve conservation of mass (1.2) in an appropriate way, with the prime candidates being mixed finite element or finite volume techniques. In this paper we discretize (1.1), (1.2) using the lowest-order mixed Raviart–Thomas–Nédélec elements on tetrahedral meshes (Nédélec [23]). Here p is approximated in the space of piecewise constant functions and \vec{u} is approximated in an appropriate subspace of the vector-valued piecewise linear functions, in which the normal component of \vec{u} is required to be continuous across the element boundaries. The resulting discretization enforces mass conservation on each element of the mesh. Since the quality of the approximations is determined by the mesh width, it is usually necessary to work with very fine meshes.

As is well known this type of discretization yields symmetric indefinite systems of *saddle-point type*. Iterative methods for indefinite systems are less powerful and robust (w.r.t. a refinement of the mesh) than the methods available for definite systems. Therefore almost all approaches to solve this system efficiently contain at some point a reduction of the system to a symmetric positive definite system. At least two different strategies have been pursued.

The first strategy is to solve the saddle-point system by a preconditioned *minimum residual* (MINRES) method using a symmetric positive definite block preconditioner. The analysis for this strategy seems to be restricted to the two-dimensional case (e.g., [26, 4]), although the preconditioner described in Rusten and Winther [26] can be readily applied in three dimensions as well. There is also recent work by Wohlmuth, Toselli, and Widlund [29] on a domain decomposition preconditioner for Raviart–Thomas–Nédélec vector fields in three dimensions which can be used in the framework of Arnold, Falk, and Winther [4].

The second strategy is to decouple the vector of velocities in the saddle-point system from the vector of pressures. This can be done by (i) mixed hybridization via Lagrange multipliers [13, 11], (ii) block elimination of the velocity variable [24], or, as presented in this paper, (iii) a direct elimination of the divergence constraint (1.2) on the element level. This technique (iii) has two distinct advantages over (i) and (ii): first, no nonphysical variables are introduced, and, second, the velocity is obtained directly without (necessarily) computing the pressure, the latter advantage being particularly attractive in groundwater flow calculations. The method (iii) was first developed in the related but different case of the Stokes problem by Crouzeix and Thomasset [28] for two dimensions and by Hecht [17] for three dimensions. In connection with the solution of system (1.1), (1.2) it appears first in Chavent et al. [10]. Since the decoupling in this way leads to much smaller and nicer linear equation systems, it was subsequently possible to develop very competitive and efficient methods for the two-dimensional case (e.g., [15, 16, 21, 22, 12]). In this paper we present a very efficient iterative method for the solution of system (1.1), (1.2) that implements this idea in three dimensions.

Thus our solver is built on three essential steps. The first step decouples the velocity field in the saddle-point problem from the pressure field. This is done by writing the velocity as the curl of an appropriate discrete *vector potential*, automatically satisfying the discrete counterpart of the mass conservation law (1.2). The required discrete vector potential turns out to be a finite element approximation of the solution of a related symmetric positive *semidefinite* problem by *edge elements* (Nédélec [23]) which can be found independently of the pressure. Because of the large kernel of the curl-operator, this system is singular. In Hiptmair and Hoppe [19] a multilevel method is constructed that solves this singular problem approximately. In

this paper we make the problem positive definite by eliminating the degrees of freedom associated with a *spanning tree* of the graph made up of the nodes and edges of the mesh. This also corresponds to finding a local basis for the divergence-free Raviart–Thomas–Nédélec elements. Such algebraic techniques have already been successfully used in other fields (e.g., [2, 20, 6] for Maxwell’s equations, [17, 18, 14] for incompressible flow (Stokes)); in the context of (1.1), (1.2) there is only one paper by Cai et al. [8], and that is restricted to uniform rectangular meshes, in which case the special spanning tree can be written down a priori.

The second step in the solver is the application of a preconditioned *conjugate gradient* (CG) method to solve for the discrete velocity. Since the system is symmetric positive definite, it is easier to solve than the original saddle-point system. It also turns out to have the added advantage of being about three times smaller than the original system. In section 7 we have included some results using a simple ILU preconditioner to illustrate this advantage. We compare our method with the (also ILU-based) block preconditioner of Rusten and Winther [26] for the original saddle-point system, and we find that our method is faster by a factor of at least 4.5 in all cases, with the greatest improvement occurring in the nonuniform mesh case (where an improvement factor of 9.0 is observed on the finest mesh (36864 freedoms)).

The third and final step in the solver is the recovery of the pressure (if it is required). The decoupled pressure system turns out to be particularly simple. By an appropriate numbering of the freedoms it can be made triangular and solved in optimal time by simple *back substitutions*. We will prove this rigorously.

The layout of this paper is as follows. In section 2 we shall describe the mathematical setting for (1.1), (1.2) and its discretization by mixed finite elements. In section 3 we present the decoupling of the vector of velocities from the vector of pressures as a general algebraic procedure. In section 4 we construct the basis for the divergence-free Raviart–Thomas–Nédélec elements, and we show in section 5 how this basis is used to implement the decoupled velocity system. In section 6 we show how the pressure is recovered, and we finish the paper with some numerical results in section 7.

2. Mixed finite element discretization. In this section we describe the mathematical setting for (1.1), (1.2) together with its discretization by mixed finite elements. Since this is a standard procedure (see, e.g., [7]), we shall be brief.

Let Ω denote an open *polyhedron*, i.e. a simply connected open domain without cavities in \mathbb{R}^3 with a connected boundary Γ composed of plane faces, which is assumed partitioned into $\Gamma_D \cup \Gamma_N$. Each of Γ_D and Γ_N is assumed to consist of a finite union of planar polygonal subsets of Γ , and Γ_D and Γ_N are both assumed to be connected. Additionally, Γ_N is assumed to be closed. Let $\vec{\nu}(\vec{x})$ denote the outward unit normal from Ω at $\vec{x} \in \Gamma$. In general we assume that K is a bounded, symmetric, and uniformly positive definite 3×3 matrix-valued function on Ω . The system (1.1), (1.2) is to be solved on Ω subject to mixed boundary conditions:

$$(2.1) \quad p = p_D \quad \text{on} \quad \Gamma_D \quad \text{and} \quad \vec{u} \cdot \vec{\nu} = 0 \quad \text{on} \quad \Gamma_N.$$

Throughout, we shall assume that $\Gamma_D \neq \emptyset$, a condition which is generically satisfied in groundwater flow applications, where some inflow and outflow must occur. The extension to the case when $\Gamma_D = \emptyset$ (when p is nonunique) can easily be made by imposing an extra condition on p in the weak form below (e.g., that p should have a prescribed mean value; see, e.g., [15]).

To discretize (1.1), (1.2), (2.1) we put it in weak form. Let $(\cdot, \cdot)_{L_2(\Omega)^d}$ denote the

usual inner product in $L_2(\Omega)^d$ for $d = 1, 2, 3$. Then introduce the Hilbert space

$$H(\operatorname{div}, \Omega) := \{\vec{v} \in L_2(\Omega)^3 : \operatorname{div} \vec{v} \in L_2(\Omega)\},$$

with the inner product

$$(\vec{u}, \vec{v})_{H(\operatorname{div}, \Omega)} := (\vec{u}, \vec{v})_{L_2(\Omega)^3} + (\operatorname{div} \vec{u}, \operatorname{div} \vec{v})_{L_2(\Omega)},$$

and its subspace

$$H_{0,N}(\operatorname{div}, \Omega) := \{\vec{v} \in H(\operatorname{div}, \Omega) : \vec{v} \cdot \vec{\nu}|_{\Gamma_N} = 0\}$$

(see [7] for details). Introduce the bilinear forms

$$m(\vec{u}, \vec{v}) := (K^{-1}\vec{u}, \vec{v})_{L_2(\Omega)^3}, \quad b(\vec{v}, w) := -(\operatorname{div} \vec{v}, w)_{L_2(\Omega)},$$

and the linear functional

$$G(\vec{v}) := (K^{-1}\vec{g}, \vec{v})_{L_2(\Omega)^3} - \int_{\Gamma_D} p_D \vec{v} \cdot \vec{\nu} dF.$$

Then the weak form of (1.1), (1.2), (2.1) is to find $(\vec{u}, p) \in H_{0,N}(\operatorname{div}, \Omega) \times L_2(\Omega)$ such that

$$(2.2) \quad \begin{cases} m(\vec{u}, \vec{v}) + b(\vec{v}, p) = G(\vec{v}) & \text{for all } \vec{v} \in H_{0,N}(\operatorname{div}, \Omega), \\ b(\vec{u}, w) = 0 & \text{for all } w \in L_2(\Omega). \end{cases}$$

The mixed finite element discretization of (2.2) is obtained by choosing finite-dimensional subspaces $\mathcal{V} \subset H_{0,N}(\operatorname{div}, \Omega)$ and $\mathcal{W} \subset L_2(\Omega)$ and seeking $(\vec{U}, P) \in \mathcal{V} \times \mathcal{W}$ such that

$$(2.3) \quad \begin{cases} m(\vec{U}, \vec{V}) + b(\vec{V}, P) = G(\vec{V}) & \text{for all } \vec{V} \in \mathcal{V}, \\ b(\vec{U}, W) = 0 & \text{for all } W \in \mathcal{W}. \end{cases}$$

In practice this is implemented by choosing bases $\{\vec{v}_i : i = 1, \dots, n_{\mathcal{V}}\}$ and $\{w_j : j = 1, \dots, n_{\mathcal{W}}\}$ for \mathcal{V} and \mathcal{W} . By writing

$$\vec{U} = \sum_{i=1}^{n_{\mathcal{V}}} u_i \vec{v}_i, \quad P = \sum_{j=1}^{n_{\mathcal{W}}} p_j w_j,$$

problem (2.3) is then reduced to the indefinite system of linear equations

$$(2.4) \quad \begin{pmatrix} M & B \\ B^T & 0 \end{pmatrix} \begin{pmatrix} \mathbf{u} \\ \mathbf{p} \end{pmatrix} = \begin{pmatrix} \mathbf{g} \\ \mathbf{0} \end{pmatrix} \quad \text{in } \mathbb{R}^{n_{\mathcal{V}}} \times \mathbb{R}^{n_{\mathcal{W}}},$$

where $M_{i,i'} := m(\vec{v}_i, \vec{v}_{i'})$ is the “mass matrix,” $B_{i,j} := b(\vec{v}_i, w_j)$ is the “discrete gradient,” and $g_i := G(\vec{v}_i)$.

In this paper we restrict attention to the (most practically important) case when \mathcal{V} is the lowest-order Raviart–Thomas–Nédélec space on tetrahedra [23]. To define this, let \mathcal{T} denote a triangulation of Ω into conforming tetrahedra $T \in \mathcal{T}$. We assume that all lines along which the boundary condition changes (i.e., the boundaries of the components of Γ_N) are edges of tetrahedra in \mathcal{T} . Let \mathcal{F} denote the set of all faces of the tetrahedra in \mathcal{T} . It is convenient to think of these faces as open so that, for

$F \in \mathcal{F}$, \overline{F} denotes the closure of F (including its boundary). For any $F \in \mathcal{F}$, we let $\vec{\nu}_F$ denote the unit normal to the face F which, for convenience, is assumed to be orientated so that

$$(2.5) \quad \vec{\nu}_F \in \{\vec{x} \in \mathbb{R}^3 : x_1 > 0\} \cup \{(0, x_2, x_3)^T \in \mathbb{R}^3 : x_2 > 0\} \cup \{(0, 0, 1)^T\}.$$

Let \mathcal{F}_I , \mathcal{F}_D , and \mathcal{F}_N denote the faces $F \in \mathcal{F}$ which lie in Ω , Γ_D , and Γ_N , respectively.

