

LOCALLY CONSERVATIVE FLUXES FOR THE CONTINUOUS GALERKIN METHOD

BERNARDO COCKBURN ^{*}, JAYADEEP GOPALAKRISHNAN [†], AND HAIYING WANG [‡]

Abstract. The standard continuous Galerkin (CG) finite element method for second order elliptic problems suffers from its inability to provide conservative flux approximations, a much needed quantity in many applications. We show how to overcome this shortcoming by using a two step postprocessing. The first step is the computation of a numerical flux trace defined on element interfaces and is motivated by the structure of the numerical traces of discontinuous Galerkin methods. This computation is non-local in that it requires the solution of a symmetric positive definite system, but the system is well conditioned independently of mesh size, so it can be solved at asymptotically optimal cost. The second step is a local element by element postprocessing of the CG solution incorporating the result of the first step. This leads to a conservative flux approximation with continuous normal components. This postprocessing applies for the CG method in its standard form or for a hybridized version of it. We present the hybridized version since it allows easy handling of variable-degree polynomials and hanging nodes. Furthermore, we provide an *a priori* analysis of the error in the postprocessed flux approximation and display numerical evidence suggesting that the approximation is competitive with the approximation provided by the Raviart-Thomas mixed method of corresponding degree.

Key words. continuous Galerkin methods, conforming finite element method, hybridization, elliptic problems, conservation

AMS subject classifications. 65M60,65N30,35L65.

1. Introduction. In this paper, we revisit the classical finite element method [13, 20], otherwise known as the continuous Galerkin (CG) method, for second-order elliptic problems, with the intention of showing how to overcome what is perhaps its main disadvantage, namely, the discontinuity of the normal component of the approximate flux across element interfaces. We show how to achieve this by means of an efficient postprocessing of the approximate solution provided by the CG method. We also show that the postprocessed flux is competitive with the flux provided by the Raviart-Thomas mixed method of corresponding degree.

We illustrate our technique in the framework of the model second order elliptic boundary value problem

$$-\nabla \cdot (a\nabla u) = f \quad \text{on } \Omega, \tag{1.1a}$$

$$u = g \quad \text{on } \Gamma_D, \tag{1.1b}$$

$$-a\nabla u \cdot \mathbf{n} = \mathbf{q}_N \quad \text{on } \Gamma_N. \tag{1.1c}$$

Here $\Omega \subset \mathbb{R}^N$ is a polyhedral domain ($N \geq 2$) with boundary $\partial\Omega$, $f \in L^2(\Omega)$, and $a = a(\mathbf{x})$ is a symmetric $N \times N$ matrix function that is uniformly positive definite on Ω with components in $L^\infty(\Omega)$. The boundary conditions are given by functions g and \mathbf{q}_N on disjoint subsets Γ_D and Γ_N of $\partial\Omega$, upon which further assumptions will

^{*}School of Mathematics, University of Minnesota, Vincent Hall, Minneapolis, MN 55455, USA, email: cockburn@math.umn.edu. Supported in part by the National Science Foundation (Grant DMS-0411254) and by the University of Minnesota Supercomputing Institute.

[†]Department of Mathematics, University of Florida, Gainesville, FL 32611-8105, USA, email: jayg@math.ufl.edu. Supported in part by the National Science Foundation under grants DMS-0410030 and SCREMS-0619080.

[‡]School of Mathematics, University of Minnesota, Vincent Hall, Minneapolis, MN 55455, USA, email: hywang@math.umn.edu.

be placed shortly. Here and elsewhere we use \mathbf{n} to denote the unit outward normal on the boundary of some domain – the domain will be clear from the context, e.g., in (1.1c) it is Ω . As is well known, this boundary value problem models a wide range of problems of practical interest from electromagnetics to heat dissipation and flow in porous media.

To facilitate the discussion of the results, let us introduce our notations for the CG method right away. Let \mathcal{T}_h denote a triangulation of the domain Ω , which for simplicity we assume to consist of simplices. Define the space

$$\mathbf{V}_h = \{v \in \mathcal{C}^0(\Omega) : v|_K \in \mathcal{P}_k(K) \text{ for } K \in \mathcal{T}_h\}, \quad (1.2)$$

where $\mathcal{C}^0(D)$ denotes the space of continuous functions on a domain D . We assume that Γ_D is the union of some mesh faces (edges if $N = 2$) lying on $\partial\Omega$ and that $\Gamma_N = \partial\Omega \setminus \Gamma_D$. We assume that g is in the space of traces on Γ_D of functions in \mathbf{V}_h and set $\mathbf{V}_h(g) = \{v \in \mathbf{V}_h : v = g \text{ on } \Gamma_D\}$. If a Dirichlet data that is not polynomial is given, one can proceed by approximating it as usual, but we shall not consider this case. As is well known, the approximate solution u_h of the CG method is the function in $\mathbf{V}_h(g)$ determined by

$$(a\nabla u_h, \nabla v)_\Omega = (f, v)_\Omega - \langle \mathbf{q}_N, v \rangle_{\Gamma_N} \quad \text{for all } v \in \mathbf{V}_h(0). \quad (1.3)$$

Here we have used common notations for innerproducts: For scalar functions w and v on some domain $\mathcal{D} \subset \mathbb{R}^N$, $(w, v)_\mathcal{D} = \int_\mathcal{D} w v \, dx$, for vector functions $(\mathbf{p}, \mathbf{q})_\mathcal{D} = \int_\mathcal{D} \mathbf{p} \cdot \mathbf{q} \, dx$, and for functions on domains B formed by lower dimensional objects like union of a few mesh faces, $\langle \eta, \zeta \rangle_B = \int_B \eta \zeta \, d\gamma$.

It is well known that the CG approximation given by $-a\nabla u_h$ to the flux $\mathbf{q} = -a\nabla u$ is not conservative. The root of the problem is evident once we write (1.1a) in conservation form as $\text{div } \mathbf{q} = f$. While the flux approximations from mixed and DG methods satisfy a discrete analogue of this equation, the CG flux $-a\nabla u_h$ does not. We say that a discrete flux \mathbf{q}_h approximating the exact flux \mathbf{q} is *conservative* if the total outward flux across any “discrete subdomain” as measured by \mathbf{q} and \mathbf{q}_h coincides, or more precisely,

$$\int_{\partial D_h} \mathbf{q} \cdot \mathbf{n} \, ds = \int_{\partial D_h} \mathbf{q}_h \cdot \mathbf{n} \, ds, \quad (1.4)$$

for any domain D_h formed by the union of some mesh elements in \mathcal{T}_h (where \mathbf{n} is unit outward normal on the boundary of D_h). Conservative flux approximations are very important in many applications, e.g., in oil recovery simulations, more generally in flows through porous media, and indeed in computational fluid dynamics in general. The same is true in computational structural mechanics, where mixed and hybrid methods were devised to cope with its absence in the so-called one-field displacement method for linear elasticity (which is the CG method for elasticity) – see, e.g. the first paragraph of §3.3 in [31].

Many researchers have attempted to overcome the lack of conservativity of the CG flux by generating a better flux through postprocessing. However, a conservative $\mathbf{H}(\text{div}, \Omega)$ -conforming flux approximation has eluded their efforts for more than three decades. Let us briefly review what has been achieved up to date. In [33], J. Wheeler showed how to postprocess the CG solution to obtain approximations to the normal component of \mathbf{q} at the boundary of the computational domain. In one space dimension, this procedure can be extended to compute approximations to \mathbf{q}

at all the nodes. In fact, such approximations were proven by M. Wheeler in [34] to superconverge with order $2k$ when using polynomial approximations of degree k . This solves the problem in the one dimensional case. In the multidimensional case however, the situation is rather different and no $\mathbf{H}(\text{div}, \Omega)$ -conforming approximation of \mathbf{q} has been constructed so far. Moreover, there are only a few theoretical and numerical studies of the approximation given by J. Wheeler's procedure. In [23], it was shown that such a procedure provides an approximation that superconverges in the $L^2(\partial\Omega)$ -norm with order $k + 1$, for $a \equiv 1$, and with order $k + 1/2$, when a is smooth (under the assumption that Ω is a square endowed with a Cartesian mesh). In [3], the integral of the normal component of the flux on the whole boundary was proven to superconverge with order $k + 1$ when Ω is a curved domain and isoparametric elements are used, and with order $2k$ when it is a polyhedron. For numerical studies, see the references cited in [25]. More importantly, in [25] the CG method was argued to have the property of local conservativity; see also [26] for an extension of this approach to the advection-diffusion and incompressible Navier-Stokes equations. In [9], the so-called *superconvergent integral flux post-processing formula* was revisited. The conservation property was proven and a relation to a Lagrange multiplier mixed formulation and the associated consistency implications were established. See also [10] for further work on conservative projections involving multipliers in a different context. However, none of the approaches used in [25, 9] can be employed to construct an $\mathbf{H}(\text{div}, \Omega)$ -conforming approximation of the flux rendering the CG method locally conservative. (The precise relation between this approach and ours is displayed right before subsection 3.2.) In [7], this approach was used (for $a = 1$ and $N = 2$) to obtain an approximation of the integral of the normal component of the flux on an internal boundary which splits the domain in two; an order of convergence of $2k$ was proven for such approximation.

In this paper, we show how to obtain a conservative flux approximation \mathbf{q}_h in $\mathbf{H}(\text{div}, \Omega)$ that renders the CG method locally conservative. This is done by post-processing the CG solution u_h in two steps. The objective of the first is to compute a numerical trace $\hat{\mathbf{q}}_h$ of the flux whose normal component is single-valued on the interelement boundaries and *renders locally conservative the CG method*, that is, it satisfies

$$-\sum_{K \in \mathcal{T}_h} (a \nabla u_h, \nabla v)_K + \sum_{K \in \mathcal{T}_h} \langle \hat{\mathbf{q}}_h \cdot \mathbf{n}, v \rangle_{\partial K} = (f, v)_\Omega$$

for all v such that $v|_K \in \mathcal{P}_\ell(K)$ for all $K \in \mathcal{T}_h$. The form of this numerical trace is similar to that of the corresponding numerical traces of the DG methods. However, unlike the DG numerical traces, the crucial *stabilization* term cannot have the form of a parameter times the jump of the u_h , since in our case such jump is identically equal to zero. Instead, it is a quantity that belongs to a certain non-standard space of *jumps* and that depends globally on the CG approximation u_h . While the need for this term is far from obvious when approaching from the standard CG formulation, it becomes clearer from the hybridized form of the CG method, which uses a space of discontinuous functions that generate the above mentioned space of jumps on mesh faces. Because of this, we now face difficulties not encountered in DG methods: the computation of $\hat{\mathbf{q}}_h$ requires (i) a local basis representation of the space of jumps, and (ii) the solution of a global system in that space. We are able to overcome the former difficulty by extending some techniques developed in [17, 18]. Although the latter difficulty persists, it turns out that the stiffness matrix of the global system is

symmetric, positive definite, and well conditioned. In particular, we prove that its condition number is bounded independently of meshsize, so it can be solved iteratively at asymptotically optimal cost. In [27] a similar but different way of computing a numerical trace has been proposed; see the discussion before subsection 3.3.

The second step in the postprocessing is the local element by element recovery of a conservative flux approximation \mathbf{q}_h throughout the computational domain by a variation of the so-called Raviart-Thomas projection [29]. Similar techniques have been used by [4] in the framework of DG methods for Darcy's law and by [19] in the context of DG methods for the Navier-Stokes equations. The flux approximation \mathbf{q}_h coincides with the numerical trace $\widehat{\mathbf{q}}_h$ on element boundaries supplied by the previous step and is lifted to the interior of each element by using the $a \nabla u_h$ in such a way that

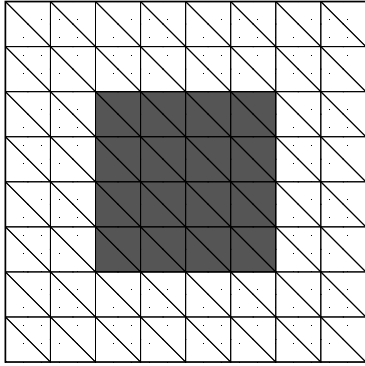
$$-\sum_{K \in \mathcal{T}_h} (a \nabla u_h, \nabla v)_K + \sum_{K \in \mathcal{T}_h} \langle \widehat{\mathbf{q}}_h \cdot \mathbf{n}, v \rangle_{\partial K} = \sum_{K \in \mathcal{T}_h} (\nabla \cdot \mathbf{q}_h, v)_K = (f, v)_\Omega$$

for all v such that $v|_K \in \mathcal{P}_\ell(K)$ for all $K \in \mathcal{T}_h$. We prove that the resulting approximation \mathbf{q}_h converges to the exact flux \mathbf{q} at the *same* order of convergence as the approximation provided by the RT mixed method of corresponding order. Moreover, since the computation of the CG solution requires solving a system that is smaller in size than the corresponding RT system, our flux computation becomes a competitive alternative.