The space \mathcal{V} is defined to be the space of all functions $\vec{v} \in H_{0,N}(\text{div}, \Omega)$ such that for all $T \in \mathcal{T}$, there exist $\vec{\alpha}_T \in \mathbb{R}^3$ and $\gamma_T \in \mathbb{R}$ such that

$$(2.6) \quad \vec{v}(\vec{x}) = \vec{\alpha}_T + \gamma_T \vec{x} \quad \text{for all } \vec{x} \in T.$$

Equivalently, we can define \mathcal{V} to be the space of all $\vec{v} : \Omega \rightarrow \mathbb{R}^3$ which satisfy (2.6) for each $T \in \mathcal{T}$, and also

$$(2.7) \quad \begin{aligned} \text{(i)} \quad & \vec{v} \cdot \vec{\nu}_F \text{ is continuous across each face } F \in \mathcal{F}_I, \\ \text{(ii)} \quad & \vec{v} \cdot \vec{\nu}_F = 0 \quad \text{for all } F \in \mathcal{F}_N. \end{aligned}$$

Because of the special form of (2.6) it is easily shown that $\vec{v}(\vec{x}) \cdot \vec{\nu}_F$ is constant for $\vec{x} \in F$ on any face F of T . Thus $\vec{v} \in \mathcal{V}$ can be completely determined by specifying the constant value of $\vec{v} \cdot \vec{\nu}_F$ for each $F \in \mathcal{F}_I \cup \mathcal{F}_D$. This leads us to introduce the standard basis for \mathcal{V} which is constructed by associating with each face $F \in \mathcal{F}_I \cup \mathcal{F}_D$, a function $\vec{v}_F \in \mathcal{V}$ with the property that

$$(2.8) \quad \vec{v}_F \cdot \vec{\nu}_{F'} = \delta_{F,F'},$$

with δ denoting the Kronecker delta.

We also have to specify the space \mathcal{W} . To fulfill the discrete inf-sup condition which is necessary for existence and uniqueness (see, e.g., [25]), \mathcal{W} is chosen as the space of piecewise constant functions on Ω , with the basis consisting of the characteristic functions w_T of each of the tetrahedra $T \in \mathcal{T}$. Thus

$$(2.9) \quad n_{\mathcal{V}} = (\#\mathcal{F}_I + \#\mathcal{F}_D), \quad n_{\mathcal{W}} = (\#\mathcal{T}),$$

where, throughout, $\#A$ denotes the number of elements of a (finite) set A .

3. Decoupled iterative method for mixed problems. In this section we formulate our method for decoupling the vector of velocities \mathbf{u} from the vector of pressures \mathbf{p} in system (2.4). We have already presented this procedure for the two-dimensional case in [12]. Recall [7] that (2.4) has a unique solution $(\mathbf{u}, \mathbf{p}) \in \mathbb{R}^{n_{\mathcal{V}}} \times \mathbb{R}^{n_{\mathcal{W}}}$ for all $\mathbf{g} \in \mathbb{R}^{n_{\mathcal{V}}}$, and clearly \mathbf{u} is in $\ker B^T$.

REMARK 3.1. *The case of $\mathbf{u} \notin \ker B^T$ (or, equivalently, of a more general right-hand side $(\mathbf{g}^T, \mathbf{h}^T)^T$ of (2.4)) arises when $\vec{\nabla} \cdot \vec{u} \neq 0$ in (1.2). This problem can be reduced to problem (2.4) following Ewing and Wang [15] (see also [8, 19] for three dimensions): in a preprocessing step based on domain decomposition (static condensation) a vector \mathbf{u}^* is calculated such that $B^T \mathbf{u}^* = \mathbf{h}$; this can be done in $O(n)$ steps (where $n = n_{\mathcal{V}} + n_{\mathcal{W}}$); the remainder $\tilde{\mathbf{u}} = \mathbf{u} - \mathbf{u}^*$ fulfills (2.4) with right-hand side $((\mathbf{g} - M\mathbf{u}^*)^T, \mathbf{0}^T)^T$ and can be calculated with the method described in this paper. See [27] for details.*

To describe our decoupling procedure, first consider (2.4) as an abstract system. The decoupling of \mathbf{u} from \mathbf{p} can be achieved by finding

$$(3.1) \quad \text{a basis } \{\mathbf{z}_1, \dots, \mathbf{z}_{\tilde{n}}\} \text{ of } \ker B^T.$$

(Since B^T has full rank, $\mathring{n} = n_{\mathcal{V}} - n_{\mathcal{W}}$.) If we have such a basis, then the solution \mathbf{u} of (2.4) can be written

$$(3.2) \quad \mathbf{u} = \sum_{j=1}^{\mathring{n}} \mathring{u}_j \mathbf{z}_j = Z^T \mathring{\mathbf{u}},$$

for some $\mathring{\mathbf{u}} \in \mathbb{R}^{\mathring{n}}$, where Z denotes the $\mathring{n} \times n_{\mathcal{V}}$ matrix with rows $\mathbf{z}_1^T, \dots, \mathbf{z}_{\mathring{n}}^T$. Also, since $ZB = (B^T Z^T)^T = 0$, multiplying the first (block) row of (2.4) by Z shows that $\mathring{\mathbf{u}}$ is a solution of the linear system

$$(3.3) \quad \mathring{A} \mathring{\mathbf{u}} = \mathring{\mathbf{g}},$$

where

$$(3.4) \quad \mathring{A} = ZMZ^T \quad \text{and} \quad \mathring{\mathbf{g}} = Z\mathbf{g}.$$

Since M is symmetric positive definite, so is \mathring{A} , and $\mathring{\mathbf{u}}$ is the unique solution of (3.3). Thus if the basis (3.1) can be found, then the velocity \mathbf{u} in (2.4) can be computed by solving the decoupled positive definite system (3.3) rather than the indefinite coupled system (2.4).

In applications to groundwater flow, where one is primarily interested in the velocity $\vec{\mathbf{u}}$ in (1.1), (1.2), the method described above is of great relevance. Even when the pressure p is also of interest our method may still be highly competitive, provided we can also compute a *complementary basis* $\{\mathbf{z}_{\mathring{n}+1}, \dots, \mathbf{z}_{n_{\mathcal{V}}}\}$ with the property that

$$(3.5) \quad \text{span}\{\mathbf{z}_1, \dots, \mathbf{z}_{\mathring{n}}, \mathbf{z}_{\mathring{n}+1}, \dots, \mathbf{z}_{n_{\mathcal{V}}}\} = \mathbb{R}^{n_{\mathcal{V}}}.$$

If this is known and if Z' denotes the matrix with rows $\mathbf{z}_{\mathring{n}+1}^T, \dots, \mathbf{z}_{n_{\mathcal{V}}}^T$, then multiplying the first (block) row of (2.4) by Z' shows that \mathbf{p} is the solution of the $n_{\mathcal{W}} \times n_{\mathcal{W}}$ system

$$(3.6) \quad (Z'B)\mathbf{p} = Z'(\mathbf{g} - M\mathbf{u}).$$

An elementary argument shows that $Z'B$ is nonsingular, and so the unique solution \mathbf{p} of (3.6) also determines the pressure in (2.4) once the velocity \mathbf{u} is known.

We show in the next three sections that in the particular case of the mixed finite element system (2.4),

- (i) it is always easy to find the basis (3.1);
- (ii) the resulting symmetric positive definite matrix \mathring{A} in the reduced problem (3.3) can be obtained by simple algebraic techniques from the stiffness matrix of an associated symmetric positive semidefinite problem in the space $H(\text{curl}, \Omega)$ discretized by Nédélec's edge elements;
- (iii) the system (3.3) is about 3 times smaller than (2.4);
- (iv) a simple choice of complementary basis can be made so that the coefficient matrix $Z'B$ in the system (3.6) is lower triangular.

To establish conclusions (i)–(iv) we need to exploit the particular properties of (2.4). In particular, note that finding the basis $\mathbf{z}_1, \dots, \mathbf{z}_{\mathring{n}}$ in (3.1) is equivalent to finding a basis $\vec{\mathbf{v}}_1, \dots, \vec{\mathbf{v}}_{\mathring{n}}$ of the finite element space

$$\mathring{\mathcal{V}} := \{\vec{\mathbf{V}} \in \mathcal{V} : b(\vec{\mathbf{V}}, W) = 0 \text{ for all } W \in \mathcal{W}\}.$$

To see why, suppose $\mathbf{z}_1, \dots, \mathbf{z}_{\mathring{n}}$ are known and let $Z = (Z_{i,j})$ be the matrix with rows $\mathbf{z}_1^T, \dots, \mathbf{z}_{\mathring{n}}^T$. Then the formulae

$$(3.7) \quad \vec{\mathbf{v}}_i = \sum_{j=1}^{n_{\mathcal{V}}} Z_{i,j} \vec{\mathbf{v}}_j, \quad i = 1, \dots, \mathring{n},$$

(where $\{\vec{v}_j\}$ is the basis of \mathcal{V}) determine the basis $\{\vec{v}_i\}$. Conversely, if the basis $\{\vec{v}_i\}$ of $\mathring{\mathcal{V}}$ is known, then the matrix Z (and hence the basis $\mathbf{z}_1, \dots, \mathbf{z}_{\tilde{n}}$ of $\ker B^T$) is determined by (3.7).

We now turn our attention to finding a basis of $\mathring{\mathcal{V}}$.

4. Construction of a divergence-free basis. As a first step, recall that \mathcal{T} is the set of all tetrahedra in the mesh and that $\mathcal{F} = \mathcal{F}_I \cup \mathcal{F}_D \cup \mathcal{F}_N$ is the set of all faces of the mesh (assumed to be open triangles) which lie in Ω , Γ_D , and Γ_N , respectively. Analogously, we can write $\mathcal{E} = \mathcal{E}_I \cup \mathcal{E}_D \cup \mathcal{E}_N$, with \mathcal{E}_I , \mathcal{E}_D , and \mathcal{E}_N denoting the edges in Ω , Γ_D , and Γ_N ; and $\mathcal{N} = \mathcal{N}_I \cup \mathcal{N}_D \cup \mathcal{N}_N$, with \mathcal{N}_I , \mathcal{N}_D , and \mathcal{N}_N denoting the nodes in Ω , Γ_D , and Γ_N . Recall that the boundaries of each of the components of Γ_N belong to Γ_N , and, since the lines between Neumann and Dirichlet boundaries are edges of the mesh, these edges lie in \mathcal{E}_N . For $E \in \mathcal{E}$, let $\vec{\tau}_E$ denote the unit tangent on edge E which, as in (2.5), is assumed to be orientated so that

$$(4.1) \quad \vec{\tau}_E \in \{\vec{x} \in \mathbb{R}^3 : x_1 > 0\} \cup \{(0, x_2, x_3)^T \in \mathbb{R}^3 : x_2 > 0\} \cup \{(0, 0, 1)^T\}.$$

Through this convention we associate an orientation with each edge of the mesh.

To construct a basis of $\mathring{\mathcal{V}}$ it is useful to introduce the following space of finite elements introduced by Nédélec in [23]. Let

$$H(\vec{\text{curl}}, \Omega) := \{\vec{\Phi} \in L_2(\Omega)^3 : \vec{\text{curl}} \vec{\Phi} \in L_2(\Omega)^3\}$$

and let \mathcal{U} be the finite-dimensional space of all functions $\vec{\Phi} \in H(\vec{\text{curl}}, \Omega)$ such that for all $T \in \mathcal{T}$, there exist $\vec{\alpha}_T, \vec{\beta}_T \in \mathbb{R}^3$ such that

$$(4.2) \quad \vec{\Phi}(\vec{x}) = \vec{\alpha}_T + \vec{\beta}_T \times \vec{x} \quad \text{for all } \vec{x} \in T.$$

In fact, \mathcal{U} is the lowest-order member of the family of spaces introduced by Nédélec in [23]. The standard basis of \mathcal{U} consists of the set of functions $\{\vec{\Phi}_E \in \mathcal{U} : E \in \mathcal{E}\}$ which are required to have the property

$$(4.3) \quad \int_{E'} \vec{\Phi}_E \cdot \vec{\tau}_{E'} ds = \delta_{E, E'} \quad \text{for all } E' \in \mathcal{E}.$$

This choice of basis functions accounts for the widely used term *edge elements*.

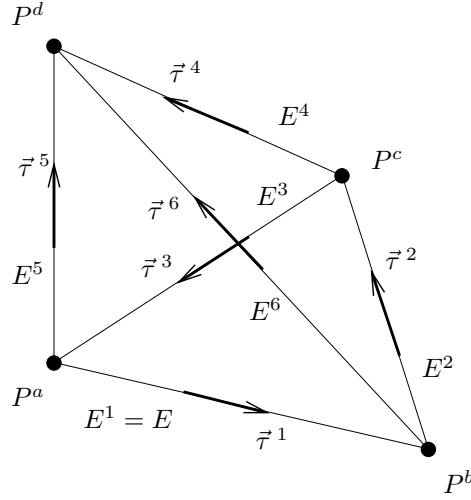
The basis for $\mathring{\mathcal{V}}$ will now be constructed from the fundamental functions $\vec{\Psi}_E$ defined by

$$(4.4) \quad \vec{\Psi}_E = \vec{\text{curl}} \vec{\Phi}_E$$

(so that $\vec{\Phi}_E$ is the *vector potential* of $\vec{\Psi}_E$). The functions (4.4) clearly satisfy $\text{div} \vec{\Psi}_E = 0$ on each tetrahedron of the mesh, and a subset of them lie in $\mathring{\mathcal{V}}$ as the following proposition shows.

PROPOSITION 4.1. *For each $E \in \mathcal{E}_I \cup \mathcal{E}_D$, $\vec{\Psi}_E \in \mathring{\mathcal{V}}$.*

Proof. Consider a general edge $E \in \mathcal{E}$. Conditions (4.2) and (4.3) clearly imply that $\text{supp } \vec{\Psi}_E$ consists only of the tetrahedra touching edge E . A typical such tetrahedron T with edges $E := E^1, E^2, \dots, E^6$, and nodes P^a, \dots, P^d , is depicted in the figure below:



with $\vec{\tau}^\alpha$, $\alpha = 1, \dots, 6$, denoting unit tangent vectors in the directions shown and \vec{r}^β denoting the position vector of P^β , $\beta = a, b, c, d$. The faces are denoted by F^a, \dots, F^d , where F^β is opposite P^β , $\beta = a, b, c, d$, and the unit outward normal on each face is denoted by $\vec{\nu}^\beta$.

Note first that for all $\vec{x} \in T$

$$\vec{\Psi}_E(\vec{x}) = \text{curl} \vec{\Phi}_E(\vec{x}) = \vec{\nabla} \times (\vec{\beta}_T \times \vec{x}) = 2\vec{\beta}_T,$$

which is easily seen to be of the form (2.6) (in fact, with $\gamma_T = 0$).

Since $\vec{\Psi}_E(\vec{x})$ is constant on T , we can write for each $\vec{x} \in T$ and for each $\beta = a, b, c, d$,

$$\vec{\Psi}_E(\vec{x}) \cdot \vec{\nu}^\beta = \frac{1}{|F^\beta|} \int_{F^\beta} \vec{\Psi}_E(\vec{x}) \cdot \vec{\nu}^\beta dF = \frac{1}{|F^\beta|} \int_{F^\beta} \text{curl} \vec{\Phi}_E(\vec{x}) \cdot \vec{\nu}^\beta dF,$$

and using Stokes's integral theorem we get

$$(4.5) \quad \vec{\Psi}_E(\vec{x}) \cdot \vec{\nu}^\beta = \frac{1}{|F^\beta|} \oint_{\partial F^\beta} \vec{\Phi}_E(\vec{x}) \cdot d\vec{s} = \begin{cases} \frac{1}{|F^\beta|} & \text{for } \beta = c, \\ -\frac{1}{|F^\beta|} & \text{for } \beta = d, \\ 0 & \text{otherwise,} \end{cases}$$

where in the last step we used (4.3) to evaluate the line integral (respecting the right-hand rule and the specific orientation of $\vec{\tau}^1$ and $\vec{\nu}^\beta$, $\beta = a, b, c, d$, as depicted in the figure above).

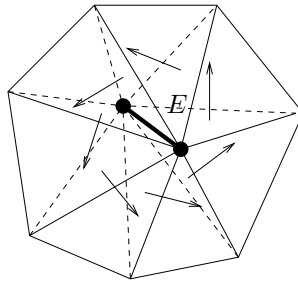
Now to obtain the result observe that, since $\text{div} \vec{\Psi}_E = 0$ on each tetrahedron, it is sufficient to show that

$$(4.6) \quad \vec{\Psi}_E \in \mathcal{V} \quad \text{for all } E \in \mathcal{E}_I \cup \mathcal{E}_D.$$

To show this we shall verify criterion (2.7). First consider $E \in \mathcal{E}_I$. Let $F \in \mathcal{F}_I$. If $F \not\subset \text{supp } \vec{\Psi}_E$, then we have trivially

$$(4.7) \quad \vec{\Psi}_E \cdot \vec{\nu}_F \quad \text{is continuous across } F.$$

Now take a general tetrahedron $T \subset \text{supp } \vec{\Psi}_E$, as pictured above. If $F = F^c$ or F^d , then performing the computation (4.5) in the other tetrahedron adjoining F and

FIG. 1. Divergence-free basis function $\vec{\Psi}_E$.

combining with (4.5) establishes (4.7). On the other hand, when $F = F^a$ or F^b , (4.7) also holds since $\vec{\Psi}_E \cdot \vec{\nu}_F|_T = 0$ and since the other tetrahedron adjoining F lies outside $\text{supp } \vec{\Psi}_E$. Altogether, we have established that $\vec{\Psi}_E$ satisfies criterion (2.7)(i).

To establish (2.7)(ii), let $F \in \mathcal{F}_N$. If $F \not\subset \text{supp } \vec{\Psi}_E$, then $\vec{\Psi}_E \cdot \vec{\nu}_F = 0$ trivially. If $F \subset \overline{T} \subset \text{supp } \vec{\Psi}_E$, then (since $E \in \mathcal{E}_I$) with the above notation F has to be either F^a or F^b and again $\vec{\Psi}_E \cdot \vec{\nu}_F = 0$, proving (2.7)(ii).

Thus we have shown that $\vec{\Psi}_E \in \mathcal{V}$ for all $E \in \mathcal{E}_I$. Similar arguments establish that $\vec{\Psi}_E \in \mathcal{V}$ for all $E \in \mathcal{E}_D$, proving (4.6). \square

Note that each $\vec{\Psi}_E$ can be expressed as a local linear combination of the basis functions $\vec{\nu}_F$ of \mathcal{V} satisfying (2.8); in fact, only those $\vec{\nu}_F$ corresponding to faces F that contain edge E appear in the expansion of $\vec{\Psi}_E$ (see Figure 1).

To find a basis for $\mathring{\mathcal{V}}$, let us first look at the pure Dirichlet case, $\Gamma_N = \emptyset$. The functions introduced in Proposition 4.1 are sufficient to span $\mathring{\mathcal{V}}$, but there are too many of them. The following theorem identifies a linearly independent subset of the functions in Proposition 4.1 that constitutes a basis of $\mathring{\mathcal{V}}$. A similar statement for the pure Neumann case, $\Gamma_D = \emptyset$, has already been proved by Dubois [14] (see Remark 4.7 for a more extensive literature survey).

The proof involves some fundamental notions and results from graph theory and algebraic topology (see Appendices A and B for a brief introduction). In particular we need the notion of a *spanning tree* of a graph (see Theorem A.4). Let $\mathbf{G} := (\mathcal{N}, \mathcal{E})$ be the graph formed by the nodes and (orientated) edges of the triangulation \mathcal{T} .

THEOREM 4.2. *Let $\Gamma_N = \emptyset$ and let $\mathcal{H} \subset \mathcal{E}$ be such that if $\mathbf{H} := (\mathcal{N}, \mathcal{H})$ is a spanning tree of \mathbf{G} , then*

$$(4.8) \quad \{\vec{\Psi}_E : E \in \mathcal{E} \setminus \mathcal{H}\} \text{ is a basis of } \mathring{\mathcal{V}}.$$

Before proving Theorem 4.2, we will first prove two lemmas. Let $\mathcal{V}(\mathbf{G})$ denote the vector space over \mathbb{Z} generated by the cycles of \mathbf{G} as defined in Definition A.1(e). Furthermore, for each face $F \in \mathcal{F}$ let μ^F be the *elementary cycle* of \mathbf{G} formed by the edges E of F . We fix the orientation of this cycle w.r.t. $\vec{\nu}_F$ by applying the right-hand rule. The associated vector $\boldsymbol{\mu}^F := [\mu_E^F]_{E \in \mathcal{E}} \in \mathcal{V}(\mathbf{G})$ is given by

$$(4.9) \quad \mu_E^F = \begin{cases} 1 & \text{if } E \text{ is an edge of } F \text{ and } \vec{\tau}_E \text{ is positively orientated w.r.t. } \vec{\nu}_F, \\ -1 & \text{if } E \text{ is an edge of } F \text{ and } \vec{\tau}_E \text{ is negatively orientated w.r.t. } \vec{\nu}_F, \\ 0 & \text{otherwise.} \end{cases}$$

LEMMA 4.3. Let $\boldsymbol{\mu} := [\mu_E]_{E \in \mathcal{E}} \in \mathcal{V}(\mathbf{G})$. Then there exist $\{\alpha_F \in \mathbb{Z} : F \in \mathcal{F}\}$ such that

$$(4.10) \quad \boldsymbol{\mu} := \sum_{F \in \mathcal{F}} \alpha_F \boldsymbol{\mu}^F.$$

Proof. Let K be the simplicial complex underlying our simplicial triangulation \mathcal{T} . In the notation of algebraic topology (see Appendix B) the vector $\boldsymbol{\mu} \in \mathcal{V}(\mathbf{G})$ can be identified with the vector of coefficients of a cycle μ of K (with orientation of its edges defined by the tangent vectors $\vec{\tau}_E$).

Since $|K| = \bar{\Omega}$ is simply connected, we know from Corollary B.4 that each cycle of K is a bounding cycle and can therefore be written as a linear combination of the boundaries of the orientated triangles of K . In particular, there exist $\{\tilde{\alpha}_F \in \mathbb{Z} : F \in \mathcal{F}\}$ such that

$$(4.11) \quad \mu = \sum_{F \in \mathcal{F}} \tilde{\alpha}_F \partial F.$$

The boundary ∂F of an orientated triangle F of K is a special cycle $\tilde{\mu}^F$ of K . As above it can therefore be identified with a vector $\tilde{\boldsymbol{\mu}}^F \in \mathcal{V}(\mathbf{G})$. Depending on the orientation of $\vec{\nu}_F$ we either have $\tilde{\boldsymbol{\mu}}^F = \boldsymbol{\mu}^F$ or $\tilde{\boldsymbol{\mu}}^F = -\boldsymbol{\mu}^F$, and we can write (4.11) in vector notation:

$$\boldsymbol{\mu} = \sum_{F \in \mathcal{F}} \alpha_F \boldsymbol{\mu}^F \quad \text{with} \quad \alpha_F = \begin{cases} \tilde{\alpha}_F & \text{if } \tilde{\boldsymbol{\mu}}^F = \boldsymbol{\mu}^F, \\ -\tilde{\alpha}_F & \text{if } \tilde{\boldsymbol{\mu}}^F = -\boldsymbol{\mu}^F. \end{cases} \quad \square$$

LEMMA 4.4. Let $\boldsymbol{\mu} \in \mathcal{V}(\mathbf{G})$ and let $\{\alpha_F \in \mathbb{Z} : F \in \mathcal{F}\}$ be such that $\boldsymbol{\mu} := \sum_{F \in \mathcal{F}} \alpha_F \boldsymbol{\mu}^F$. Then

$$(4.12) \quad \sum_{F \in \mathcal{F}} \alpha_F \int_F \vec{\Psi}_E \cdot \vec{\nu}_F dF = \mu_E \quad \text{for all } E \in \mathcal{E}.$$

Proof. Let $F \in \mathcal{F}$. Using (4.5) we get

$$\int_F \vec{\Psi}_E \cdot \vec{\nu}_F dF = \begin{cases} 1 & \text{if } E \subset \overline{F} \text{ and } \vec{\tau}_E \text{ positively orientated w.r.t. } \vec{\nu}_F, \\ -1 & \text{if } E \subset \overline{F} \text{ and } \vec{\tau}_E \text{ negatively orientated w.r.t. } \vec{\nu}_F, \\ 0 & \text{otherwise,} \end{cases}$$

and, therefore, recalling the definition (4.9), we have $\int_F \vec{\Psi}_E \cdot \vec{\nu}_F dF = \mu_E^F$. Multiplying this by α_F and summing over $F \in \mathcal{F}$ we obtain (4.12). \square

We can now prove Theorem 4.2.