In [32], a technique is proposed for computing a locally conservative flux approximation in the domain Ω from its exact divergence in Ω and an approximation of its normal component on the inter-element boundaries. It also proceeds in two steps. In the first, a locally conservative approximation to the normal component is obtained by solving a global constrained minimization problem. Then, on each element, the data on the border is lifted to the interior to obtain the desired flux; a local mixed element method is used to achieve this. The application of this technique to the CG method differs from ours in several respects. First of all, the resulting numerical trace does not render locally conservative the CG method, in the sense defined above. Moreover, to obtain it, a global constrained minimization problem is to be solved; this has to be contrasted with our unconstrained minimization problem whose stiffness matrix has a condition number bounded independently of the meshsize. Finally, to obtain what we call \mathbf{q}_h , the approximation u_h given by the CG method is not used.

Let us compare our flux \mathbf{q}_h with the RT flux obtained for the model problem (1.1) with $f = 0$, $\Omega = (0, 1) \times (0, 1)$, and boundary conditions as indicated in Fig. 1.1. Here $a = 0.001Id$ in the region $(.25, .75) \times (.25, .75)$ and $a = Id$ elsewhere (Id denotes the identity matrix) – see Fig. 1.1. We can think of this problem as modeling the steady state flow of a fluid through a porous medium with permeability given by a . In Fig. 1.2 we display the streamlines of the approximations to the velocity field $-a \nabla u$ for the approximation given by RT mixed method of order 1 (left) as well as that given by a postprocessing of the CG method of order 2 (right). The results are very similar. Notice that the singularity of the flow around the corners of the low permeability region $(.25, .25) \times (.75, .75)$ makes this a hard test problem.

We discuss the postprocessing procedure for a hybridized version of the CG method, although it can be applied directly to the standard CG formulation. This is not only because it is easier to understand the first step of the postprocessing using the hybridized formulation (as mentioned previously), but also because the hybridized method has interesting features in its own right. The hybridized CG method is obtained as a natural extension of the new perspective introduced in [15] for hybridizing



Boundary conditions:

$$u = 0 \quad \text{on } \{(1, y) : y \in [0, 1]\},$$

$$\mathbf{q} \cdot \mathbf{n} = \begin{cases} -1 & \{(0, y) : y \in [0, .5]\}, \\ 0 & \{(0, y) : y \in [.5, 1]\}, \\ 0 & \{(x, 1) : x \in [0, 1]\}, \\ 0 & \{(x, 0) : x \in [0, 1]\}. \end{cases}$$

FIG. 1.1. The computational domain $\Omega = (0, 1)^2$ with a uniform 8×8 mesh. The region of low permeability is indicated in dark gray.

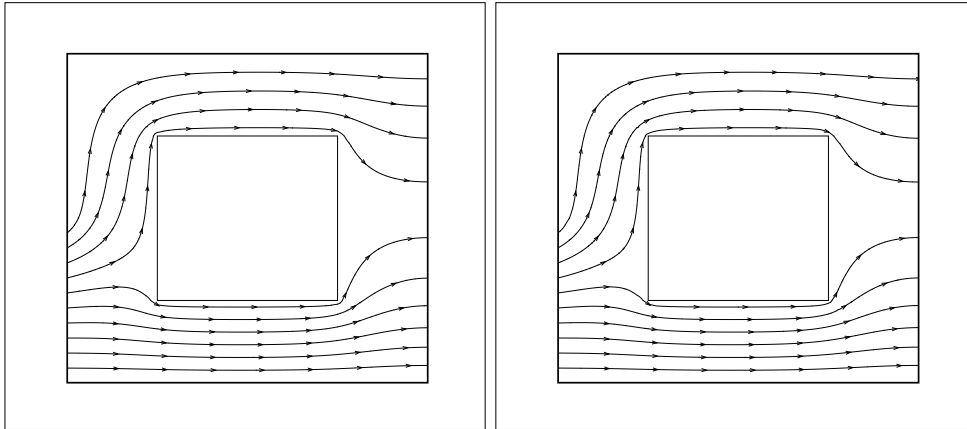


FIG. 1.2. Streamlines of the approximate fluxes for the RT_1 method (left) and the RT_1 -postprocessed CG_2 method (right) obtained using a uniform 32×32 mesh.

mixed methods. It can be briefly described in two steps. First, we express the approximate solution of the CG method u_h in terms of the data components (g and f) and a *Lagrange multiplier* λ_h . It turns out that for the CG method, λ_h is nothing but the restriction of u_h to the faces of the elements of the triangulation. The second step consists in showing that λ_h can be characterized as the only element of certain set $M_h(g)$ satisfying a weak formulation of the form

$$a_h(\lambda_h, \mu) = b_h(\mu) \quad \text{for all } \mu \in M_h(0). \quad (1.5)$$

This formulation was also obtained in [5] with the purpose of devising efficient substructuring preconditioners for the CG method.

Hybridization in the context of mixed methods is different from what goes by the name of static condensation in the engineering literature, because the former gives extra information through the Lagrange multiplier, a solution component absent in static condensation. However, in the hybridized CG case, the fact that the Lagrange multiplier λ_h equals u_h on the element interfaces implies that hybridization and static

condensation coincides, except when we have variable degree elements and hanging nodes. In the static condensation approach, the degrees of freedom of the approximate solution u_h must be very carefully chosen in order to ensure the required continuity across interelement boundaries. The data structures needed to enforce such continuity for variable-degree approximations and hanging nodes have attained a high degree of sophistication – see, for example, [21, 22]. On the other hand, if we use the hybridized version of the CG method (1.5), there is no need to enforce any continuity constraint at all. We apply CG on each element without caring about continuity restrictions, as the continuity is automatically enforced by the equations of the method provided we pick a suitable Lagrange multiplier space $M_h(0)$.

The paper is organized as follows. In Section 2, we present the hybridized CG method and briefly discuss the result characterizing λ_h as the unique solution of (1.5). We also discuss extensions to the variable-degree case and hanging nodes. In Section 3, we describe the construction of $\mathbf{H}(\text{div}, \Omega)$ -conforming approximation to the flux. We state the error estimates of the flux approximation and the results on the relationship between our method and the corresponding RT mixed method. We explain how to explicitly construct a local basis for the space required to compute a single-valued numerical flux trace. An estimate of the conditioning of the global system that arises also appears in this section. In Section 4, we give all the proofs of the theorems. A numerical study of the approximation properties of these approximations is presented in Section 5. We end with some concluding remarks in Section 6.

2. Characterization of the Lagrange multiplier. We begin this section by hybridizing the CG method. We then state, discuss and prove the main result of this section which characterizes the Lagrange multiplier, Theorem 2.1.

2.1. The hybridized CG method. To hybridize the CG method, we relax the continuity restriction and impose it back through suitably chosen new equations. Since the continuity restriction is enforced in the sets $V_h(\cdot)$, to relax it means to work instead with the space

$$V_h = \{v \in L^2(\Omega) : v|_K \in \mathcal{P}_k(K) \quad \text{for all } K \in \mathcal{T}_h\}. \quad (2.1a)$$

The new approximation U_h in V_h must however coincide with u_h , which means, in particular, that it has to be continuous. To enforce the continuity of U_h across interelement boundaries, we force U_h to be equal to the Lagrange multiplier λ_h which we take in

$$M_h(g) = \{\mu \in \mathcal{C}^0(\mathcal{E}_h) : \mu|_e \in \mathcal{P}_k(e) \text{ for all } e \in \mathcal{E}_h, \mu = g \text{ on } \Gamma_D\}, \quad (2.1b)$$

where

$$\mathcal{E}_h = \{e : e \text{ is a face of } K \text{ for all } K \in \mathcal{T}_h\}. \quad (2.1c)$$

Notice that we are implicitly assuming that the triangulation \mathcal{T}_h does not have hanging nodes. To ensure that $U_h = u_h$, we are going to use an auxiliary variable which approximates $\mathbf{q} \cdot \mathbf{n} = -a\nabla u \cdot \mathbf{n}$ on ∂K , for each element K . This additional variable is denoted by $q_{n,h}$ and will be taken in the space

$$W_h = \{p \in L^2(\{\partial K : K \in \mathcal{T}_h\}) : p|_{\partial K} = v|_{\partial K} \quad \text{for } v \in V_h\}. \quad (2.1d)$$

Note that $p \in W_h$ is double-valued in the interior faces of the elements $K \in \mathcal{T}_h$. Thus the hybridized method seeks an approximation to $(u|_{K \in \mathcal{T}_h}, u|_{\mathcal{E}_h}, \mathbf{q} \cdot \mathbf{n}|_{\partial K, K \in \mathcal{T}_h})$, $(U_h, \lambda_h, q_{n,h})$ in the space $V_h \times M_h(g) \times W_h$. It is defined by

$$\sum_{K \in \mathcal{T}_h} (a \nabla U_h, \nabla v)_K + \sum_{K \in \mathcal{T}_h} \langle q_{n,h}, v \rangle_{\partial K} = (f, v)_\Omega \quad \text{for all } v \in V_h, \quad (2.2a)$$

$$U_h = \lambda_h \quad \text{on } \mathcal{E}_h, \quad (2.2b)$$

$$\sum_{K \in \mathcal{T}_h} \langle q_{n,h}, \mu \rangle_{\partial K} = \langle \mathbf{q}_N, \mu \rangle_{\Gamma_N} \quad \text{for all } \mu \in M_h(0). \quad (2.2c)$$

Notice that, by the definition of the space $M_h(0)$, (2.1b), $\mu = v|_{\mathcal{E}_h}$ belongs to $M_h(0)$ whenever $v \in \mathbf{V}_h(0)$. This implies that the last equation can be rewritten as

$$\langle \llbracket q_{n,h} \rrbracket, v \rangle_{\mathcal{E}_h} = \langle \mathbf{q}_N, v \rangle_{\Gamma_N} \quad \text{for all } v \in \mathbf{V}_h(0),$$

where the *jump* of the approximate normal component of the flux is

$$\llbracket q_{n,h} \rrbracket := \begin{cases} q_{n,h}|_{\partial K^+} + q_{n,h}|_{\partial K^-} & \text{on the face } e = \partial K^+ \cap \partial K^-, \\ q_{n,h} & \text{on the face } e = \partial K \cap \partial \Omega. \end{cases}$$

We thus see that it enforces a weak continuity of the inter-element boundary of the jump of this variable; this is why we call it the *jump condition*. Next, we see that this condition ensures that $U_h = u_h$.

PROPOSITION 2.1. *There exists a unique function $(U_h, \lambda_h, q_{n,h})$ in the space $V_h \times M_h(g) \times W_h$ satisfying the formulation (2.2). Moreover,*

$$U_h = u_h \text{ on } \Omega \quad \text{and} \quad \lambda_h = u_h \text{ on } \mathcal{E}_h.$$

Proof. Since $\lambda_h \in M_h(g)$ and $U_h \in V_h$, we have that $U_h \in \mathbf{V}_h(g)$. Moreover, since $\mathbf{V}_h(0) \subset V_h$, by the equation (2.2a), we have

$$(a \nabla U_h, \nabla v)_\Omega + \sum_{K \in \mathcal{T}_h} \langle q_{n,h}, v \rangle_{\partial K} = (f, v)_\Omega \quad \text{for all } v \in \mathbf{V}_h(0),$$

and, by the jump condition (2.2c),

$$(a \nabla U_h, \nabla v)_\Omega = (f, v)_\Omega - \langle \mathbf{q}_N, v \rangle_{\Gamma_N} \quad \text{for all } v \in \mathbf{V}_h(0).$$

By the uniqueness of the approximate of the CG method, we immediately obtain that $U_h = u_h$ on Ω and, as a consequence, that $\lambda_h = u_h$ on \mathcal{E}_h .