Proof. Let us first check that the number of basis functions in (4.8) coincides with $\hat{n} = \dim \hat{\mathcal{V}}$. Since Ω is simply connected without cavities, we can apply *Euler's polyhedron theorem* [9]

$$(4.13) \quad \#\mathcal{N} - \#\mathcal{E} + \#\mathcal{F} - \#\mathcal{T} = 1$$

to the triangulation \mathcal{T} . Now observe that \mathbf{H} is a tree, and therefore $\#\mathcal{H} = \#\mathcal{N} - 1$ (cf. Theorem A.3(iii)). Using this fact together with (4.13) we get

$$(4.14) \quad \#(\mathcal{E} \setminus \mathcal{H}) = \#\mathcal{E} - \#\mathcal{N} + 1 = \#\mathcal{F} - \#\mathcal{T}.$$

Now recalling that $\mathring{n} = n_{\mathcal{V}} - n_{\mathcal{W}} = \#\mathcal{F} - \#\mathcal{T}$, it follows from (4.14) that the number of functions in (4.8) is \mathring{n} , as required.

To establish linear independency of the functions in (4.8), suppose $\{\beta_{E'} : E' \in \mathcal{E} \setminus \mathcal{H}\}$ are scalars such that

$$\vec{0} = \sum_{E' \in \mathcal{E} \setminus \mathcal{H}} \beta_{E'} \vec{\Psi}_{E'}.$$

Now let $E \in \mathcal{E} \setminus \mathcal{H}$ and let $\boldsymbol{\mu}^E$ denote the vector associated with the unique cycle μ^E generated by taking edge E into the tree \mathbf{H} , which has the property that $\mu_{E'}^E := \delta_{E,E'}$ for all $E' \in \mathcal{E} \setminus \mathcal{H}$ (cf. Theorem A.6). Then using Lemma 4.3 we can find $\{\alpha_F \in \mathbb{Z} : F \in \mathcal{F}\}$ such that $\boldsymbol{\mu}^E := \sum_{F \in \mathcal{F}} \alpha_F \boldsymbol{\mu}^F$, and so by Lemma 4.4

$$0 = \sum_{F \in \mathcal{F}} \alpha_F \int_F \left(\sum_{E' \in \mathcal{E} \setminus \mathcal{H}} \beta_{E'} \vec{\Psi}_{E'} \right) \cdot \vec{\nu}_F dF = \sum_{E' \in \mathcal{E} \setminus \mathcal{H}} \beta_{E'} \mu_{E'}^E = \beta_E,$$

which establishes the linear independency of the functions in (4.8). \square

Now let us look at mixed boundary conditions, $\Gamma_N \neq \emptyset$. In the following corollary we will see that the results of Theorem 4.2 extend to this case provided each component of Γ_N is simply connected. Our proof of this result makes use of the methods of Hecht [17] developed for the nonconforming P1-P0 elements for the approximation of solenoidal vector fields in $H^1(\Omega)^3$ (see Remark 4.7 for a more extensive discussion).

Therefore let n_C denote the number of connected components in Γ_N and write

$$\Gamma_N = \Gamma_N^1 \cup \Gamma_N^2 \cup \cdots \cup \Gamma_N^{n_C}, \quad \Gamma_N^\ell \cap \Gamma_N^{\ell'} = \emptyset \quad \text{for all } \ell \neq \ell' \in \{1, \dots, n_C\}.$$

For $\ell = 1, \dots, n_C$, let $\mathcal{N}_N^\ell \subset \mathcal{N}$, $\mathcal{E}_N^\ell \subset \mathcal{E}$, and $\mathcal{F}_N^\ell \subset \mathcal{F}$ denote the set of mesh nodes, edges, and faces on Γ_N^ℓ , respectively.

COROLLARY 4.5. *Suppose $n_C \neq 0$ and suppose that Γ_N^ℓ is simply connected for each $\ell = 1, \dots, n_C$. Let $\mathcal{H} \subset \mathcal{E}$ such that $\mathbf{H} = (\mathcal{N}, \mathcal{H})$ is a spanning tree of \mathbf{G} and such that for each $\ell = 1, \dots, n_C$, the restriction $\mathbf{H}_N^\ell := (\mathcal{N}_N^\ell, \mathcal{H} \cap \mathcal{E}_N^\ell)$ of \mathbf{H} to Γ_N^ℓ is also a tree. Then*

$$(4.15) \quad \{\vec{\Psi}_E : E \in (\mathcal{E}_I \cup \mathcal{E}_D) \setminus \mathcal{H}\} \quad \text{is a basis of } \mathring{\mathcal{V}}.$$

REMARK 4.6. *The general case, when Γ_N^ℓ is not simply connected for some ℓ , involves the introduction of a small number of additional nonlocal basis functions. In order not to complicate this paper, we omit the details for this case, but they will be given in [27]. Thus, from now on we will assume that Γ_N^ℓ is simply connected for all $\ell = 1, \dots, n_C$.*

Proof. Since $(\mathcal{E}_I \cup \mathcal{E}_D) \subset \mathcal{E}$, we also have $(\mathcal{E}_I \cup \mathcal{E}_D) \setminus \mathcal{H} \subset \mathcal{E} \setminus \mathcal{H}$, and therefore following the proof of Theorem 4.2 the functions $\vec{\Psi}_E$ in (4.15) are linearly independent.

We have only to check that the number of basis functions in (4.15) coincides with $\mathring{n} = \dim \mathring{\mathcal{V}}$. To do this, we need two elementary formulae (Cauchy [9]): *Euler's polyhedron theorem* (4.13) and the *Euler–Cauchy formula* for planar networks of polygons (4.16). Consider a typical Neumann boundary segment Γ_N^ℓ . Since Γ_N^ℓ is simply connected, we have

$$(4.16) \quad \#\mathcal{N}_N^\ell - \#\mathcal{E}_N^\ell + \#\mathcal{F}_N^\ell = 1.$$

Now observe that \mathbf{H}_N^ℓ is a tree, and therefore (again by virtue of Theorem A.3(iii)) $\#(\mathcal{H} \cap \mathcal{E}_N^\ell) = \#\mathcal{N}_N^\ell - 1$. Using (4.16) and summing over $\ell = 1, \dots, n_C$, we obtain

$$(4.17) \quad \#(\mathcal{H} \cap \mathcal{E}_N) = \sum_{\ell=1}^{n_C} (\#\mathcal{N}_N^\ell - 1) = \sum_{\ell=1}^{n_C} (\#\mathcal{E}_N^\ell - \#\mathcal{F}_N^\ell) = (\#\mathcal{E}_N - \#\mathcal{F}_N).$$

Since the sets \mathcal{E}_I , \mathcal{E}_D , and \mathcal{E}_N partition \mathcal{E} , we also have

$$(\mathcal{E}_I \cup \mathcal{E}_D) \setminus \mathcal{H} = (\mathcal{E} \setminus \mathcal{E}_N) \setminus \mathcal{H} = \mathcal{E} \setminus (\mathcal{E}_N \cup \mathcal{H}),$$

and, therefore, the number of functions in (4.15) is $\#\mathcal{E} - (\#\mathcal{E}_N + \#\mathcal{H} - \#(\mathcal{H} \cap \mathcal{E}_N))$. Combining this with (4.17), and using the fact that \mathbf{H} is a tree (and therefore $\#\mathcal{H} = \#\mathcal{N} - 1$), we finally get

$$(4.18) \quad \#((\mathcal{E}_I \cup \mathcal{E}_D) \setminus \mathcal{H}) = (\#\mathcal{E} - \#\mathcal{N} + 1) - \#\mathcal{F}_N = (\#\mathcal{F} - \#\mathcal{T}) - \#\mathcal{F}_N,$$

where in the last step we have used Euler's polyhedron theorem (4.13). Now recalling that $\mathring{n} = n_V - n_W$ (and from section 3 we have $n_V = \#\mathcal{F}_I + \#\mathcal{F}_D = \#\mathcal{F} - \#\mathcal{F}_N$ and $n_W = \#\mathcal{T}$), it follows from (4.18) that the number of functions in (4.15) is \mathring{n} , as required. \square

REMARK 4.7. *The idea of spanning trees in the context of finite element methods first appears in the context of the Stokes problem in a paper by Hecht [17], where it is used in the same way as here to find a basis for the space of divergence-free nonconforming P1-P0 elements for the approximation of solenoidal vector fields in $H^1(\Omega)^3$.*

In an unpublished manuscript [18], Hecht extends these results to a wider family of finite elements in $H^1(\Omega)^3$, including the (nonconforming) Raviart–Thomas–Nédélec elements. The published literature on divergence-free Raviart–Thomas–Nédélec elements in $H(\text{div}, \Omega)$ considered here is restricted to the pure Neumann case, $\Gamma_D = \emptyset$, in a paper by Dubois [14], where he uses it to solve model incompressible flow problems with prescribed vorticity.

In the context of the three-dimensional problem (1.1), (1.2) considered in this paper, the only other work which we are aware of is the recent paper [8], but this is restricted to uniform rectangular meshes and a special spanning tree which can be constructed a priori.

Independently, spanning trees also appear as a technique for computing a discrete gauge condition in eddy-current calculations in computational electromagnetism (e.g., in Albanese and Rubinacci [2] or Kettunen and Turner [20]). A thorough presentation of the theoretical foundation of those techniques using homology theory (related to our Appendix B) can be found in Bossavit [6, Ch. 5].

5. Implementation. To implement the decoupled system (3.3) for determining $\hat{\mathbf{u}}$ (and hence \mathbf{u}) we must work with the matrix $\hat{\mathbf{A}}$ and the right-hand side $\hat{\mathbf{g}}$ specified in (3.4). We observe that these are formally defined in terms of multiplications with the matrix Z which, through (3.7), represents the basis $\{\vec{v}_i\}$ of $\mathring{\mathcal{V}}$ in terms of the basis $\{\vec{v}_j\}$ of \mathcal{V} .

In the specific system (2.4) the $\{\vec{v}_j\}$ are the Raviart–Thomas velocity basis functions $\{\vec{v}_F : F \in \mathcal{F}_I \cup \mathcal{F}_D\}$ given in section 2, whereas the $\{\vec{v}_i\}$ are the basis functions $\{\vec{\Psi}_E : E \in (\mathcal{E}_I \cup \mathcal{E}_D) \setminus \mathcal{H}\}$ specified in Corollary 4.5 (or Theorem 4.2, if $\Gamma_N = \emptyset$). Thus we can identify the columns of Z with the indices $F \in \mathcal{F}_I \cup \mathcal{F}_D$, whereas the rows of Z correspond to $E \in (\mathcal{E}_I \cup \mathcal{E}_D) \setminus \mathcal{H}$.

Using this identification of Z we can rewrite (3.7) as

$$(5.1) \quad \vec{\Psi}_E = \sum_{F \in \mathcal{F}_I \cup \mathcal{F}_D} Z_{E,F} \vec{v}_F, \quad E \in (\mathcal{E}_I \cup \mathcal{E}_D) \setminus \mathcal{H}.$$

Note that the matrix Z is sparse; in fact, $Z_{E,F} \neq 0$ only when edge E is an edge of the face F .

The set $\mathcal{H} \subset \mathcal{E}$ of edges that form a spanning tree in the graph $\mathbf{G} = (\mathcal{N}, \mathcal{E})$ can be found in optimal time (proportional to the number of edges) using Algorithm A.7 presented in Appendix A. The introduction of the special spanning tree for the mixed boundary problem in Corollary 4.5 does not pose an extra problem to our method. We need only to modify Algorithm A.7 slightly: We choose $x_1 \in \mathcal{N}_N$ and at first consider only nodes $y_1 \in \mathcal{N}_N$ in the function “*recursive(.)*” to find a spanning tree \mathbf{H}_N^ℓ for each Γ_N^ℓ ; then (without resetting the array “*mark[.]*”) we call the function “*recursive(\tilde{x})*” with argument $\tilde{x} \in \mathcal{N}_I \cup \mathcal{N}_D$ to find the rest of the spanning tree.

With the same convention as above we can write the elements of the matrix M appearing in (2.4) as

$$(5.2) \quad M_{F,F'} = m(\vec{v}_F, \vec{v}_{F'}), \quad F, F' \in \mathcal{F}_I \cup \mathcal{F}_D.$$

With these observations, it is simple to write \mathring{A} as a sum of element matrices. To be precise, recalling that \mathcal{T} is the set of tetrahedral elements, we can write

$$M = \sum_{T \in \mathcal{T}} M_T, \quad \text{where } (M_T)_{F,F'} = \int_T K^{-1} \vec{v}_F \cdot \vec{v}_{F'} d\vec{x}.$$

Then we can similarly write \mathring{A} as

$$(5.3) \quad \mathring{A} = \sum_{T \in \mathcal{T}} \mathring{A}_T,$$

where $\mathring{A}_T = Z_T M_T Z_T^T$, and Z_T denotes the matrix whose entries equal the entries of Z for columns and rows corresponding to T (i.e., faces $F \subset \overline{T}$ and edges $E \subset \overline{T}$) and are zero elsewhere. The representation (5.3) may be important if iterative methods are being used to solve (3.3). A similar elementwise representation can be given for the computation of \mathring{g} in (3.4).

Alternatively, \mathring{A} can be determined (elementwise or globally) from an approximation of a related bilinear form by Nédélec’s edge elements, without the assembly of any Raviart–Thomas stiffness matrix entries, as the following calculation shows.

Introduce the bilinear form

$$(5.4) \quad a(\vec{\Phi}, \vec{\Phi}') := (K^{-1} \operatorname{curl} \vec{\Phi}, \operatorname{curl} \vec{\Phi}')_{L^2(\Omega)^3} \quad \text{for all } \vec{\Phi}, \vec{\Phi}' \in H(\operatorname{curl}, \Omega),$$

and, for $E, E' \in \mathcal{E}$, set

$$\mathcal{A}_{E,E'} := a(\vec{\Phi}_E, \vec{\Phi}_{E'}),$$

where $\{\vec{\Phi}_E\}$ are the basis functions of the piecewise linear Nédélec’s edge elements defined in (4.2) and (4.3). Thus (after specifying an ordering of the edges in \mathcal{E}), \mathcal{A} is the stiffness matrix corresponding to the bilinear form $a(\cdot, \cdot)$ discretized by Nédélec’s edge elements, with a natural boundary condition on all of Γ . Because of the nontrivial

kernel of $a(\cdot, \cdot)$, this bilinear form is degenerate and therefore not elliptic on $H(\text{curl}, \Omega)$. In fact, let $v \in H^1(\Omega)$; then $a(\vec{\nabla} v, \vec{\Phi}') = 0$ for all $\vec{\Phi}' \in H(\text{curl}, \Omega)$. Consequently, \mathcal{A} is singular.

The following result shows that the minor of this matrix obtained by restricting to $E, E' \in (\mathcal{E}_I \cup \mathcal{E}_D) \setminus \mathcal{H}$, where $\mathcal{H} \subset \mathcal{E}$ as defined in Corollary 4.5, determines the matrix \mathring{A} in (3.3). (This corresponds to imposing an essential boundary condition on Γ_N and restricting to the orthogonal complement of the kernel of $a(\cdot, \cdot)$.)

THEOREM 5.1. *Let $\mathcal{H} \subset \mathcal{E}$ as defined in Corollary 4.5 (or Theorem 4.2, if $\Gamma_N = \emptyset$). Then*

$$\mathring{A}_{E,E'} = \mathcal{A}_{E,E'} \quad \text{for all } E, E' \in (\mathcal{E}_I \cup \mathcal{E}_D) \setminus \mathcal{H}.$$

Proof. Let $\mathcal{H} \subset \mathcal{E}$ as defined in Corollary 4.5 (or Theorem 4.2, if $\Gamma_N = \emptyset$) and let $E, E' \in (\mathcal{E}_I \cup \mathcal{E}_D) \setminus \mathcal{H}$. Then using the definition (3.4) of \mathring{A} together with (5.2) and (5.1) we get

$$\mathring{A}_{E,E'} = \sum_{F,F' \in \mathcal{F}_I \cup \mathcal{F}_D} Z_{E,F} M_{F,F'} Z_{E',F'} = m \left(\sum_F Z_{E,F} \vec{v}_F, \sum_{F'} Z_{E',F'} \vec{v}_{F'} \right) = m(\vec{\Psi}_E, \vec{\Psi}_{E'}).$$

Now using the definitions of $m(\cdot, \cdot)$ and $\vec{\Psi}_E$, we finally get

$$\begin{aligned} \mathring{A}_{E,E'} &= (K^{-1} \vec{\Psi}_E, \vec{\Psi}_{E'})_{L^2(\Omega)^3} = (K^{-1} \text{curl} \vec{\Phi}_E, \text{curl} \vec{\Phi}_{E'})_{L^2(\Omega)^3} \\ &= a(\vec{\Phi}_E, \vec{\Phi}_{E'}) = \mathcal{A}_{E,E'} \quad \square \end{aligned}$$

REMARK 5.2. *Hiptmair and Hoppe [19] solve the singular symmetric positive semidefinite system with stiffness matrix \mathcal{A} by multilevel preconditioned CG without explicitly eliminating columns and rows corresponding to edges $E \in \mathcal{H}$. In their multilevel splitting they eliminate the kernel of $a(\cdot, \cdot)$ only approximately by relaxing the orthogonality condition and thus avoid the construction of a basis. Here we eliminate the kernel a priori, which allows us to then apply the CG algorithm with a range of possible preconditioners.*

REMARK 5.3. *Observe that the decoupled system (3.3) is about three times smaller than the original indefinite system (2.4). More precisely, the dimension of (3.3) is smaller than that of (2.4) by a factor*

$$C := \frac{\#\mathcal{F}_I + \#\mathcal{F}_D + \#\mathcal{T}}{\#\mathcal{F}_I + \#\mathcal{F}_D - \#\mathcal{T}}.$$

Since $4(\#\mathcal{T}) = 2(\#\mathcal{F}_I) + \#\mathcal{F}_D + \#\mathcal{F}_N$ we have

$$C = 3 \left\{ \frac{\#\mathcal{F}_I + \frac{5}{6}(\#\mathcal{F}_D) + \frac{1}{6}(\#\mathcal{F}_N)}{\#\mathcal{F}_I + \frac{3}{2}(\#\mathcal{F}_D) - \frac{1}{2}\#\mathcal{F}_N} \right\}.$$

Under reasonable mesh regularity assumptions $\#\mathcal{F}_I$ is the dominant part of $\#\mathcal{F}$ as $\#\mathcal{T} \rightarrow \infty$, and so $C \rightarrow 3$ as $\#\mathcal{T} \rightarrow \infty$.

6. Pressure computations. In this section we present a procedure for the efficient recovery of the pressure \mathbf{p} from the decoupled system (3.6).

In the general situation described in section 3, the assembly of (3.6) requires the computation of a complementary basis $\{\mathbf{z}_{\hat{n}+1}, \dots, \mathbf{z}_{n_\nu}\}$ satisfying (3.5). This is again equivalent to finding a complementary basis $\{\vec{v}_{\hat{n}+1}^c, \dots, \vec{v}_{n_\nu}^c\}$ to $\{\vec{v}_1, \dots, \vec{v}_{\hat{n}}\}$ such that

$$(6.1) \quad \text{span} \left\{ \vec{v}_1, \dots, \vec{v}_{\hat{n}}, \vec{v}_{\hat{n}+1}^c, \dots, \vec{v}_{n_\nu}^c \right\} = \mathcal{V}.$$

In the context of the specific system (2.4) (using lowest-order Raviart–Thomas–Nédélec elements) this can be done by finding a distinguished subset of faces

$$\mathcal{F}^c \subset \mathcal{F}_I \cup \mathcal{F}_D$$

such that the corresponding subset of Raviart–Thomas–Nédélec basis functions, $\{\vec{v}_F : F \in \mathcal{F}^c\}$, constitutes a complementary basis. Note that this set must contain $n_{\mathcal{W}} := n_{\mathcal{V}} - \hat{n} = \#\mathcal{T}$ elements. The following simple Algorithm 6.1 chooses $n_{\mathcal{W}}$ appropriate faces, yielding a complementary basis, in such a way that the system (3.6) has a particularly simple form. We have already presented this algorithm to find a subset of faces $\mathcal{F}^c \subset \mathcal{F}_I \cup \mathcal{F}_D$ for the two-dimensional case in [12], but here we will give a rigorous proof that the corresponding subset of Raviart–Thomas–Nédélec basis functions, $\{\vec{v}_F : F \in \mathcal{F}^c\}$, constitutes a complementary basis to (4.15) (or to (4.8), if $\Gamma_N = \emptyset$) in \mathcal{V} .

ALGORITHM 6.1.

1. Choose $T_1 \in \mathcal{T}$ to be any tetrahedron with a face $F_1 \in \mathcal{F}_D$ and set $\mathcal{F}^c = \{F_1\}$.
2. For $j = 2, \dots, n_{\mathcal{W}}$,
 - choose $T_j \in \mathcal{T} \setminus \{T_\ell : \ell = 1, \dots, j-1\}$ with the property that there exists $F_j \in \mathcal{F}_I$ such that

$$(6.2) \quad F_j \subset \overline{T}_j \cap \left\{ \bigcup_{\ell=1}^{j-1} \overline{T}_\ell \right\}$$

- update $\mathcal{F}^c = \mathcal{F}^c \cup \{F_j\}$.
- End of loop over j .
3. Assemble $Z'B$ as

$$(Z'B)_{i,j} = b(\vec{v}_{F_i}, w_{T_j}), \quad i, j = 1, \dots, n_{\mathcal{W}}.$$

THEOREM 6.2. Algorithm 6.1 is well defined and the matrix $Z'B$ given in step 3 is lower triangular.

Proof. Since $\Gamma_D \neq \emptyset$, there exists a $T \in \mathcal{T}$ with a face $F \in \mathcal{F}_D$. Let T be T_1 . Now, assume we have found $j-1 < n_{\mathcal{W}} = \#\mathcal{T}$ tetrahedra T_ℓ in step 2 that fulfill property (6.2). Since Ω is connected, there exists a tetrahedron $T \in \mathcal{T}$ that has a face in common with $\bigcup_{\ell=1}^{j-1} \overline{T}_\ell$. Let T be T_j . The existence of a set \mathcal{F}^c therefore follows by an inductive argument.

Second, let $i, j = 1, \dots, n_{\mathcal{W}}$ with $i < j$. By the algorithm $F_i \subset \overline{T}_i \cap \{\bigcup_{\ell=1}^{i-1} \overline{T}_\ell\}$ and therefore $F_i \not\subset \overline{T}_j$. Since T_j is not in $\text{supp } \vec{v}_{F_i}$, it follows that $(Z'B)_{i,j} = b(\vec{v}_{F_i}, w_{T_j}) = 0$, and the matrix $Z'B$ given in step 3 is lower triangular. \square

To show that Algorithm 6.1 yields a complementary basis, let us first consider the pure Dirichlet case, $\Gamma_N = \emptyset$.

THEOREM 6.3. Let $\Gamma_N = \emptyset$. The functions

$$(6.3) \quad \{\vec{v}_F : F \in \mathcal{F}^c\}$$

form a complementary basis to (4.8) in \mathcal{V} .

Before proving this theorem we will first prove two lemmas again.