It only remains to prove that the function $q_{n,h}$ exists and is unique. This is equivalent to proving that the trivial solution is the only solution of

$$\sum_{K \in \mathcal{T}_h} \langle q_{n,h}, v \rangle_{\partial K} = 0 \quad \text{for all } v \in V_h.$$

Since $q_{n,h} \in W_h$, there is a $w \in V_h$ such that $q_{n,h} = w$. Taking $v = w$ in the above equation, we conclude that $q_{n,h} \equiv 0$, as wanted. This completes the proof. \square

2.2. Characterization of the Lagrange multiplier λ_h . Next, we show how to eliminate the unknowns U_h and $q_{n,h}$ from the equations (2.2) and obtain a formulation only for λ_h . The discussion here is a straightforward generalization of well known results in domain decomposition [5] to the case when subdomains are reduced to elements. Analogous to the discrete harmonic extensions of [5], we now define a local lifting.

The lifting associates to each $\mathbf{m} \in M_h(\cdot)$ the pair of functions $(\mathcal{U}\mathbf{m}, \mathcal{Q}_n\mathbf{m}) \in \mathcal{P}_k(K) \times \{v|_{\partial K} : v \in \mathcal{P}_k(K)\}$ defined by requiring that

$$(a\nabla\mathcal{U}\mathbf{m}, \nabla v)_K + \langle \mathcal{Q}_n\mathbf{m}, v \rangle_{\partial K} = 0 \quad \text{for all } v \in \mathcal{P}_k(K) \quad (2.3a)$$

$$\mathcal{U}\mathbf{m} = \mathbf{m} \quad \text{on } \partial K. \quad (2.3b)$$

In addition, we define a second local mapping that associates to the function $f \in L^2(\Omega)$ the pair of functions $(\mathcal{U}f, \mathcal{Q}_nf) \in \mathcal{P}_k(K) \times \{v|_{\partial K}, v \in \mathcal{P}_k(K)\}$ defined by

$$(a\nabla\mathcal{U}f, \nabla v)_K + \langle \mathcal{Q}_nf, v \rangle_{\partial K} = (f, v)_K \quad \text{for all } v \in \mathcal{P}_k(K) \quad (2.4a)$$

$$\mathcal{U}f = 0 \quad \text{on } \partial K. \quad (2.4b)$$

Notice that $(\mathcal{U}\mathbf{m}, \mathcal{Q}_n\mathbf{m})$ and $(\mathcal{U}f, \mathcal{Q}_nf)$ are approximations to the solutions of

$$-\operatorname{div}(a\nabla u) = 0, \quad -\operatorname{div}(a\nabla u) = f, \quad \text{on } K, \quad (2.5a)$$

$$u = \mathbf{m}, \quad u = 0, \quad \text{on } \partial K. \quad (2.5b)$$

We are now ready to state the characterization of the CG solution in terms of the Lagrange multiplier, whose proof is at the end of this section.

THEOREM 2.1. *Let $(U_h, \lambda_h, q_{n,h})$ be the solution of the hybridized version of the CG method. Then*

$$U_h = \mathcal{U}\lambda_h + \mathcal{U}f \quad \text{and} \quad q_{n,h} = \mathcal{Q}_n\lambda_h + \mathcal{Q}_nf.$$

Moreover, the Lagrange multiplier $\lambda_h \in M_h(g)$ is the unique solution of

$$\sum_{K \in \mathcal{T}_h} (a\nabla\mathcal{U}\lambda_h, \nabla\mathcal{U}\mu)_K = (f, \mathcal{U}\mu)_\Omega - \langle \mathbf{q}_N, \mu \rangle_{\Gamma_N} \quad \text{for all } \mu \in M_h(0).$$

Like other hybridized formulations, the utility of such a result lies in its ease of computation of a ‘‘stiffness matrix’’ for the Lagrange multiplier. Furthermore, once λ_h has been obtained, U_h and $q_{n,h}$ can be easily computed element by element using the local mappings (2.3) and (2.4).

It is interesting to note that $q_{n,h}|_{\partial K}$ is strongly related to what was denoted by $H^h(K)$ in [25]; in fact, when the element K does not have a face lying on the boundary, these two quantities are identical. However, in [25] they are used to uncover a local conservativity property of the CG method whereas here we use them as an auxiliary means to hybridize it.

Finally, notice that Theorem 2.1 states that the functions $q_{n,h}|_{\partial K}$ need *not* be actually computed to construct the matrix equations for the multiplier λ_h . Indeed, from the definition of the lifting (2.3), we see that we can independently compute $\mathcal{U}\mathbf{m}$ on the element K by solving

$$(a\nabla\mathcal{U}\mathbf{m}, \nabla v)_K = 0 \quad \text{for all } v \in \mathcal{P}_k(K) \text{ such that } v = 0 \text{ on } \partial K,$$

$$\mathcal{U}\mathbf{m} = \mathbf{m} \quad \text{on } \partial K.$$

This implies that \mathcal{U}_m can be written as a linear combination of

$$\dim \mathcal{P}_k(K) - \dim \mathcal{P}_{k-3}(K) = \binom{k+N}{N} - \binom{k-3+N}{N}$$

basis functions, when $k \geq 3$, of course. In two space dimension ($N = 2$), this means that instead of working with a basis of $(k+2)(k+1)/2$ functions, we can work with a basis of only $3k$ functions. In three space dimensions, it means that instead of working with $(k+3)(k+2)(k+1)/6$ basis functions, we only have to work with $(3k^2 + 3k + 2)/2$. Thus, the computation of \mathcal{U}_m can be rendered extremely efficient, especially for high polynomial degrees k . This is especially true if the exact solution is harmonic, that is, if $f = 0$.

2.3. Variable-degree approximations and hanging nodes. The hybridized CG formulation is particularly attractive for variable-degree approximate spaces and meshes with hanging nodes.

We begin by briefly showing how to extend our previous results to the variable-degree case, that is, to the case in which the approximate solution u_h belongs to

$$V_h(s) = \{v \in \mathcal{C}^0(\Omega) : v|_K \in \mathcal{P}_{k(K)}(K), v = s \text{ on } \Gamma_D\},$$

where the polynomial degree $k(K)$ now varies with as K varies within \mathcal{T}_h . We can then hybridize the resulting CG method, just as we hybridized the uniform-degree CG method, if we take

$$\begin{aligned} M_h(g) &= \{\mu \in \mathcal{C}^0(\mathcal{E}_h) : \mu|_e \in \mathcal{P}_{k(e)}(e) \text{ for all } e \in \mathcal{E}_h, \mu = g \text{ on } \Gamma_D\}, \\ W_h &= \{w \in L^2(\mathcal{E}_h) : w = v|_{\partial K}, v \in \mathcal{P}_{k(K)}(K) \text{ for all } K \in \mathcal{T}_h\}, \\ V_h &= \{v \in L^2(\mathcal{T}_h) : v|_K \in \mathcal{P}_{k(K)}(K) \text{ for all } K \in \mathcal{T}_h\}. \end{aligned} \quad (2.6)$$

With this, the burden of enforcing the continuity constraint is automatically dealt with by the local mappings which are defined *exactly* as before with k replaced by $k(K)$. While the current practice for implementing variable degree methods is via transitional basis functions and the minimum degree rule [21], the above hybridization approach removes the continuity matching considerations from the design of shape functions.

To end this subsection, let us briefly address the case of hanging nodes, which is also surprisingly simple to handle by hybridization, even in three dimensions. We only have to define the multiplier space $M_h(g)$ in a suitable way. In fact, we can continue to define $M_h(g)$ by (2.6) provided we redefine the set \mathcal{E}_h there. To do this, we need to introduce the notion of a *maximal face*. A face e of an element $K \in \mathcal{T}_h$ is said to be a maximal face of the triangulation \mathcal{T}_h if it lies on $\partial\Omega$ or whenever there is another element $K' \in \mathcal{T}_h$ such that $e \cap \partial K'$ has non-zero $(N-1)$ -Lebesgue measure, $e \cap \partial K'$ is a face of K' . An illustration is given in Fig. 2.1. The new definition of \mathcal{E}_h is simply

$$\mathcal{E}_h = \{e : e \text{ is a maximal face of the triangulation } \mathcal{T}_h\}. \quad (2.7)$$

2.4. Proof of Theorem 2.1. To prove this result, we need the following lemma.

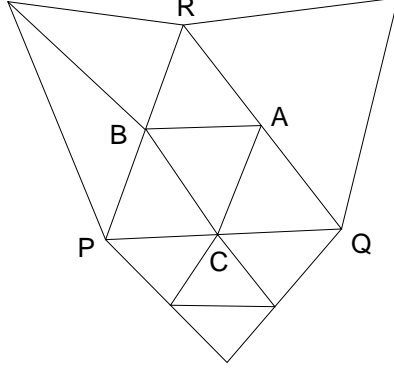


FIG. 2.1. Detail of a triangulation. The faces RQ , BR and AB are maximal whereas the faces RA and AQ are not.

LEMMA 2.2 (elementary identities). We have, for any $\mathbf{m} \in M_h(\cdot)$, $\mu \in M_h(0)$, and $f \in L^2(\Omega)$,

$$\begin{aligned} \text{(i)} \quad & - \sum_{K \in \mathcal{T}_h} \langle \mathcal{Q}_n \mathbf{m}, \mu \rangle_{\partial K} = \sum_{K \in \mathcal{T}_h} (a \nabla \mathcal{U} \mathbf{m}, \nabla \mathcal{U} \mu)_K, \\ \text{(ii)} \quad & - \sum_{K \in \mathcal{T}_h} \langle \mathcal{Q}_n f, \mu \rangle_{\partial K} = -(f, \mathcal{U} \mu)_\Omega. \end{aligned}$$

Proof. We have

$$\begin{aligned} - \sum_{K \in \mathcal{T}_h} \langle \mathcal{Q}_n \mathbf{m}, \mu \rangle_{\partial K} &= - \sum_{K \in \mathcal{T}_h} \langle \mathcal{Q}_n \mathbf{m}, \mathcal{U} \mu \rangle_{\partial K} \quad \text{by (2.3b)}, \\ &= \sum_{K \in \mathcal{T}_h} (a \nabla \mathcal{U} \mathbf{m}, \nabla \mathcal{U} \mu)_K \end{aligned}$$

by (2.3a). This proves the first identity.

Let us prove the second identity. We have

$$\begin{aligned} - \sum_{K \in \mathcal{T}_h} \langle \mathcal{Q}_n f, \mu \rangle_{\partial K} &= - \sum_{K \in \mathcal{T}_h} \langle \mathcal{Q}_n f, \mathcal{U} \mu \rangle_{\partial K} \quad \text{by (2.3b)}, \\ &= -(f, \mathcal{U} \mu)_\Omega + \sum_{K \in \mathcal{T}_h} (a \nabla \mathcal{U} f, \nabla \mathcal{U} \mu)_K \quad \text{by (2.4a)}, \\ &= -(f, \mathcal{U} \mu)_\Omega + \sum_{K \in \mathcal{T}_h} \langle \mathcal{Q}_n \mu, \mathcal{U} f \rangle_{\partial K} \quad \text{by (2.3a)}, \\ &= -(f, \mathcal{U} \mu)_\Omega \end{aligned}$$

by (2.4b). This completes the proof. \square

Proof of Theorem 2.1. By the definition of the local mappings, we have that

$$U_h = \mathcal{U} \lambda_h + \mathcal{U} f \quad \text{and} \quad q_{n,h} = \mathcal{Q}_n \lambda_h + \mathcal{Q}_n f.$$

This implies that the third equation in the definition of the hybridized version of the CG method (2.2c) can be rewritten as

$$\sum_{K \in \mathcal{T}_h} \langle \mathcal{Q}_n \lambda_h + \mathcal{Q}_n f, \mu \rangle_{\partial K} = \langle \mathbf{q}_N, v \rangle_{\Gamma_N} \quad \text{for all } \mu \in M_h(0),$$

or, by Lemma 2.2, as

$$-\sum_{K \in \mathcal{T}_h} (a \nabla \mathcal{U}_{\lambda_h}, \nabla \mathcal{U}_{\mu})_K + (f, \mathcal{U}_{\mu})_{\Omega} = \langle \mathbf{q}_{\mathbf{N}}, v \rangle_{\Gamma_N}.$$

This completes the proof. \square

3. An $\mathbf{H}(\text{div}, \Omega)$ -conforming approximation of the flux. In this section, we define an $\mathbf{H}(\text{div}, \Omega)$ -conforming approximation, \mathbf{q}_h , to the flux $\mathbf{q} = -a \nabla u$. Then we state, discuss and prove a theorem about the quality of the resulting approximation as well as the complexity of the algorithm needed to compute it. Although all considerations in this section hold for the variable degree case, for simplicity we restrict ourselves to the uniform degree case spaces defined in (2.1) with no hanging nodes.

3.1. The new approximation to the flux. The key step in the construction of an $\mathbf{H}(\text{div}, \Omega)$ -conforming approximation \mathbf{q}_h is the definition of its normal component on the element interfaces. The function $q_{n,h}$ represents an approximation to the normal component of the flux, but unfortunately it is not a single-valued function in general. Notice however that by equation (2.2a), we have

$$\sum_{K \in \mathcal{T}_h} (a \nabla U_h, \nabla v)_K + \sum_{K \in \mathcal{T}_h} \langle q_{n,h}, v \rangle_{\partial K} = (f, v)_{\Omega} \quad \text{for all } v \in V_h,$$

so the possibility of constructing a single-valued function $\widehat{\mathbf{q}}_h$ satisfying

$$\sum_{K \in \mathcal{T}_h} \langle \widehat{\mathbf{q}}_h \cdot \mathbf{n}, v \rangle_{\partial K} = \sum_{K \in \mathcal{T}_h} \langle q_{n,h}, v \rangle_{\partial K} \quad \text{for all } v \in V_h \quad (3.1)$$

opens up. If such a $\widehat{\mathbf{q}}_h$ could be constructed, we could then define the approximate flux \mathbf{q}_h as follows: On any simplicial element K , we can set \mathbf{q}_h in the Raviart-Thomas space

$$V_{RT\ell}(K) := \mathcal{P}_{\ell}(K)^N + \mathbf{x} \mathcal{P}_{\ell}(K), \quad (3.2a)$$

by requiring that

$$\langle \mathbf{q}_h \cdot \mathbf{n}, v \rangle_e = \langle \widehat{\mathbf{q}}_h \cdot \mathbf{n}, v \rangle_e \quad \forall v \in \mathcal{P}_{\ell}(e), \text{ for any face } e \subset \partial K, \quad (3.2b)$$

$$(\mathbf{q}_h, \mathbf{v})_K = -(a \nabla U_h, \mathbf{v})_K \quad \forall \mathbf{v} \in \mathcal{P}_{\ell-1}(K)^N. \quad (3.2c)$$

Note that the definition (3.2) is a modification of the well known Raviart-Thomas projection, (see (3.12) later). A similar projection was suggested in [4] in the framework of the interior penalty method for Darcy's law and in [19] in the framework of LDG methods for the Navier-Stokes equations. It is not difficult to show that a \mathbf{q}_h constructed by (3.2b) belongs to $\mathbf{H}(\text{div}, \Omega)$, thanks to the single-valuedness of the normal component of the numerical trace $\widehat{\mathbf{q}}_h$.

Such a construction will yield a flux \mathbf{q}_h that is conservative whenever $\ell \leq k$. Indeed, we can rewrite the equation (2.2a)

$$-\sum_{K \in \mathcal{T}_h} (\mathbf{q}_h, \nabla v)_K + \sum_{K \in \mathcal{T}_h} \langle \mathbf{q}_h \cdot \mathbf{n}, v \rangle_{\partial K} = (f, v)_{\Omega}$$

for all v such that $v|_K \in \mathcal{P}_\ell(K)$ for all $K \in \mathcal{T}_h$. Hence, if we take v to be the characteristic function of a discrete subdomain D_h formed by the union of some elements $K \in \mathcal{T}_h$, we obtain

$$\langle \mathbf{q}_h \cdot \mathbf{n}, 1 \rangle_{\partial D_h} = (f, 1)_{D_h},$$

which is the same as the exact conservation property (1.4).

It is interesting to see that there is an extremely simple relation between the normal component of \mathbf{q}_h and the approximation $H^h(\cdot)$ defined in [25] or, equivalently, what is called $\bar{\sigma}_h$ in [9]. Indeed, if D is any union of elements $K \in \mathcal{T}_h$, then from the definition of \mathbf{q}_h and that of $H^h(D)$, see equations (47) and (57) in [25], we have that

$$\langle \mathbf{q}_h \cdot \mathbf{n} - H^h(D), v \rangle_{\partial D} = 0 \quad \forall v \in V_{h,D} := \{v \in V_h \cap \mathcal{C}^0(D)\}.$$

Since $H^h(D)|_{\partial D}$ belongs to the space of traces on ∂D of the functions in $V_{h,D}$, we see that $H^h(D)$ is the L^2 -projection of $\mathbf{q}_h \cdot \mathbf{n}$ into such space.

3.2. The numerical trace $\hat{\mathbf{q}}_h$. It remains to find the numerical trace $\hat{\mathbf{q}}_h$. To do that, we first notice that if $\hat{\mathbf{q}}_h$ is single valued, the equation (3.1) takes the form

$$\langle \hat{\mathbf{q}}_h, \llbracket v \mathbf{n} \rrbracket \rangle_{\mathcal{E}_h} = \sum_{K \in \mathcal{T}_h} \langle q_{n,h}, v \rangle_{\partial K} \quad \text{for all } v \in V_h.$$

Since the flux on Γ_N is given to be \mathbf{q}_N , incorporating this information into the above equation, we get

$$\langle \hat{\mathbf{q}}_h, \llbracket v \mathbf{n} \rrbracket \rangle_{\mathcal{E}_h \setminus \Gamma_N} = \sum_{K \in \mathcal{T}_h} \langle q_{n,h}, v \rangle_{\partial K} - \langle \mathbf{q}_N, v \rangle_{\Gamma_N} \quad \text{for all } v \in V_h. \quad (3.3)$$

In the one dimensional case $\Omega = (0, 1)$, this equation can be readily solved. Indeed, we have that

$$\hat{\mathbf{q}}_h(x_i) = \begin{cases} q_{n,h}(1^-) & \text{if } x_i = 1, \\ q_{n,h}(x_i^-) = -q_{n,h}(x_i^+) & \text{if } x_i \text{ is an interior node,} \\ -q_{n,h}(0^+) & \text{if } x_i = 0, \end{cases}$$

where we have used the fact that, by equation (2.2c), $q_{n,h}(x_i^-) + q_{n,h}(x_i^+) = 0$ on all interior nodes x_i . Let us find expressions for $q_{n,h}$ in terms of the data f and u_h . By (2.2c), $q_{n,h} = \mathbf{q}_N$ on Γ_N , and we get that

$$\hat{\mathbf{q}}_h = \mathbf{q}_N \quad \text{on } \Gamma_N.$$

To find $q_{n,h}$ in the remaining nodes, we simply use equations (2.2a). Thus, if we let x_i be any node not lying on Γ_N , and let φ_i^+ (resp., φ_i^-) be the linear function with support the interval $I_i^+ = (x_i, x_{i+1})$ (resp., $I_i^- = (x_{i-1}, x_i)$) such that $\varphi_i^+(x_i) = 1$ and $\varphi_i^+(x_{i+1}) = 0$ (resp., $\varphi_i^-(x_i) = 1$ and $\varphi_i^-(x_{i-1}) = 0$), we obtain that

$$\hat{\mathbf{q}}_h(x_i) = \mp \left(a \frac{d}{dx} u_h, \frac{d}{dx} \varphi_i^\pm \right)_{I_i^\pm} \pm (f, \varphi_i^\pm)_{I_i^\pm}.$$

These expressions have been known for a long time – see the work by J. Wheeler [33] and M.F. Wheeler [34]. Moreover, in [34], it was shown that the approximation $\hat{\mathbf{q}}_h$

TABLE 3.1

The L^2 norm of the error $\mathbf{q} - \mathbf{q}_h$ when the wrong flux trace is used. The parameters are the same as that described in Fig. 3.1

h	$k = 1$		$k = 2$		$k = 3$	
	error	order	error	order	error	order
1/8	0.11E+00	0.46	0.62E-01	0.41	0.40E-01	0.46
1/16	0.77E-01	0.48	0.45E-01	0.46	0.29E-01	0.48
1/32	0.55E-01	0.49	0.32E-01	0.48	0.21E-01	0.49
1/64	0.39E-01	0.50	0.23E-01	0.49	0.15E-01	0.50

superconverges with order $2k$ if the CG method uses polynomials of degree k and is exact, that is,

$$\widehat{\mathbf{q}}_h(x_i) = -a \frac{d}{dx} u(x_i),$$

whenever a is a constant.

Extensions of the above approach to the multi-dimensional case for obtaining approximations to the normal component of $\widehat{\mathbf{q}}_h$ have been explored by many authors. See [9] for an overview and recent developments, [8] for early computational tricks, [11] for a fully developed technique, and [14, 28] for rigorous error estimates. Here, we do not use this approach.

Instead, we begin by noting that, from the formulation (3.3), it is clear that we can only obtain a projection of $\widehat{\mathbf{q}}_h$ into the space of jumps

$$\mathfrak{J}_h = \{ \llbracket w \mathbf{n} \rrbracket |_{\mathcal{E}_h \setminus \Gamma_N} : w \in V_h \}. \quad (3.4)$$

This may seem to suggest choosing $\widehat{\mathbf{q}}_h$ in \mathfrak{J}_h . We have experimented with such a choice. The results of one such experiment are reported in Fig. 3.1 and Table 3.1. We found that such a flux approximation is often reasonable away from the boundary, but near $\partial\Omega$ the degradation of the approximation is clearly evident for some problems. Furthermore, from a theoretical standpoint, such a choice appears dubious as the space \mathfrak{J}_h does not contain the constant function. For these reasons, we do *not* advocate it.

The solution we found practically acceptable as well as theoretically sound proceeds by borrowing ideas from the development of the discontinuous Galerkin (DG) method. We select the following form for the numerical trace:

$$\widehat{\mathbf{q}}_h = \begin{cases} \mathbf{q}_N \mathbf{n} & \text{on } \Gamma_N, \\ -a \nabla U_h + \alpha \mathbf{J}_h & \text{on } \Gamma_D, \\ -\llbracket a \nabla U_h \rrbracket - \beta \llbracket a \nabla U_h \cdot \mathbf{n} \rrbracket + \alpha \mathbf{J}_h & \text{on } \mathcal{E}_h \setminus \partial\Omega, \end{cases} \quad (3.5)$$

where α and β are single valued bounded (scalar and vector, resp.) functions on $\mathcal{E}_h \setminus \partial\Omega$, $\alpha > 0$, and \mathbf{J}_h is an element of the space of jumps \mathfrak{J}_h to be determined. A typical choice of the parameters that we have found adequate in our numerical experiments (on uniform meshes) is $\beta \equiv \mathbf{0}$ and $\alpha \equiv 1$ (also see Theorem 3.2 for better choices of α on highly non-uniform meshes). Here, we have used the now standard DG notation (cf., e.g., [2]),

$$\llbracket v \rrbracket = \begin{cases} \frac{1}{2} (v^+ + v^-) & \text{on } \mathcal{E}_h^\circ, \\ v & \text{on } \partial\Omega, \end{cases} \quad (3.6a)$$