LEMMA 6.4. Let $\mu := (\mu_E)_{E \in \mathcal{E}} \in \mathcal{V}(\mathbf{G})$. Then there exist $\{\tilde{\alpha}_F \in \mathbb{Z} : F \in \mathcal{F} \setminus \mathcal{F}^c\}$ such that

$$(6.4) \quad \mu := \sum_{F \in \mathcal{F} \setminus \mathcal{F}^c} \tilde{\alpha}_F \mu^F.$$

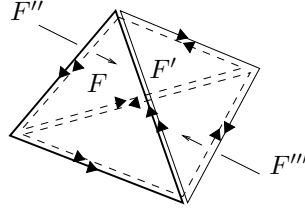
Proof. Let $F \in \mathcal{F}^c$ and let $\boldsymbol{\mu}^F$ be the vector associated with the elementary cycle μ^F formed by the edges E of F in the graph $\mathbf{G} := (\mathcal{N}, \mathcal{E})$. We will first show that there exist $\{\alpha_{F'}^F \in \mathbb{Z} : F' \in \mathcal{F} \setminus \mathcal{F}^c\}$ such that

$$(6.5) \quad \boldsymbol{\mu}^F = \sum_{F' \in \mathcal{F} \setminus \mathcal{F}^c} \alpha_{F'}^F \boldsymbol{\mu}^{F'}.$$

Let $j \in \{1, \dots, n_{\mathcal{W}}\}$ be such that $F = F_j$ in Algorithm 6.1 and let F' , F'' , and F''' be the other faces of T_j ; then there exist $\alpha_{F'}^F, \alpha_{F''}^F, \alpha_{F'''}^F \in \{-1, 1\}$ such that

$$\boldsymbol{\mu}^F = \alpha_{F'}^F \boldsymbol{\mu}^{F'} + \alpha_{F''}^F \boldsymbol{\mu}^{F''} + \alpha_{F'''}^F \boldsymbol{\mu}^{F'''}$$

as depicted below:



If $F', F'', F''' \in \mathcal{F} \setminus \mathcal{F}^c$, the proof of (6.5) is complete. Otherwise assume, without loss of generality, that $F' \in \mathcal{F}^c$. By construction there has to be a $j' \in \{j+1, \dots, n_{\mathcal{W}}\}$ such that $F' = F_{j'}$. Let $\tilde{F}', \tilde{F}'',$ and \tilde{F}''' be the other faces of $T_{j'}$. As before, there exist $\alpha_{\tilde{F}'}^{F'}, \alpha_{\tilde{F}''}^{F'}, \alpha_{\tilde{F}'''}^{F'} \in \{-1, 1\}$ such that $\boldsymbol{\mu}^{F'} = \alpha_{\tilde{F}'}^{F'} \boldsymbol{\mu}^{\tilde{F}'} + \alpha_{\tilde{F}''}^{F'} \boldsymbol{\mu}^{\tilde{F}''} + \alpha_{\tilde{F}'''}^{F'} \boldsymbol{\mu}^{\tilde{F}'''}$ and therefore

$$\boldsymbol{\mu}^F = \alpha_{F'}^F (\alpha_{\tilde{F}'}^{F'} \boldsymbol{\mu}^{\tilde{F}'} + \alpha_{\tilde{F}''}^{F'} \boldsymbol{\mu}^{\tilde{F}''} + \alpha_{\tilde{F}'''}^{F'} \boldsymbol{\mu}^{\tilde{F}'''}) + \alpha_{F''}^F \boldsymbol{\mu}^{F''} + \alpha_{F'''}^F \boldsymbol{\mu}^{F'''}$$

If $F'', F''', \tilde{F}', \tilde{F}'', \tilde{F}''' \in \mathcal{F} \setminus \mathcal{F}^c$, the proof of (6.5) is complete. Otherwise, we can repeat the above procedure for the faces $F'', F''', \tilde{F}', \tilde{F}'',$ and \tilde{F}''' , and since the set $\{1, \dots, n_{\mathcal{W}}\}$ is finite, the procedure will terminate in a finite number of steps. Altogether, we have established that there exist $\{\alpha_{F'}^F \in \mathbb{Z} : F' \in \mathcal{F} \setminus \mathcal{F}^c\}$ such that (6.5) holds.

Now let $\boldsymbol{\mu} \in \mathcal{V}(\mathbf{G})$. Substituting (6.5) into (4.10), we find that there exist $\{\tilde{\alpha}_F \in \mathbb{Z} : F \in \mathcal{F} \setminus \mathcal{F}^c\}$ such that $\boldsymbol{\mu} = \sum_{F \in \mathcal{F} \setminus \mathcal{F}^c} \tilde{\alpha}_F \boldsymbol{\mu}^F$. \square

LEMMA 6.5. Let $\boldsymbol{\mu} \in \mathcal{V}(\mathbf{G})$ and $\{\tilde{\alpha}_F \in \mathbb{Z} : F \in \mathcal{F} \setminus \mathcal{F}^c\}$ be such that $\boldsymbol{\mu} := \sum_{F \in \mathcal{F} \setminus \mathcal{F}^c} \tilde{\alpha}_F \boldsymbol{\mu}^F$. Then

$$(6.6) \quad \sum_{F \in \mathcal{F} \setminus \mathcal{F}^c} \tilde{\alpha}_F \int_F \vec{v}_{F'} \cdot \vec{v}_F dF = 0 \quad \text{for all } F' \in \mathcal{F}^c.$$

Proof. Let $F' \in \mathcal{F}^c$; then $\vec{v}_{F'} \cdot \vec{v}_F = 0$ for all $F \in \mathcal{F} \setminus \mathcal{F}^c$, which implies (6.6). \square

We can now prove Theorem 6.3.

Proof. Since $\#\mathcal{F}^c = n_{\mathcal{W}}$, we merely need to show that the union of the sets of functions (6.3) and (4.8) is a linearly independent set. Therefore suppose $\{\beta_{E'} : E' \in \mathcal{E} \setminus \mathcal{H}\}$ and $\{\gamma_F : F \in \mathcal{F}^c\}$ are scalars such that

$$\vec{0} = \sum_{E' \in \mathcal{E} \setminus \mathcal{H}} \beta_{E'} \vec{\Psi}_{E'} + \sum_{F \in \mathcal{F}^c} \gamma_F \vec{v}_F.$$

Let $E \in \mathcal{E} \setminus \mathcal{H}$ and let $\boldsymbol{\mu}^E$ denote the vector associated with the unique cycle μ^E generated by taking edge E into the tree $\mathbf{H} = (\mathcal{N}, \mathcal{H})$, which has the property that $\mu_{E'}^E := \delta_{E, E'}$ for all $E' \in \mathcal{E} \setminus \mathcal{H}$ (cf. Theorem A.6 and the proof to Theorem 4.2). Now, using Lemma 6.4 we can find $\{\tilde{\alpha}_F \in \mathbb{Z} : F \in \mathcal{F} \setminus \mathcal{F}^c\}$ such that $\boldsymbol{\mu} := \sum_{F \in \mathcal{F} \setminus \mathcal{F}^c} \tilde{\alpha}_F \boldsymbol{\mu}^F$, and so by Lemmas 4.4 and 6.5

$$\begin{aligned} 0 &= \sum_{F \in \mathcal{F} \setminus \mathcal{F}^c} \tilde{\alpha}_F \int_F \left(\sum_{E \in \mathcal{E} \setminus \mathcal{H}} \beta_E \tilde{\Psi}_E + \sum_{F' \in \mathcal{F}^c} \gamma_{F'} \vec{v}_{F'} \right) \cdot \vec{v}_F dF \\ &= \sum_{E' \in \mathcal{E} \setminus \mathcal{H}} \beta_{E'} \sum_{F \in \mathcal{F} \setminus \mathcal{F}^c} \tilde{\alpha}_F \int_F \tilde{\Psi}_E \cdot \vec{v}_F dF + \sum_{F' \in \mathcal{F}^c} \gamma_{F'} \sum_{F \in \mathcal{F} \setminus \mathcal{F}^c} \tilde{\alpha}_F \int_F \vec{v}_{F'} \cdot \vec{v}_F dF \\ &= \sum_{E' \in \mathcal{E} \setminus \mathcal{H}} \beta_{E'} \mu_{E'}^E = \beta_E. \end{aligned}$$

Since the functions $\{\vec{v}_F : F \in \mathcal{F}^c\}$ form a subset of the Raviart–Thomas–Nédélec basis functions, they have to be linearly independent. Therefore we also have $\gamma_F = 0$, for all $F \in \mathcal{F}^c$, which establishes the linear independence of the functions in (6.3) and (4.8). \square

COROLLARY 6.6. *Let $n_{\mathcal{C}} \neq 0$ and let Γ_N^ℓ be simply connected for each $\ell = 1, \dots, n_{\mathcal{C}}$. The functions*

$$(6.7) \quad \{\vec{v}_F : F \in \mathcal{F}^c\}$$

form a complementary basis to (4.15) in \mathcal{V} .

Proof. Since $\#\mathcal{F}^c = n_{\mathcal{W}}$, the result follows directly from Corollary 4.5 and Theorem 6.3. \square

Using this complementary basis (6.7) and applying the general theory presented in section 3, we can therefore find the unique solution \mathbf{p} from (3.6) by simple *back substitutions*.

The matrix $Z'B$ is obtained from the original matrix B in (2.4) by deleting some rows and reordering the rows and columns. Equivalently, the right-hand side $Z'(\mathbf{g} - M\mathbf{u})$ of (3.6) is obtained from $\mathbf{g} - M\mathbf{u}$ by deleting some rows and reordering the rows.

7. Numerical results. In this section we want to demonstrate the performance of the proposed method in two very simple test cases. Let Ω be the unit cube $(0, 1)^3$. We will consider only the constant coefficient case $K \equiv 1$, of system (1.1), (1.2), with zero right-hand side $\vec{g} = \vec{0}$. The two experiments are induced by different choices of boundary conditions. In Figure 2 we illustrate the Neumann boundary Γ_N in each case.

In the *first experiment* (Figure 2 (left)) we choose

$$\begin{aligned} p_D(x, y, z) &= 1 - x && \text{on } \Gamma_D = \{0, 1\} \times (0, 1) \times (0, 1) \cup [0, 1] \times (0, 1) \times \{1\}, \\ \vec{u} \cdot \vec{\nu} &= 0 && \text{on } \Gamma_N = \Gamma \setminus \Gamma_D, \end{aligned}$$

where $\vec{\nu}(\vec{x})$ denotes the outward unit normal from Ω at $\vec{x} \in \Gamma$ as before.

In the *second experiment* (Figure 2 (right)) we choose

$$\begin{aligned} p_D(x, y, z) &= 1 - x && \text{on } \Gamma_D = (0, 1) \times (0, 1) \times \{1\}, \\ \vec{u} \cdot \vec{\nu} &= 0 && \text{on } \Gamma_N = \Gamma \setminus \Gamma_D. \end{aligned}$$

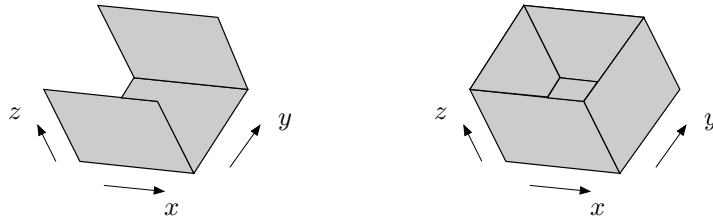


FIG. 2. The “no-flux” boundary Γ_N for Experiments 1 and 2, respectively.

We discretize these problems using the mixed finite element discretization (2.3) with lowest-order Raviart–Thomas–Nédélec elements on a sequence of uniform and nonuniform meshes of different refinement levels L . The uniform mesh is constructed from a uniform rectangular mesh of $L * L * L$ cubes which are each subdivided themselves into six tetrahedra. The nonuniform mesh is constructed from a nonuniform hexahedral mesh of $L * L * L$ hexahedra which are each subdivided themselves into 24 tetrahedra.

To solve the resulting saddle-point system (2.4) we use the decoupled iterative method described above: The construction of the matrix \tilde{A} in the decoupled velocity system (3.3) is carried out in an elementwise fashion as presented in (5.3); the resulting symmetric positive definite system (3.3) is solved with ILU(0)-preconditioned conjugate gradients (PCG); the matrix $Z'B$ in the decoupled pressure system (3.6) is obtained from the original matrix B in (2.4) by deleting some rows and reordering the rows and columns (as mentioned at the end of section 6); the resulting triangular system (3.6) is solved by simple back substitutions. The convergence criterion in the PCG method is the relative reduction of the residual by a factor of 10^{-5} .

Tables 1 and 2 show the performance of our method for Experiments 1 and 2, respectively. First of all we observe in both cases the reduction in size from the full mixed system (2.4) to the decoupled velocity system (3.3). (Compare rows 3 and 4 in Tables 1 and 2.) It is approximately 3 as claimed in Remark 5.3.

Second, let us look at the number of *floating point operations* (Flops) needed in the decoupling process (row 7). This process is asymptotically optimal ($\# \text{ Flops} = O(\# \text{ Freedoms})$), since the matrix \tilde{A} and the right-hand side in (3.3) are constructed in an elementwise fashion and since the matrix $Z'B$ and the right-hand side in (3.6) are obtained from the original system (2.4) by simple reordering. While this decoupling process is still the dominant part in the total Flop count (row 9) on very small problems (e.g., column 2), its effect gets less important and will eventually vanish for larger problems (e.g., columns 5 and 8).

The core part of the calculation is the solution of the decoupled velocity system (3.3). In the uniform mesh case (columns 2–5) the condition number of \tilde{A} (row 5) grows as expected with $O(\tilde{n}^{2/3})$, where \tilde{n} is the dimension of \tilde{A} (row 4). The growth in the nonuniform mesh case (columns 6–8) is (naturally) slightly worse. This growth of the condition number is reflected in the iteration and Flop count for the PCG method to solve (3.3). The number of iterations (row 6) grows with the square root of the condition number of \tilde{A} ; in the uniform mesh case (columns 2–5) this is a factor of about 2 from one refinement level to the next.

REMARK 7.1. *Although the effect of the ILU(0) preconditioner deteriorates as the grid size decreases, it is extremely cheap to invert and remains a cost effective way of preconditioning this system. Diagonal scaling, for example, leads to a similar*

TABLE 1
Performance of the decoupled iterative method for Experiment 1.

	Uniform mesh				Nonuniform mesh		
Refinement level L	2	4	8	16	2	4	8
Freedoms (full mixed)	144	1152	9216	73728	576	4608	36864
Freedoms (decoupled)	48	384	3072	24576	192	1536	12288
Condition # (decoupled)	150	610	2370	9230	1240	10670	89700
PCG-iterations	14	26	45	97	33	90	232
MFlops (decoupling)	0.032	0.34	3.0	25	0.17	1.5	13
MFlops (PCG)	0.031	0.53	8.0	142	0.35	8.0	171
MFlops (total)	0.065	0.89	11.2	169	0.53	9.6	185

TABLE 2
Performance of the decoupled iterative method for Experiment 2.

	Uniform mesh				Nonuniform mesh		
Refinement level L	2	4	8	16	2	4	8
Freedoms (full mixed)	128	1088	8960	72704	544	4480	36352
Freedoms (decoupled)	32	320	2816	23552	160	1408	11776
Condition # (decoupled)	87	450	2080	8740	360	2780	21300
PCG-iterations	9	18	35	75	24	80	187
MFlops (decoupling)	0.024	0.29	2.8	25	0.15	1.4	12
MFlops (PCG)	0.012	0.31	5.8	106	0.20	6.5	132
MFlops (total)	0.037	0.61	8.7	132	0.36	8.0	145

asymptotic behavior but is more expensive in terms of iterations as well as Flops. For Experiment 1 the number of iterations using diagonal scaling as the preconditioner in the PCG method in the uniform mesh case is 38, 109, 245, and 494 for $L = 2, 4, 8$, and 16, respectively. The number of MFlops is 0.086, 1.73, 29.7, and 47. Future work will involve a more detailed investigation of various preconditioners like algebraic multigrid (AMG).

Finally, let us look at the solution of the decoupled pressure system (3.6). Since the matrix $Z'B$ in (3.6) is triangular (as shown in Theorem 6.2) it can be solved in an asymptotically optimal number of operations by simple back substitutions. Therefore this part of the calculation does not affect the overall cost of the method significantly. The number of Flops is included in the total Flop count (row 9) reported in Tables 1 and 2, and it accounts for less than 1% of the total cost of the method for larger problems (e.g., columns 5 and 8).

In Tables 3 and 4 we compare the performance of our method with the performance of a preconditioned MINRES method for the original (full mixed) saddle-point system (2.4) for Experiments 1 and 2, respectively. The convergence criterion for MINRES is again the relative reduction of the residual by a factor of 10^{-5} . To precondition this MINRES method we take an optimal symmetric positive definite block diagonal preconditioner presented and analyzed in Rusten and Winther [26] (using an ILU(0) factorization of $B^T B$ for the pressure block).

REMARK 7.2. *As for the decoupled system it would be possible to employ other, more efficient preconditioners like AMG for the pressure block. However, in order to evaluate the performance of the decoupling procedure, our prime incentive, we wanted to compare “similar” iterative methods for the full mixed and the decoupled system.*

Comparing columns 6 and 7 in Table 3 we observe that for the first experiment our decoupled method is about 4.5 times faster than preconditioned MINRES on the uniform meshes (rows 3–6), and about 6 times faster on the nonuniform meshes

TABLE 3

Comparison of the decoupled iterative method with a full mixed method [26] for Experiment 1.

		Freedoms		Iterations		MFlops	
Mesh	L	Mixed	Decoupled	Mixed	Decoupled	Mixed	Decoupled
Uniform	2	144	48	37	14	0.29	0.065
	4	1152	384	56	26	3.8	0.89
	8	9216	3072	89	45	49	11
	16	73728	24576	175	97	790	169
Nonuniform	2	576	192	100	33	3.4	0.53
	4	4608	1536	204	90	57	9.6
	8	36864	12288	456	232	1030	185

TABLE 4

Comparison of the decoupled iterative method with a full mixed method [26] for Experiment 2.

		# Freedoms		Iterations		MFlops	
Mesh	L	Mixed	Decoupled	Mixed	Decoupled	Mixed	Decoupled
Uniform	2	128	32	37	9	0.25	0.038
	4	1088	320	57	18	3.6	0.62
	8	8960	2816	109	35	59	8.7
	16	72704	23552	217	75	960	132
Nonuniform	2	544	160	110	24	3.5	0.35
	4	4480	1408	254	80	68	8.0
	8	36352	11776	582	187	1300	145

(rows 7–9). The advantage of our method over preconditioned MINRES is even more impressive for the second experiment (columns 6 and 7 in Table 4). On the uniform meshes (rows 3–6) it is about 7 times faster, on the nonuniform meshes (rows 7–9) about 9 times faster.

In conclusion we have found a very competitive and practicable iterative method to solve saddle-point problems of the form (2.4). The decoupling procedure and the recovery of the pressure are asymptotically optimal. The decoupled velocity system, on the other hand, is symmetric positive definite and of second order, and there may still be room for improvement in solving this system by employing the more sophisticated preconditioning techniques which are available for such problems.

Appendix A. Some results from graph theory.

DEFINITION A.1 (Berge [5]).

- (a) A graph (or more precisely a 1-graph) \mathbf{G} is defined to be a pair $(\mathcal{X}, \mathcal{U})$, where \mathcal{X} is a set $\{x_1, x_2, \dots, x_n\}$ of elements called vertices (or nodes), and \mathcal{U} is a subset $\{u_1, u_2, \dots, u_m\}$ of $\mathcal{X} \times \mathcal{X}$ of elements called arcs (or orientated edges). For an arc $u = (x, y) \in \mathcal{X} \times \mathcal{X}$, the vertex x is called its initial endpoint, and the vertex y is called its terminal endpoint. A vertex $y \in \mathcal{X}$ is called a neighbor of $x \in \mathcal{X}$ if either $(x, y) \in \mathcal{U}$ or $(y, x) \in \mathcal{U}$. The set of all neighbors of a vertex x in the graph \mathbf{G} will be denoted by $\Gamma_{\mathbf{G}}(x)$.
- (b) A partial graph of a graph $\mathbf{G} = (\mathcal{X}, \mathcal{U})$ is a graph $\mathbf{H} = (\mathcal{X}, \mathcal{V})$ with $\mathcal{V} \subset \mathcal{U}$.
- (c) A chain is a sequence $\mu = (u_{i_1}, u_{i_2}, \dots, u_{i_q})$ of arcs of a graph \mathbf{G} such that each arc in the sequence has one endpoint in common with its predecessor and its other endpoint in common with its successor. A chain that does not encounter the same vertex twice is called elementary. A chain that does not use the same arc twice is called simple.

- (d) For two vertices x and y of a graph \mathbf{G} let us define the equivalence relation $x \equiv y$ by

$$[x = y, \text{ or } x \neq y \text{ and there exists a chain in } \mathbf{G} \text{ connecting } x \text{ and } y].$$

The equivalence classes of \equiv are called the connected components of \mathbf{G} . A connected graph is a graph that consists only of one connected component.

- (e) A cycle is a simple chain whose terminal endpoint coincides with its initial endpoint. Let m be the number of arcs in \mathbf{G} . With each cycle μ of \mathbf{G} we can associate a vector $\boldsymbol{\mu} \in \mathbb{Z}^m$ with

$$\mu_i = \begin{cases} 0 & \text{if } u_i \text{ is not in } \mu, \\ 1 & \text{if } u_i \text{ is in } \mu \text{ and shares initial endpoint with its predecessor,} \\ -1 & \text{if } u_i \text{ is in } \mu \text{ and shares terminal endpoint with its predecessor.} \end{cases}$$

The set of all those vectors $\boldsymbol{\mu} \in \mathbb{Z}^m$ generates a vector space over \mathbb{Z} . We denote this vector space by $\mathcal{V}(\mathbf{G})$.

- (f) A forest is defined to be a graph without cycles. A tree is defined to be a connected graph without cycles.

THEOREM A.2. Let \mathbf{G} be a graph with n vertices, m arcs, and p connected components. The dimension of $\mathcal{V}(\mathbf{G})$ is $m - n + p$.

Proof. The proof is seen in Berge [5, p. 16]. \square

THEOREM A.3. Let $\mathbf{H} = (\mathcal{X}, \mathcal{U})$ be a graph with $n > 2$ vertices. The following properties are equivalent and each characterizes a tree:

- (i) \mathbf{H} is connected and has no cycles.
- (ii) \mathbf{H} has $n - 1$ arcs and has no cycles.
- (iii) \mathbf{H} is connected and contains $n - 1$ arcs.
- (iv) \mathbf{H} has no cycles and adding an arc creates a unique cycle.
- (v) \mathbf{H} is connected and removing an arc leaves the remaining graph disconnected.
- (vi) Every pair of vertices x, y of \mathbf{H} is connected by a unique chain.

Proof. The proof is seen in Berge [5, p. 24]. \square

THEOREM A.4. Let $\mathbf{G} = (\mathcal{X}, \mathcal{U})$ be a connected graph. There exists a partial graph $\mathbf{H} = (\mathcal{X}, \mathcal{V})$ such that \mathbf{H} is a tree.

Proof. The proof is seen in Berge [5, p. 25]. \square

The tree \mathbf{H} obtained from \mathbf{G} as above is called a *spanning tree*. An optimal algorithm to find a spanning tree \mathbf{H} of a connected graph \mathbf{G} is presented in Algorithm A.7.

THEOREM A.5. Let \mathbf{G} be a graph with n vertices and $m \geq n$ arcs. The time spent on Algorithm A.7 is proportional to the number of arcs, i.e., $O(m)$.

Proof. The proof is seen in Aho, Hopcroft, and Ullman [1]. \square

THEOREM A.6. Let $\mathbf{G} = (\mathcal{X}, \mathcal{U})$ be a connected graph with n vertices and m arcs, let $\mathbf{H} = (\mathcal{X}, \mathcal{V})$ be a spanning tree of \mathbf{G} , and let $u_i \in \mathcal{U}$ be an arc of \mathbf{G} not in tree \mathbf{H} , i.e., $u_i \notin \mathcal{V}$. Adding u_i to \mathbf{H} creates a unique cycle μ^i , and its associated vector $\boldsymbol{\mu}^i$ satisfies $\mu_i^i = 1$. The set $\{\boldsymbol{\mu}^i : u_i \in \mathcal{U} \setminus \mathcal{V}\}$ forms a basis of $\mathcal{V}(\mathbf{G})$.

Proof. The existence of $\boldsymbol{\mu}^i$ for all $u_i \in \mathcal{U} \setminus \mathcal{V}$ is guaranteed by virtue of Theorem A.3(iv). The vectors are linearly independent, since $\mu_i^j = \delta_{i,j}$, for all $u_i, u_j \in \mathcal{U} \setminus \mathcal{V}$. Moreover,

$$\dim\{\boldsymbol{\mu}^i : u_i \in \mathcal{U} \setminus \mathcal{V}\} = \#\mathcal{U} - \#\mathcal{V} = m - (n - 1) = \dim \mathcal{V}(\mathbf{G}),$$

where in the last step we used Theorem A.2 with $p = 1$. \square

ALGORITHM A.7.

```

variables
   $n$  – number of vertices;
   $\text{mark}[1:n]$  – array of flags;
begin
   $\mathcal{V} = \emptyset$ ;
  for  $i := 1$  to  $n$  do  $\text{mark}[x_i] := \text{unvisited}$ ;
   $\text{recursive}(x_1)$ ;
end;

procedure  $\text{recursive}(x - \text{vertex})$ ;
  variables
     $y - \text{vertex}$ ;
  begin
     $\text{mark}[x] := \text{visited}$ ;
    for each  $\text{vertex } y \in \Gamma_{\mathbf{G}}(x)$  do
      if  $\text{mark}[y] = \text{unvisited}$  then
         $\mathcal{V} = \mathcal{V} \cup \{u\}$  – where  $u \in \mathcal{U}$  with endpoints  $x$  and  $y$ ;
         $\text{recursive}(y)$ ;
    end;

```

Appendix B. A topological result on simplicial triangulations.