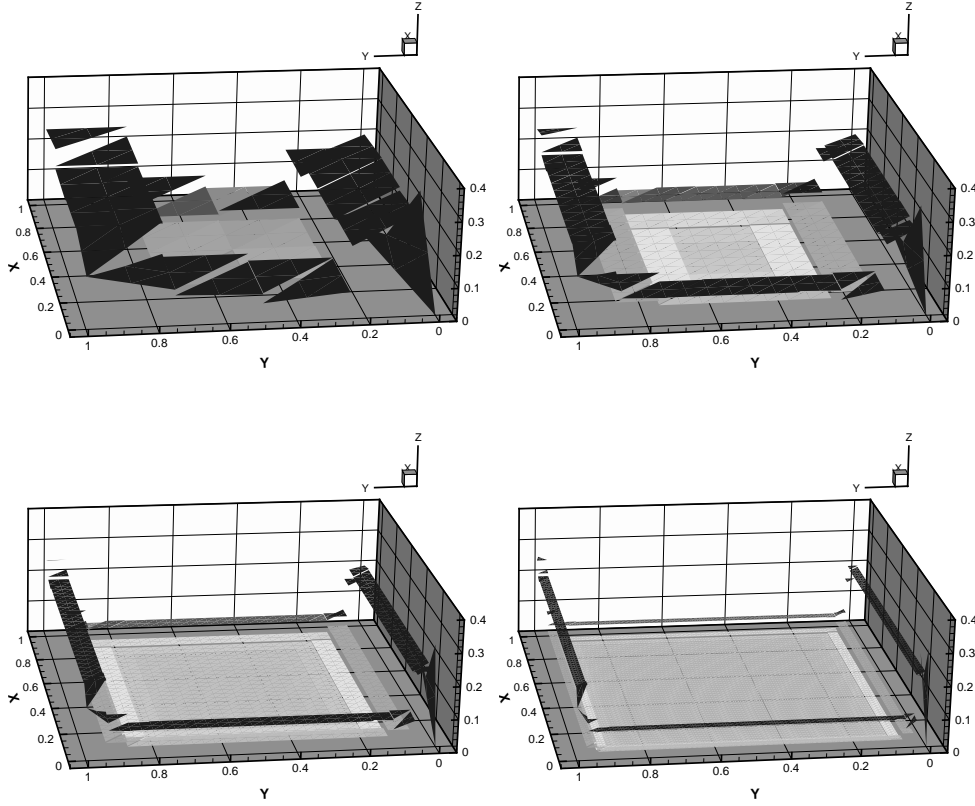


FIG. 3.1. Plots of the error $|q - q_h|$ with the wrong flux trace choice on successively refined meshes. Computational details: Here q_h is obtained by (3.2) with \hat{q}_h chosen as the unique function in \mathfrak{J}_h satisfying (3.3). The parameters are $a = Id$, $f = 0$, $\Omega = (0, 1) \times (0, 1)$, $\Gamma_D = \{0\} \times (0, 1)$, the polynomial degrees are $k = 1$ and $\ell = k - 1$ (for postprocessing), and the boundary conditions are set in such a way that the exact solution is $u(x, y) = 1 + x$. We see that while the error is small far from the boundary, near the boundary the error remains of order one. Therefore, we expect to see an order of convergence of $1/2$ in the L^2 -norm. This is confirmed in Table 3.1.

and

$$[[v\mathbf{n}]] = \begin{cases} v^+ \mathbf{n}^+ + v^- \mathbf{n}^- & \text{on } \mathcal{E}_h^\circ, \\ v \mathbf{n} & \text{on } \partial\Omega, \end{cases} \quad (3.6b)$$

where for a piecewise smooth function v , the traces from either side of a mesh face (edge) e are denoted by $v^\pm(\mathbf{x}) = \lim_{\epsilon \downarrow 0} v(\mathbf{x} - \epsilon \mathbf{n}^\pm)$ for all \mathbf{x} in e (and \mathbf{n}^\pm denotes the corresponding unit outward normal on e from either side).

Next, we insert the expression we have selected in (3.5) for the numerical flux \hat{q}_h into (3.3). This gives us an equation for \mathbf{J}_h :

$$\begin{aligned} \langle \alpha \mathbf{J}_h, [[v\mathbf{n}]] \rangle_{\mathcal{E}_h \setminus \Gamma_N} &= \langle \{ \{ a \nabla U_h \} \} + \beta [[a \nabla U_h \cdot \mathbf{n}]], [[v\mathbf{n}]] \rangle_{\mathcal{E}_h \setminus \partial\Omega} \\ &+ \langle a \nabla U_h, v \mathbf{n} \rangle_{\Gamma_D} + \sum_{K \in \mathcal{T}_h} \langle q_{n,h}, v \rangle_{\partial K} - \langle q_N, v \rangle_{\Gamma_N} \end{aligned} \quad (3.7)$$

The computation of \mathbf{J}_h from this equation requires solving a global, but well conditioned, system. The details involved are discussed in the next subsection. For the moment, observe that if we are using the hybridized form of the CG method and have already computed U_h and $q_{n,h}$, the right hand side of (3.7) can be computed using integrations only on element boundaries.

On the other hand, if we have computed U_h using a standard CG implementation without hybridization (and so do not have access to $q_{n,h}$), we can still use the above postprocessing. Indeed, by using the equation (2.2a), we can transform (3.7) to an equation that is more convenient for this case:

$$\begin{aligned} \langle \alpha \mathbf{J}_h, \llbracket v\mathbf{n} \rrbracket \rangle_{\mathcal{E}_h \setminus \Gamma_N} &= \langle \{ \{ a \nabla U_h \} \} + \beta \llbracket a \nabla U_h \cdot \mathbf{n} \rrbracket, \llbracket v\mathbf{n} \rrbracket \rangle_{\mathcal{E}_h \setminus \partial\Omega} + \\ &\quad \langle a \nabla U_h, v\mathbf{n} \rangle_{\Gamma_D} - \sum_{K \in \mathcal{T}_h} (a \nabla U_h, \nabla v)_K + (f, v)_\Omega - \langle \mathbf{q}_N, v \rangle_{\Gamma_N}. \end{aligned} \quad (3.8)$$

Observe that whenever $\llbracket v\mathbf{n} \rrbracket = 0$ on $\mathcal{E}_h \setminus \Gamma_N$, i.e., whenever $v \in \mathbf{V}_h(0)$, the right-hand side of the above equation is equal to zero by the definition of the CG method, equation (1.3), and Proposition 2.1. Therefore, this equation defines \mathbf{J}_h uniquely. When using spaces of high polynomial degrees, it is preferable to use (3.7) instead of (3.8) as the former involves faster quadratures.

This completes the definition of the numerical trace $\widehat{\mathbf{q}}_h$. *To summarize*, $\widehat{\mathbf{q}}_h$ is defined by (3.3) wherein \mathbf{J}_h is the unique function in \mathfrak{J}_h satisfying (3.7) or (3.8). Let us point out that this definition of the numerical trace reproduces constant fluxes. More precisely, if $-a \nabla U_h$ is a constant vector, say \mathbf{c} , then $\widehat{\mathbf{q}}_h$ is also \mathbf{c} . To see this, note that in this case we must have $\mathbf{q}_N = \mathbf{c} \cdot \mathbf{n}$ and $f = 0$, so that (3.8) becomes

$$\langle \alpha \mathbf{J}_h, \llbracket v\mathbf{n} \rrbracket \rangle_{\mathcal{E}_h \setminus \Gamma_N} = -\langle \mathbf{c}, \llbracket v\mathbf{n} \rrbracket \rangle_{\mathcal{E}_h \setminus \Gamma_N} + \sum_{K \in \mathcal{T}_h} (\mathbf{c}, \nabla v)_K - \langle \mathbf{c} \cdot \mathbf{n}, v \rangle_{\Gamma_N} = 0.$$

This implies $\mathbf{J}_h \equiv \mathbf{0}$, and hence $\widehat{\mathbf{q}}_h = \mathbf{c}$, as claimed.

Let us end this subsection by relating our approach to compute $\widehat{\mathbf{q}}_h$ to that proposed in [27]. In such approach, the numerical trace $\widehat{\mathbf{q}}_h$ is taken as in (3.5) with $\beta = 0$ and $\alpha = 1/h$, where \mathbf{J}_h is taken in the space

$$\mathfrak{J}_{h,0} = \{ \llbracket w\mathbf{n} \rrbracket \}_{\mathcal{E}_h \setminus \Gamma_N} : w|_K \in \mathcal{P}_0(K) \quad \forall K \in \mathcal{T}_h \},$$

and is defined by requiring that

$$\langle \alpha \mathbf{J}_h, \llbracket v\mathbf{n} \rrbracket \rangle_{\mathcal{E}_h \setminus \Gamma_N} = \langle \{ \{ a \nabla U_h \} \}, \llbracket v\mathbf{n} \rrbracket \rangle_{\mathcal{E}_h \setminus \partial\Omega} + \langle a \nabla U_h, v\mathbf{n} \rangle_{\Gamma_D} + (f, v)_\Omega - \langle \mathbf{q}_N, v \rangle_{\Gamma_N},$$

be satisfied for all $v \in \mathfrak{J}_{h,0}$. Note that our \mathbf{J}_h also satisfies this formulation, in the case in which $\beta = 0$ and $\alpha = 1/h$, since the formulation (3.8) reduces to the one under consideration when $v \in \mathfrak{J}_h$ is restricted to $v \in \mathfrak{J}_{h,0}$.

3.3. The computation of \mathbf{J}_h . Next, we discuss the computation of \mathbf{J}_h through solution of (3.7) or (3.8). First, in order to represent \mathbf{J}_h in computations we need a basis for the space \mathfrak{J}_h of jumps. We construct a local basis for \mathfrak{J}_h extending a similar construction carried out in [17, 18] in the context of Stokes flow. Second, we need to solve for \mathbf{J}_h from (3.7) or (3.8). We show that this can be accomplished by solving a square system whose matrix is well conditioned. Thus, we conclude that the computational complexity needed to solve for \mathbf{J}_h is negligible with respect to that required to solve for the multiplier λ_h . Proofs of all results here are given in Section 4.

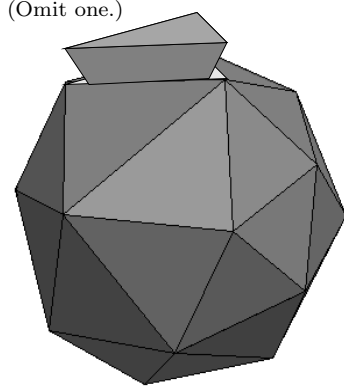


FIG. 3.2. Illustration of the elements connected to one vertex \mathbf{z} and the omission of one $\phi_{\mathbf{z},K}$ to construct a basis.

The basis is easiest to see in the lowest order case (i.e., when $k = 1$). In two dimensions, this basis is closely related to the “wedge” basis functions obtained in [17]. However, in three dimensions, it is different from that given in [18], so let us begin by describing our lowest order basis in three dimensions. For a mesh vertex \mathbf{z} and a mesh element K having \mathbf{z} as a vertex, let $\lambda_{\mathbf{z},K}$ denote the linear function on K which equals one on \mathbf{z} and zero at all other vertices of K and let $\phi_{\mathbf{z},K}$ denote its extension by zero from K to all Ω . Then clearly the restrictions of $[[\phi_{\mathbf{z},K}\mathbf{n}]]$ on $\mathcal{E}_h \setminus \Gamma_N$ are in \mathfrak{F}_h . However, they are not all linearly independent, because for any vertex \mathbf{z} not on Γ_D , the sum of the functions $[[\phi_{\mathbf{z},K}\mathbf{n}]]|_{\mathcal{E}_h \setminus \Gamma_N}$ over all K sharing \mathbf{z} vanishes. Therefore, we must omit one function per vertex to get a basis: For each \mathbf{z} , we define $V_{\mathbf{z}}$ as the set of functions $\phi_{\mathbf{z},K}$ for all K having \mathbf{z} as a vertex. Then for vertices \mathbf{z} not on $\overline{\Gamma}_D$, we define $V_{\mathbf{z}}^*$ as the set obtained by omitting (any) one member of $V_{\mathbf{z}}$ (see Fig. 3.2), while for vertices on $\overline{\Gamma}_D$, we define $V_{\mathbf{z}}^* = V_{\mathbf{z}}$. Then, by a straightforward generalization of the arguments in [17, Proposition 4.2], one can prove that the set

$$\mathfrak{B}^1 = \{ [[\phi\mathbf{n}]]|_{\mathcal{E}_h \setminus \Gamma_N} : \phi \in V_{\mathbf{z}}^* \text{ for all mesh vertices } \mathbf{z} \}$$

is linearly independent, so it forms a basis for \mathfrak{F}_h .

Next, we describe one possible extension of this basis construction to the higher order case. For any given simplex $S \in \mathbb{R}^N$ with vertices $\mathbf{x}_{i,S}$, $i = 1, \dots, N+1$, we define the points in its principal lattice (of order k) [13] by

$$\mathbf{x}_{\alpha,S} = \sum_{j=1}^{N+1} \alpha_j \mathbf{x}_{j,S},$$

where α is taken in $\mathcal{A}_N^k = \{(\alpha_1, \dots, \alpha_{N+1}) : k\alpha_j \in \{0, 1, \dots, k\} \text{ and } \sum_{j=1}^{N+1} \alpha_j = 1\}$. We associate to each point $\mathbf{x}_{\alpha,S}$, the standard Lagrange finite element basis function $v_{\alpha,S}$ defined as the unique function in $\mathcal{P}_k(S)$ satisfying

$$v_{\alpha,S}(\mathbf{x}_{\beta,S}) = \begin{cases} 1 & \text{if } \alpha = \beta, \\ 0 & \text{otherwise,} \end{cases}$$

for all α and β in \mathcal{A}_N^k . Let $\phi_{\alpha,S}$ be the extension by zero to Ω of $v_{\alpha,S}$. Since the basis will be constructed using the jumps of these functions across element interfaces,

we will need to separate the functions associated to the points on element interfaces, which we collect in

$$\mathcal{G}_h^k = \{\mathbf{x}_{\alpha,S} \in \partial K : K \in \mathcal{T}_h, \alpha \in \mathcal{A}_N^k\}. \quad (3.9)$$

To any \mathbf{z} in \mathcal{G}_h^k , we associate more than one $\phi_{\alpha,K}$ if more than one simplex shares \mathbf{z} . We collect these functions in $V_{\mathbf{z}} = \{\phi_{\alpha,K} : \mathbf{x}_{\alpha,K} = \mathbf{z}\}$ and define

$$V_{\mathbf{z}}^* = \begin{cases} V_{\mathbf{z}} & \text{if } \mathbf{z} \in \bar{\Gamma}_D, \\ V_{\mathbf{z}} \setminus \{\phi_{\alpha^*,K^*} \text{ for some } \mathbf{x}_{\alpha^*,K^*} = \mathbf{z}\} & \text{otherwise,} \end{cases}$$

where, like in the lowest order case, we have selected (arbitrarily) one degree of freedom (represented by the multi-index α^* and K^*) for every \mathbf{z} in $\mathcal{G}_h^k \setminus \bar{\Gamma}_D$, and omitted the corresponding Lagrange function ϕ_{α^*,K^*} . With this notation, we have the following result, whose proof follows by generalizing the above mentioned arguments for the lowest order case.

THEOREM 3.1. *The set*

$$\mathfrak{B}^k = \{[\phi \mathbf{n}]|_{\mathcal{E}_h \setminus \Gamma_N} : \phi \in V_{\mathbf{z}}^* \text{ for all } \mathbf{z} \in \mathcal{G}_h^k\}$$

is a basis for \mathfrak{J}_h .

Now that a local basis of the jump space \mathfrak{J}_h has been constructed, we can compute the representation of the jump function \mathbf{J}_h in the basis \mathfrak{B}^k by using the equation (3.7) or (3.8). E.g., to solve for \mathbf{J}_h using (3.7), we begin by introducing an extension operator T_h from the space of jumps \mathfrak{J}_h to the space V_h , constructed in such a way that we have

$$[T_h(\mathbf{J}_h) \mathbf{n}] = \mathbf{J}_h \quad \text{on } \mathcal{E}_h \setminus \Gamma_N, \text{ and} \quad (3.10a)$$

$$T_h(\mathbf{J}_h)|_{K^*, \mathbf{x}_{\alpha^*,K^*}} = 0. \quad (3.10b)$$

Here and elsewhere, we use the notation $w|_{K,\mathbf{r}}$ to denote the limit of the function $w(\mathbf{x})$ as \mathbf{x} approaches \mathbf{r} from within K . One can easily verify that the choice

$$T_h(\mathbf{J}_h) = \sum_{\mathbf{z} \in \mathcal{G}_h^k} \sum_{\phi \in V_{\mathbf{z}}^*} c_{\phi} \phi \quad \text{whenever} \quad \mathbf{J}_h = \sum_{\mathbf{z} \in \mathcal{G}_h^k} \sum_{\phi \in V_{\mathbf{z}}^*} c_{\phi} [\phi \mathbf{n}], \quad (3.11)$$

satisfies both the properties of (3.10). Then, for any $\mathbf{Y}_h \in \mathfrak{J}_h$, setting $v = T_h(\mathbf{Y}_h)$ in (3.7) and using (3.10a), we get that \mathbf{J}_h satisfies

$$\begin{aligned} \langle \alpha \mathbf{J}_h, \mathbf{Y}_h \rangle_{\mathcal{E}_h \setminus \Gamma_N} &= \langle \{a \nabla U_h\} + \beta [a \nabla U_h \cdot \mathbf{n}], [T_h(\mathbf{Y}_h) \mathbf{n}] \rangle_{\mathcal{E}_h \setminus \partial \Omega} \\ &\quad + \langle a \nabla U_h, T_h(\mathbf{Y}_h) \mathbf{n} \rangle_{\Gamma_D} + \sum_{K \in \mathcal{T}_h} \langle q_{n,h}, T_h(\mathbf{Y}_h) \rangle_{\partial K} - \langle \mathbf{q}_N, T_h(\mathbf{Y}_h) \rangle_{\Gamma_N} \\ &\equiv F(\mathbf{Y}_h). \end{aligned}$$

This shows that \mathbf{J}_h is the unique solution of a square system.

The next result shows that this square system is well conditioned. Let $[\mathbf{J}_h]$ denote the vector of coefficients in the expansion of \mathbf{J}_h in the basis \mathfrak{B}^k . We place some minimal assumptions on the mesh from now on. As per standard terminology, we say that the mesh \mathcal{T}_h is shape regular if letting ρ_K be the diameter of the largest ball

contained in K , the ratios $\gamma_K = \text{diam}(K)/\rho_K$ are uniformly bounded by some fixed constant γ for all K . If we use the parameter α on every mesh face e to scale by the measure of the face, namely $|e|$, then we obtain a well conditioned matrix as stated in the following theorem.

THEOREM 3.2. *Let M be the matrix defined by*

$$[\mathbf{Y}_h]^t M[\mathbf{Z}_h] = \langle \alpha \mathbf{Z}_h, \mathbf{Y}_h \rangle_{\mathcal{E}_h \setminus \Gamma_N} \quad \text{for all } \mathbf{Y}_h \text{ and } \mathbf{Z}_h \in \mathfrak{J}_h.$$

Then, whenever

$$\alpha|_e = \frac{\zeta|e}{|e|} \quad \text{on } e \in \mathcal{E}_h \setminus \Gamma_N,$$

for some piecewise constant function ζ on $\mathcal{E}_h \setminus \Gamma_N$ satisfying $0 < \zeta_* \leq \zeta \leq \zeta^*$, the spectral condition number of M , namely κ_M , is uniformly bounded by

$$\kappa_M \leq C_0 \zeta^* / \zeta_*,$$

where $C_0 > 0$ is independent of the number of mesh elements (but depends on the space dimension N , and the polynomial degree k , and the shape regularity constant γ).

This theorem implies that to compute the solution \mathbf{J}_h of (3.7) by the method of conjugate gradients, we need a number of iterations that is independent of the number of unknowns. In Table 3.2, we numerically verify this fact for $k = 1, 2$ and 3

TABLE 3.2

Condition number of M with $\alpha \equiv 1$ for different mesh levels using different polynomial approximations

h	1	1/2	1/4	1/8	1/16	1/32
$k = 1$	3.4	5.7	6.8	7.3	7.5	7.5
$k = 2$	2.0	3.4	5.9	7.5	8.7	10.4
$k = 3$	2.1	5.8	5.7	8.6	11.8	10.9

on a sequence of uniform meshes. Since the meshes are uniform, we have simply taken $\alpha \equiv 1$ for this computation. Notice that, as expected, the condition numbers observed do not vary significantly as the mesh-size h is reduced. For practical computations, one often uses the method of conjugate gradients to solve for the Lagrange multiplier λ_h . Since, the system for λ_h has condition number $O(h^{-2})$ (cf. [24]) without any preconditioner, it is clear that the cost of computing \mathbf{J}_h is a negligible addition to the cost of solving for λ_h .

3.4. Error analysis. In this subsection, we give *a priori* error estimates for our new postprocessed flux approximation \mathbf{q}_h . Recall that \mathbf{q}_h is computed by the following steps:

1. Compute the CG solution U_h .
2. Using this U_h in (3.7), compute the unique function \mathbf{J}_h in \mathfrak{J}_h satisfying (3.7).

3. Set the flux trace $\widehat{\mathbf{q}}_h$ by substituting the \mathbf{J}_h computed above in (3.3).
4. Solve for \mathbf{q}_h element by element using (3.2), with the data set by the U_h and $\widehat{\mathbf{q}}_h$ computed above.

Notice that the last step involves equations very similar to the well known Raviart-Thomas projection defined as follows: We denote by $\pi_\ell \mathbf{q}$ the Raviart-Thomas projection [29] of the function \mathbf{q} , which is the unique function in

$$V_{RT_\ell}(K) := \mathcal{P}_\ell(K)^N + \mathbf{x} \mathcal{P}_\ell(K), \quad (3.12a)$$

satisfying

$$\langle \pi_\ell \mathbf{q} \cdot \mathbf{n}, v \rangle_e = \langle \mathbf{q} \cdot \mathbf{n}, v \rangle_e \quad \forall v \in \mathcal{P}_\ell(e), \text{ for any face } e \subset \partial K, \quad (3.12b)$$

$$(\pi_\ell \mathbf{q}, \mathbf{v})_K = (\mathbf{q}, \mathbf{v})_K \quad \forall \mathbf{v} \in \mathcal{P}_{\ell-1}(K)^N. \quad (3.12c)$$

It is well known [6] that the domain of definition of π_ℓ is slightly smaller than $\mathbf{H}(\text{div}, \Omega)$. We shall tacitly assume that the exact flux \mathbf{q} is smooth enough so that π_ℓ can be applied to it (e.g., \mathbf{q} in $\mathbf{H}(\text{div}, \Omega) \cap L^p(\Omega)$ with $p > 2$ is enough when $N = 2$.) Because of the similarity of (3.12) with (3.2), we shall refer to our flux approximation \mathbf{q}_h as the *RT $_\ell$ -postprocessed CG $_k$ flux*.

To describe our error estimates for this flux approximation, we need the following notation. We set $\text{diam}(K) = h_K$, and h to be the maximum of h_K over all K in \mathcal{T}_h . For Sobolev norms, we denote by $\|\cdot\|_{\ell, \mathcal{D}}$ and $|\cdot|_{\ell, \mathcal{D}}$ the H^ℓ -norm and seminorm, respectively, on \mathcal{D} . We also set

$$\mathbf{V}_{h, \ell}^0 := \{\mathbf{v} \in \mathbf{H}(\text{div}, \Omega) : \nabla \cdot \mathbf{v} = 0, \mathbf{v}|_K \in \mathcal{P}_\ell(K)^N \forall K \in \mathcal{T}_h, \mathbf{v} \cdot \mathbf{n}|_{\Gamma_N} = 0\}, \quad (3.13a)$$

and denote by \mathbb{P}_ℓ the weighted L^2 -projection into $\mathbf{V}_{h, \ell}^0$ defined by

$$(a^{-1}(\mathbb{P}_\ell \mathbf{q} - \mathbf{q}), \mathbf{v})_\Omega = 0 \quad \forall \mathbf{v} \in \mathbf{V}_{h, \ell}^0. \quad (3.13b)$$

Finally, we set $\alpha^* := \max_{e \in \mathcal{E}_h \setminus \Gamma_N} \alpha|_e$ and $\alpha_* := \min_{e \in \mathcal{E}_h \setminus \Gamma_N} \alpha|_e$. With these notations, we have the following result.

THEOREM 3.3. *Let \mathbf{q}_h be the RT $_\ell$ -postprocessed CG $_k$ flux for an integer $0 \leq \ell \leq k$. Then the following statements hold:*

1. *The RT $_\ell$ -postprocessed CG $_k$ flux \mathbf{q}_h is in $\mathbf{H}(\text{div}, \Omega)$ and satisfies*

$$\text{div}(\pi_\ell \mathbf{q} - \mathbf{q}_h) = 0.$$

In particular, it satisfies the exact conservativity property (1.4).