DEFINITION B.1 (the fundamental group (Armstrong [3, Ch. 5])).

- (a) A topological space is a set S together with a collection \mathcal{U} of subsets of S satisfying the following conditions:

- (1) $\emptyset \in \mathcal{U}$, $S \in \mathcal{U}$.
- (2) If $U_1, \dots, U_n \in \mathcal{U}$, then $\bigcap_{i=1}^n U_i \in \mathcal{U}$.
- (3) If $\tilde{\mathcal{U}} \subset \mathcal{U}$, then $\bigcup_{U \in \tilde{\mathcal{U}}} U \in \mathcal{U}$.

The elements of \mathcal{U} are called open sets in S . \mathcal{U} is called a topology on S .

- (b) Let X be a topological space. A path in X from x_0 to x_1 (with origin x_0 and end x_1) is a continuous map $\alpha : [0, 1] \rightarrow X$ such that $\alpha(0) = x_0$ and $\alpha(1) = x_1$. Let α be a path in X from x_0 to x_1 and let β be a path in X from x_1 to x_2 . The product of α and β is the path $\alpha\beta$ from x_0 to x_2 defined by

$$\alpha\beta(t) = \begin{cases} \alpha(2t) & \text{for } t \in [0, 1/2], \\ \beta(2t - 1) & \text{for } t \in [1/2, 1]. \end{cases}$$

The inverse of α is the path α^{-1} from x_1 to x_0 defined by $\alpha^{-1}(t) = \alpha(1 - t)$.

- (c) Two paths α and β from x_0 to x_1 are homotopic (written $\alpha \simeq \beta$) if there exists a continuous map $F : [0, 1] \times [0, 1] \rightarrow X$ such that

$$\begin{aligned} F(0, t) = x_0 & \quad \text{and} \quad F(1, t) = x_1 & \quad \text{for all } t \in [0, 1], \\ F(s, 0) = \alpha(s) & \quad \text{and} \quad F(s, 1) = \beta(s) & \quad \text{for all } s \in [0, 1]. \end{aligned}$$

- (d) Let X be a topological space and let $x_0 \in X$. The set of \simeq equivalence classes of paths with origin x_0 and end x_0 forms a group under the operations of multiplication and inverse as defined above. This group is denoted $\pi_1(X, x_0)$ and is called the fundamental group of the pair (X, x_0) . X is called simply connected if its fundamental group is trivial.

DEFINITION B.2 (the first homology group (Armstrong [3, Ch. 8])).

- (a) Let V be a vector space over \mathbb{R} and let $\{v_0, v_1, \dots, v_k\} \subset V$ such that the set $\{v_1 - v_0, \dots, v_k - v_0\}$ is linearly independent. The smallest convex set containing $\{v_0, v_1, \dots, v_k\}$, i.e., the convex hull

$$\left\{ v := \sum_{i=0}^k \lambda_i v_i : \lambda_i \geq 0 \text{ and } \sum_{i=0}^k \lambda_i = 1 \right\},$$

is called a simplex of dimension k (or a k -simplex). The points v_0, v_1, \dots, v_k are called the vertices (or nodes) of the simplex. The simplices formed by the subsets of $\{v_0, v_1, \dots, v_k\}$ are called the faces of the simplex.

- (b) A simplicial complex K is a finite set of simplices in V such that

- (1) if $A \in K$, then the faces of A are also in K ;
- (2) if $A, B \in K$ and $A \cap B \neq \emptyset$, then $A \cap B \in K$.

The dimension of K is the maximum dimension of the simplices of K . The point set union of all simplices in K is denoted by $|K|$.

- (c) Let K be a simplicial complex. An orientated edge in K is an ordered pair (u, v) such that u and v lie in some simplex of K . An orientated triangle in K is an ordered triple (u, v, w) such that u, v, w lie in some simplex of K . Note that $(u, v, w) = (v, w, u) = (w, u, v)$. A change of orientation is denoted by a minus sign, thus $(v, u) = -(u, v)$ and $(v, u, w) = -(u, v, w)$. The boundary of the orientated edge (u, v) is defined to be

$$\partial(u, v) = v - u.$$

The boundary of the orientated triangle (u, v, w) is

$$\partial(u, v, w) = (v, w) + (w, u) + (u, v).$$

Let n be the number of all edges in K . A linear combination of orientated edges

$$\sum_{i=1}^n \lambda_i (u_i, v_i) \quad \text{with the property that} \quad \sum_{i=1}^n \lambda_i \partial(u_i, v_i) = 0$$

and $\lambda_i \in \mathbb{Z}$ for all $i = 1, \dots, n$ is called a (one-dimensional) cycle of K . A cycle β is called a bounding cycle if we can find a linear combination

$$\sum_{j=1}^k \alpha_j (u_j, v_j, w_j)$$

of orientated triangles in K such that

$$\beta = \sum_{j=1}^k \alpha_j \partial(u_j, v_j, w_j).$$

- (d) The set of all cycles of K forms an abelian group under the addition

$$\sum_{i=1}^n \lambda_i (u_i, v_i) + \sum_{i=1}^n \mu_i (u_i, v_i) = \sum_{i=1}^n (\lambda_i + \mu_i) (u_i, v_i).$$

We denote this group by $Z_1(K)$. The bounding cycles form a subgroup $B_1(K)$ of $Z_1(K)$. The quotient group

$$H_1(K) = Z_1(K) \setminus B_1(K)$$

is called the first homology group of K .

We will need only the following fundamental theorem which is a corollary to the simplicial approximation theorem (Armstrong [3, p. 128]).

THEOREM B.3. *Let K be a simplicial complex and let v be a vertex of K . If $|K|$ is connected, abelianizing $\pi_1(|K|, v)$ gives the first homology group $H_1(K)$.*

Proof. The proof is seen in Armstrong [3, p. 182]. \square

COROLLARY B.4. *If $|K|$ is simply connected, then each cycle of K is a bounding cycle.*

Proof. From Definition B.1(d) we know that if $|K|$ is simply connected, then $\pi_1(|K|, v)$ is trivial for any vertex v of K . As a consequence of Theorem B.3 this also implies that $H_1(K)$ is trivial (since abelianizing the trivial group has to result in the trivial group again). Therefore $B_1(K) = Z_1(K)$. \square

Acknowledgment. I would like to thank Professor Thomas Russell for kindly sending me a copy of Hecht's unpublished manuscript [18] and particularly thank Professor Ivan Graham and Professor Andrew Swann for useful discussions and comments.

REFERENCES

- [1] A.V. AHO, J.E. HOPCROFT, AND J.D. ULLMAN, *Data Structures and Algorithms*, Addison-Wesley, Reading, MA, 1983.
- [2] R. ALBANESE AND G. RUBINACCI, *Integral formulation for 3D eddy-current computation using edge-elements*, IEE Proceedings A, 135 (1988), pp. 457–462.
- [3] M.A. ARMSTRONG, *Basic Topology*, Springer, New York, 1983.
- [4] D.N. ARNOLD, R.S. FALK, AND R. WINTHER, *Preconditioning in $H(\text{div})$ and applications*, Math. Comp., 66 (1997), pp. 957–984.
- [5] C. BERGE, *Graphs and Hypergraphs*, North-Holland, Amsterdam, 1973.
- [6] A. BOSSAVIT, *Computational Electromagnetism: Variational Formulation, Complementarity, Edge Elements*, Academic Press Electromagnetism Series 2, Academic Press, San Diego, 1998.
- [7] F. BREZZI AND M. FORTIN, *Mixed and Hybrid Finite Element Methods*, Springer, New York, 1991.
- [8] Z. CAI, R.R. PARASHKEVOV, T.F. RUSSELL, AND X. YE, *Domain decomposition for a mixed finite element method in three dimensions*, SIAM J. Numer. Anal., submitted.
- [9] A.-L. CAUCHY, *Recherches sur les polyèdres – premier mémoire*, J. de l'École Polytechnique 9, (1813), pp. 68–86.
- [10] G. CHAVENT, G. COHEN, J. JAFFRE, M. DUPUY, AND I. RIBERA, *Simulation of two-dimensional water flooding by using mixed finite elements*, Soc. Petroleum Eng. J., 24 (1984), pp. 382–390.
- [11] Z. CHEN, R.E. EWING, R.D. LAZAROV, S. MALIASOV, AND Y.A. KUZNETSOV, *Multilevel preconditioners for mixed methods for second order elliptic problems*, Numer. Linear Algebra Appl., 3 (1996), pp. 427–453.
- [12] K.A. CLIFFE, I.G. GRAHAM, R. SCHEICHL, AND L. STALS, *Parallel computation of flow in heterogeneous media modelled by mixed finite elements*, J. Comput. Phys., 164 (2000), pp. 258–282.
- [13] L.C. COWSAR, J. MANDEL, AND M.F. WHEELER, *Balancing domain decomposition for mixed finite elements*, Math. Comp., 211 (1995), pp. 989–1015.
- [14] F. DUBOIS, *Discrete vector potential representation of a divergence-free vector field in three-dimensional domains: Numerical analysis of a model problem*, SIAM J. Numer. Anal., 27 (1990), pp. 1103–1141.
- [15] R.E. EWING AND J. WANG, *Analysis of the Schwarz algorithm for mixed finite element methods*, RAIRO Modél. Math. Anal. Numér., 26 (1992), pp. 739–756.

- [16] R.E. EWING AND J. WANG, *Analysis of multilevel decomposition iterative methods for mixed finite element methods*, RAIRO Modél. Math. Anal. Numér., 28 (1994), pp. 377–398.
- [17] F. HECHT, *Construction d’une base de fonctions P_1 non-conformes à divergence nulle dans \mathbb{R}^3* , RAIRO Anal. Numér., 15 (1981), pp. 119–150.
- [18] F. HECHT, *Construction d’une base pour des éléments finis mixtes à divergence faiblement nulle*, manuscript, 1988.
- [19] R. HIPTMAIR AND R.H.W. HOPPE, *Multilevel methods for mixed finite elements in three dimensions*, Numer. Math., 82 (1999), pp. 253–279.
- [20] L. KETTUNEN AND L.R. TURNER, *How to define the minimum set of equations for edge elements*, Int. J. Appl. Electromagnetics Mater., 3 (1992), pp. 47–53.
- [21] T.P. MATHEW, *Schwarz alternating and iterative refinement methods for mixed formulations of elliptic problems. I.*, Numer. Math., 65 (1993), pp. 445–468.
- [22] T.P. MATHEW, *Schwarz alternating and iterative refinement methods for mixed formulations of elliptic problems. II.*, Numer. Math., 65 (1993), pp. 469–492.
- [23] J.C. NÉDÉLEC, *Mixed finite elements in \mathbb{R}^3* , Numer. Math., 35 (1980), pp. 315–341.
- [24] L.F. PAVARINO AND M. RAMÉ, *Numerical experiments with an overlapping additive Schwarz solver for 3D parallel reservoir simulation*, Int. J. Supercomp. Appl., 1 (1995), pp. 3–17.
- [25] J.E. ROBERTS AND J.M. THOMAS, *Mixed and hybrid methods*, in Handbook of Numerical Analysis, Vol. 2, P.G. Ciarlet and J.L. Lions, eds., North-Holland, Amsterdam, 1991, pp. 523–639.
- [26] T. RUSTEN AND R. WINTER, *A preconditioned iterative method for saddlepoint problems*, SIAM J. Matrix Anal. Appl., 13 (1992), pp. 887–904.
- [27] R. SCHEICHL, *Iterative Solution of Saddle-Point Problems Using Divergence-Free Finite Elements with Applications to Groundwater Flow*, Ph.D. Thesis, University of Bath, Bath, UK, 2000.
- [28] F. THOMASSET, *Implementation of Finite Element Methods for Navier-Stokes Equations*, Springer, New York, 1981.
- [29] B.I. WOHLMUTH, A. TOSELLI, AND O.B. WIDLUND, *An iterative substructuring method for Raviart-Thomas vector fields in three dimensions*, SIAM J. Numer. Anal., 37 (2000), pp. 1657–1676.