2. *If $\ell > 1$ and $a(\mathbf{x})$ is constant on each mesh element, $\mathbb{P}_{\ell-1}(\pi_\ell \mathbf{q} - \mathbf{q}_h) = 0$.*
3. *The divergence of the flux approximation satisfies*

$$\|\nabla \cdot (\mathbf{q} - \mathbf{q}_h)\|_{0, \Omega} \leq C_1 h^{\min(\ell, s)+1} |f|_{s+1, \Omega}.$$

4. *If $a(\mathbf{x})|_K$ is in $W^{1, \infty}(K)$ for all mesh elements K , and the mesh \mathcal{T}_h is quasi-uniform, then the following error estimate holds:*

$$\|\mathbf{q} - \mathbf{q}_h\|_{0, \Omega} \leq C_2 h^{\min(k, \ell+1, s)} (|\mathbf{q}|_{s, \Omega} + |u|_{s, \Omega}).$$

In the inequalities above, C_1 and $C_2/(1 + \alpha^*/\alpha_*)$ are independent of \mathbf{q} and h (but dependent on k, N, β, a , and γ).

The first identity of the theorem can be interpreted as a superconvergence property for the divergence. Indeed, if the load f is a piecewise polynomial satisfying $f|_K \in \mathcal{P}_\ell(K)$, then the exact and discrete divergences coincide, i.e., $\operatorname{div}(\mathbf{q} - \mathbf{q}_h) = 0$, because of a well known commutativity property of π_ℓ . In one space dimension, it states that the difference between $\pi_\ell \mathbf{q}$ and \mathbf{q}_h is just a constant; moreover, if the Neumann boundary is not empty, $\pi_\ell \mathbf{q}$ and \mathbf{q}_h are *identical*. This implies that, at each node x_i we have that

$$\mathbf{q}_h(x_i) = \mathbf{q}(x_i),$$

by definition of the projection π_ℓ . The fact that this holds independently of how we chose the parameter β in (3.5) is remarkable, although this fits very well with similar results obtained in [12]. Notice also that when the Neumann boundary is empty but a is piecewise constant, the second identity states that $\pi_\ell \mathbf{q}$ and \mathbf{q}_h are also identical, provided $\ell \geq 1$.

Next, we compare our approximation \mathbf{q}_h to the corresponding Raviart-Thomas approximation. We begin by recalling the standard result that *all statements of Theorem 3.3 continue to hold if we replace \mathbf{q} by the approximation to the flux given by the RT_ℓ method*, i.e., both the standard RT_{k-1} method and our new RT_{k-1} -postprocessed CG_k method produce $H(\operatorname{div})$ -conforming approximations to the flux \mathbf{q} that converge at the same order. This indicates that the RT_{k-1} -postprocessed CG_k method is competitive with the RT_{k-1} method. Indeed, to compare their computational complexities, we recall the earlier observation that the cost of the computation of our approximation \mathbf{q}_h is negligible compared to that of solving for the Lagrange multiplier λ_h . The condition number of the Lagrange multiplier system in the CG case as well as the RT case [24] is $O(h^{-2})$, so in both cases, the cost of solving for λ_h dominates the cost of the computation of \mathbf{q}_h . Thus the relative size of the stiffness matrices for the Lagrange multiplier becomes the deciding factor.

It is not difficult to see that this matrix for the CG_k method has smaller size than that of the RT_{k-1} method. Let us show this in the case of a two dimensional simply connected domain Ω . We denote the number of mesh vertices, edges, and triangles, by n_e, n_v and n_t , respectively. The number of degrees of freedom of the Lagrange multipliers for the CG_k method is $(k-1)n_e + n_v$ whereas it is kn_e for the RT_{k-1} method. Since $n_v - n_e + n_t = 1$, we see that the Lagrange multipliers of the CG_k method have $(n_t - 1)$ fewer degrees of freedom than RT_{k-1} . This is a significant difference in practice. The numerical experiments of Section 5 show that, for the same mesh, the approximations given by the RT_{k-1} -postprocessed CG_k method and the RT_{k-1} method are very similar. This shows that the former method may be better than the latter. A final point reinforcing this conclusion is obtained by comparing the approximation to u given by both the methods. The u_h of the CG_k method converges in the L^2 -norm with order $k+1$ when the exact solution is smooth. However, the approximation to u given by the RT_{k-1} method converges only with order k . Of course, following [1], we can use the Lagrange multipliers to obtain a locally postprocessed approximation that also converges with order $k+1$, but such postprocessing is not available for arbitrary values of k . Moreover, our numerical results show that for $k \in \{0, 1, 2\}$, the CG_k and the postprocessed RT_{k-1} methods produce roughly similar approximations to u .

4. Proofs.

4.1. Norm equivalences. In this subsection, we will prove the condition number estimate of Theorem 3.2 using certain norm equivalences. Recall that the extension operator T_h is defined in (3.10) and that $[\mathbf{J}_h]$ is the vector representation of the function \mathbf{J}_h in the basis \mathfrak{B}^k . Define the following norms

$$\|\mathbf{J}_h\|_S = \left(\sum_{e \in \mathcal{E}_h \setminus \Gamma_N} \frac{1}{|e|} \|\mathbf{J}_h\|_{0,e}^2 \right)^{1/2} \quad \text{and} \quad \|\mathbf{J}_h\|_T = \left(\sum_{K \in \mathcal{T}_h} \frac{1}{|K|} \|T(\mathbf{J}_h)\|_{0,K}^2 \right)^{1/2},$$

where $|K|$ and $|e|$ denote the measures (in their respective dimensions) of an element K and face e , respectively. The following lemma shows that the three norms $\|\mathbf{J}_h\|_S$, $\|\mathbf{J}_h\|_T$, and $\|[\mathbf{J}_h]\|_{\ell^2}^2$ are equivalent.

LEMMA 4.1. *There is a constant C independent of the mesh size (but depending on the degree k , dimension N , and the shape regularity constant γ) such that*

$$\frac{1}{C} \|\mathbf{J}_h\|_T^2 \leq \|\mathbf{J}_h\|_S^2 \leq C \|\mathbf{J}_h\|_T^2, \quad (4.1)$$

and

$$\frac{1}{C} \|\mathbf{J}_h\|_S^2 \leq \|[\mathbf{J}_h]\|_{\ell^2}^2 \leq C \|\mathbf{J}_h\|_S^2, \quad (4.2)$$

for all $\mathbf{J}_h \in \mathfrak{J}_h$.

Proof. First, let us prove the upper bound of (4.1). Recall that by (3.10a), \mathbf{J}_h and $[T_h(\mathbf{J}_h) \mathbf{n}]$ coincide for any \mathbf{J}_h in \mathfrak{J}_h , so by standard trace inequalities,

$$\frac{1}{|e|} \|\mathbf{J}_h\|_{0,e}^2 \leq C \sum_{K \in \mathcal{K}_e} \frac{1}{|K|} \|T_h(\mathbf{J}_h)\|_{0,K}^2, \quad (4.3)$$

where \mathcal{K}_e denotes the set of elements $K \in \mathcal{T}_h$ such that e is a face of K . Since \mathcal{K}_e has at most two elements for any mesh face e , summing over all edges in $\mathcal{E}_h \setminus \Gamma_N$, we obtain the upper bound in (4.1).

Next, let us prove the lower bound of (4.1). By standard scaling arguments using the principal lattice \mathcal{A}_N^k on any mesh element K , we have

$$\begin{aligned} \frac{C}{|K|} \|T_h(\mathbf{J}_h)\|_{0,K}^2 &\leq \sum_{\alpha \in \mathcal{A}_N^k} |(T_h(\mathbf{J}_h))(\mathbf{x}_{\alpha,K})|^2 \\ &= \sum_{\mathbf{z} \in \mathfrak{S}_h^k \cap \partial K} T_h(\mathbf{J}_h)|_{K,\mathbf{z}}^2, \end{aligned} \quad (4.4)$$

where, as before, $w|_{K,\mathbf{z}}$ denote the limit of $w(\mathbf{x})$ as \mathbf{x} approaches \mathbf{z} from within K .

We need to bound each of the terms in (4.4) using norms of \mathbf{J}_h . Let us first consider the case when \mathbf{z} is not on Γ_D . For such a \mathbf{z} , recalling the way we constructed the basis of Theorem 3.1, note that there is a mesh element K^* such that $\mathbf{z} = \mathbf{x}_{\alpha^*,K^*}$ where the limit $T_h(\mathbf{J}_h)|_{K^*,\mathbf{z}}$ is zero – see (3.10b). Using this fact, it is easy to see that we can write $T_h(\mathbf{J}_h)|_{K,\mathbf{z}}$ as the telescoping sum

$$T_h(\mathbf{J}_h)|_{K,\mathbf{z}} = \sum_{i=1}^m \left[T_h(\mathbf{J}_h)|_{K_i,\mathbf{z}} - T_h(\mathbf{J}_h)|_{K_{i+1},\mathbf{z}} \right], \quad (4.5)$$

for some collection of mesh elements K_i such that \mathbf{z} is in \bar{K}_i , $K_1 = K$, $K_{m+1} = K^*$, and $K_i \cap K_{i+1}$ is a mesh face in $\mathcal{E}_h \setminus \Gamma_N$. If \mathbf{z} lies on Γ_D , we can still write a similar sum as long as we omit the last term (as there is no K^* for such \mathbf{z}) and choose K_m such that it has a face on $\partial\Omega$. By (3.10a), the absolute value of the i^{th} summand inside the square brackets in (4.5) equals the magnitude of the limit of \mathbf{J}_h as we approach \mathbf{z} from within the mesh face $K_i \cap K_{i+1}$. Expressing each of the terms in the sum in (4.4) in terms of \mathbf{J}_h this way, we obtain

$$\frac{C}{|K|} \|T_h(\mathbf{J}_h)\|_{0,K}^2 \leq \sum_{\mathbf{z} \in \mathcal{G}_h^k \cap \partial K} \sum_{e \in \mathcal{F}_{\mathbf{z}}} \frac{1}{|e|} \|\mathbf{J}_h\|_{0,e}^2, \quad (4.6)$$

where $\mathcal{F}_{\mathbf{z}}$ denotes the set of all mesh faces e in $\mathcal{E}_h \setminus \Gamma_N$ such that $\mathbf{z} \in \bar{e}$. Note that in obtaining the above inequality, we have used the fact that for every \mathbf{z} , the number m in (4.5) can be bounded uniformly in terms of the shape regularity constants. Summing over all mesh elements K , we obtain the lower bound of (4.1).

It now remains to prove (4.2). Recall that a standard norm equivalence asserts the existence of a constant C (depending on the shape regularity of K , but otherwise independent of K) such that for all $w \in \mathcal{P}_k(K)$,

$$\frac{1}{C|K|} \|w\|_{0,K}^2 \leq \sum_{\alpha \in \mathcal{A}_N^k} |w(\mathbf{x}_{\alpha,K})|^2 \leq \frac{C}{|K|} \|w\|_{0,K}^2.$$

Applying this with $w = T_h(\mathbf{J}_h)|_K$, and observing that in the expansion for $T_h(\mathbf{J}_h)$ in (3.11), the coefficients $\{c_\phi\}$ are the nonzero values of $T_h(\mathbf{J}_h)$ at the points $\mathbf{x}_{\alpha,K}$, we obtain

$$\frac{1}{C|K|} \|T_h(\mathbf{J}_h)\|_{0,K}^2 \leq \sum_{\mathbf{z} \in \mathcal{G}_h^k \cap \partial K} \sum_{\phi \in V_{\mathbf{z}}^*} c_\phi^2 \leq \frac{C}{|K|} \|T_h(\mathbf{J}_h)\|_{0,K}^2. \quad (4.7)$$

The upper inequality above implies

$$\begin{aligned} \sum_{\mathbf{z} \in \mathcal{G}_h^k \cap \partial K} \sum_{\phi \in V_{\mathbf{z}}^*} c_\phi^2 &\leq \frac{C}{|K|} \|T_h(\mathbf{J}_h)\|_{0,K}^2, && \text{by (4.7),} \\ &\leq C \sum_{\mathbf{z} \in \mathcal{G}_h^k \cap \partial K} \sum_{e \in \mathcal{F}_{\mathbf{z}}} \frac{1}{|e|} \|\mathbf{J}_h\|_{0,e}^2, && \text{by (4.6).} \end{aligned}$$

If we sum this inequality over all mesh elements K , the resulting left hand side dominates $\|[\mathbf{J}_h]\|_{\ell^2}^2$. Hence we have proven that

$$\|[\mathbf{J}_h]\|_{\ell^2}^2 \leq C \sum_{e \in \mathcal{E}_h \setminus \Gamma_N} \frac{1}{|e|} \|\mathbf{J}_h\|_{0,e}^2. \quad (4.8)$$

Returning to (4.7) and using its lower inequality, we also have

$$\begin{aligned} \frac{1}{|e|} \|\mathbf{J}_h\|_{0,e}^2 &\leq C \sum_{K \in \mathcal{K}_e} \frac{1}{|K|} \|T_h(\mathbf{J}_h)\|_{0,K}^2, && \text{by (4.3),} \\ &\leq C \sum_{K \in \mathcal{K}_e} \sum_{\mathbf{z} \in \mathcal{G}_h^k \cap \partial K} \sum_{\phi \in V_{\mathbf{z}}^*} c_\phi^2, && \text{by (4.7).} \end{aligned}$$

Summing this inequality over all edges e in $\mathcal{E}_h \setminus \Gamma_N$ and noting that the resulting number of repetitions in \mathcal{C}_ϕ^2 can be uniformly bounded, we obtain

$$\sum_{e \in \mathcal{E}_h \setminus \Gamma_N} \frac{1}{|e|} \|\mathbf{J}_h\|_{0,e}^2 \leq C \|\mathbf{J}_h\|_{\ell^2}^2. \quad (4.9)$$

Combining (4.9) and (4.8), the proof of (4.2) is finished. \square

Proof of Theorem 3.2. Since the matrix M is symmetric and positive definite, its spectral condition number κ_M is given by

$$\kappa_M = \left(\max_{\mathbf{Y}_h \in \mathfrak{J}_h} \frac{\langle \alpha \mathbf{Y}_h, \mathbf{Y}_h \rangle_{\mathcal{E}_h \setminus \Gamma_N}}{\|\mathbf{Y}_h\|_{\ell^2}^2} \right) / \left(\min_{\mathbf{Y}_h \in \mathfrak{J}_h} \frac{\langle \alpha \mathbf{Y}_h, \mathbf{Y}_h \rangle_{\mathcal{E}_h \setminus \Gamma_N}}{\|\mathbf{Y}_h\|_{\ell^2}^2} \right). \quad (4.10)$$

By the assumptions in the theorem on α ,

$$\zeta_\star \left(\sum_{e \in \mathcal{E}_h \setminus \Gamma_N} \frac{1}{|e|} \|\mathbf{Y}_h\|_{0,e}^2 \right) \leq \langle \alpha \mathbf{Y}_h, \mathbf{Y}_h \rangle_{\mathcal{E}_h \setminus \Gamma_N} \leq \zeta^\star \left(\sum_{e \in \mathcal{E}_h \setminus \Gamma_N} \frac{1}{|e|} \|\mathbf{Y}_h\|_{0,e}^2 \right).$$

Applying (4.2) of Lemma 4.1 to the above inequality, we obtain

$$\frac{\zeta_\star}{C} \|\mathbf{Y}_h\|^2 \leq \langle \alpha \mathbf{Y}_h, \mathbf{Y}_h \rangle_{\mathcal{E}_h \setminus \Gamma_N} \leq C \zeta^\star \|\mathbf{Y}_h\|^2.$$

Using this in (4.10), we find that $\kappa_M \leq C^2 \zeta^\star / \zeta_\star$. \square

4.2. Proof of the flux error estimates. This subsection is devoted to proving the error estimates of Theorem 3.3. The error in the divergence is easy to analyze, but the proof of the L^2 -estimate is more involved. Proceeding as in the analysis of [16] of the hybridized Raviart-Thomas method, we start with the error equations.

Proof of Theorem 3.3. We divide this proof into seven steps.

Step 1. *Obtaining the error equations:* If we set $\mathbf{q}_{\nabla,h} := -a \nabla U_h$, on each element $K \in \mathcal{T}_h$, from the equations (2.2a), and (2.2b) defining the hybridized continuous Galerkin method and from the equation (3.1) relating $q_{n,h}$ with the numerical trace $\widehat{\mathbf{q}}_h$, it follows that

$$\begin{aligned} (a^{-1} \mathbf{q}_{\nabla,h}, \mathbf{v})_K - (U_h, \nabla \cdot \mathbf{v})_K &= -\langle \lambda_h, \mathbf{v} \cdot \mathbf{n} \rangle_{\partial K}, \\ -(\mathbf{q}_{\nabla,h}, \nabla w)_K + \langle \widehat{\mathbf{q}}_h, w \mathbf{n} \rangle_{\partial K} &= (f, w)_K, \end{aligned}$$

for any $(\mathbf{v}, w) \in \mathcal{P}_k(K)^N \times \mathcal{P}_k(K)$. As a consequence, by the definition of \mathbf{q}_h given by equations (3.2), we obtain that, for $\ell \leq k$,

$$(a^{-1} \mathbf{q}_h, \mathbf{v})_K - (U_h, \nabla \cdot \mathbf{v})_K = -\langle \lambda_h, \mathbf{v} \cdot \mathbf{n} \rangle_{\partial K} + (a^{-1} (\mathbf{q}_h - \mathbf{q}_{\nabla,h}), \mathbf{v})_K, \quad (4.11a)$$

$$-(\mathbf{q}_h, \nabla w)_K + \langle \mathbf{q}_h, w \mathbf{n} \rangle_{\partial K} = (f, w)_K, \quad (4.11b)$$

for any $(\mathbf{v}, w) \in \mathcal{P}_k(K)^N \times \mathcal{P}_\ell(K)$. The error equations are derived by comparing these equations to the equations satisfied by the exact solution (\mathbf{q}, u) , namely

$$\begin{aligned} (a^{-1} \mathbf{q}, \mathbf{v})_K - (u, \nabla \cdot \mathbf{v})_K &= -\langle u, \mathbf{v} \cdot \mathbf{n} \rangle_{\partial K}, \\ -(\mathbf{q}, \nabla w)_K + \langle \mathbf{q}, w \mathbf{n} \rangle_{\partial K} &= (f, w)_K, \end{aligned}$$

for any $(\mathbf{v}, w) \in V_{RT_\ell}(K) \times \mathcal{P}_\ell(K)$. They imply, as a consequence of the definition of the Raviart-Thomas projection π_ℓ given in (3.12), that

$$(a^{-1}\pi_\ell \mathbf{q}, \mathbf{v})_K - (u, \nabla \cdot \mathbf{v})_K = -\langle u, \mathbf{v} \cdot \mathbf{n} \rangle_{\partial K} + (a^{-1}(\pi_\ell \mathbf{q} - \mathbf{q}), \mathbf{v})_K, \quad (4.12a)$$

$$-(\pi_\ell \mathbf{q}, \nabla w)_K + \langle \pi_\ell \mathbf{q}, w \mathbf{n} \rangle_{\partial K} = (f, w)_K, \quad (4.12b)$$

for any $(\mathbf{v}, w) \in V_{RT_\ell}(K) \times \mathcal{P}_\ell(K)$.

Thus, if we define the errors of the approximation as

$$\mathbf{e}_q = \pi_\ell \mathbf{q} - \mathbf{q}_h, \quad e_u = u - U_h, \quad e_\lambda = u - \lambda_h,$$

we see, after subtracting the equations (4.11) from the equations (4.12), that they satisfy

$$\begin{aligned} (a^{-1}\mathbf{e}_q, \mathbf{v})_K - (e_u, \nabla \cdot \mathbf{v})_K &= -\langle e_\lambda, \mathbf{v} \cdot \mathbf{n} \rangle_{\partial K} \\ &\quad - (a^{-1}(\mathbf{q}_h - \mathbf{q}_{\nabla, h}), \mathbf{v})_K \\ &\quad + (a^{-1}(\pi_\ell \mathbf{q} - \mathbf{q}), \mathbf{v})_K, \end{aligned} \quad (4.13a)$$

$$-(\mathbf{e}_q, \nabla w)_K + \langle \mathbf{e}_q \cdot \mathbf{n}, w \rangle_{\partial K} = 0. \quad (4.13b)$$

for any $(\mathbf{v}, w) \in (V_{RT_\ell}(K) \cap \mathcal{P}_k(K)^N) \times \mathcal{P}_\ell(K)$ for $\ell \leq k$.

Step 2. *Analyzing errors in the divergence of the flux:* Integrating (4.13b) by parts, we obtain

$$(\nabla \cdot \mathbf{e}_q, w)_K = 0 \quad \forall w \in \mathcal{P}_\ell(K).$$

Since $\nabla \cdot \mathbf{e}_q \in \mathcal{P}_\ell(K)$, we immediately get that

$$\nabla \cdot \mathbf{e}_q \equiv 0 \quad \text{on } K,$$

which is the first identity of Theorem 3.3. (It is obvious that \mathbf{e}_q is in $\mathbf{H}(\text{div}, \Omega)$.)

The first inequality of Theorem 3.3 from the fact that

$$\nabla \cdot (\mathbf{q} - \mathbf{q}_h) = \nabla \cdot (\mathbf{q} - \pi_\ell \mathbf{q}) = (Id - P_\ell) f,$$

where P_ℓ is the L^2 -projection into the space of functions w such that $w_K \in \mathcal{P}_\ell(K)$ for all $K \in \mathcal{T}_h$. Notice that in the last step, we used the commutativity property $\nabla \cdot \pi_\ell = P_\ell \nabla \cdot$ (see e.g. [6, 16, 22]).

Step 3. *Establishing the second identity:* If in the error equation (4.13a), we select $\mathbf{v} \in \mathcal{P}_{\ell-1}^N(K)$, we find that whenever $a(\mathbf{x})$ is constant on K ,

$$(a^{-1}\mathbf{e}_q, \mathbf{v})_K - (e_u, \nabla \cdot \mathbf{v})_K = -\langle e_\lambda, \mathbf{v} \cdot \mathbf{n} \rangle_{\partial K},$$

where we used the equation (3.2c) of the definition of \mathbf{q}_h and the equation (3.12c) of the definition of π_ℓ . This readily implies that

$$(a^{-1}\mathbf{e}_q, \mathbf{v})_\Omega - \sum_{K \in \mathcal{T}_h} (e_u, \nabla \cdot \mathbf{v})_K = - \sum_{e \in \mathcal{E}_h} \langle e_\lambda, \llbracket \mathbf{v} \cdot \mathbf{n} \rrbracket \rangle_e,$$

and so

$$(a^{-1}\mathbf{e}_q, \mathbf{v})_\Omega = -\langle e_\lambda, \mathbf{v} \cdot \mathbf{n} \rangle_{\Gamma_D} = 0 \quad \forall \mathbf{v} \in \mathbf{V}_{h, \ell-1}^0,$$

since we are assuming that $g|_e \in \mathcal{P}_k(e)$ for each face e on Γ_D . The second identity of Theorem 3.3 immediately follows from this and the definition of the projection $\mathbb{P}_{\ell-1}$ given by (3.13).

Step 4. *Splitting errors in the flux:* It remains to prove the second inequality of Theorem 3.3. To do this, we begin by noting that since $\mathbf{e}_q \in \mathbf{V}_{h,\ell}^0$, we can choose $\mathbf{v} = \mathbf{e}_q$ in the error equation (4.13a). Doing this and summing over all the elements $K \in \mathcal{T}_h$, we obtain

$$(a^{-1}e_q, e_q)_\Omega = - (a^{-1}(\mathbf{q}_h - \mathbf{q}_{\nabla,h}), \mathbf{e}_q)_\Omega + (a^{-1}(\pi_\ell \mathbf{q} - \mathbf{q}), \mathbf{e}_q)_\Omega.$$

Introducing $\pi_\ell \mathbf{q}_{\nabla,h}$ into the right hand side,

$$\begin{aligned} (a^{-1}e_q, e_q)_\Omega &= - (a^{-1}(\mathbf{q}_h - \pi_\ell \mathbf{q}_{\nabla,h}), \mathbf{e}_q)_\Omega - (a^{-1}(\pi_\ell \mathbf{q}_{\nabla,h} - \mathbf{q}_{\nabla,h}), \mathbf{e}_q)_\Omega \\ &\quad + (a^{-1}(\pi_\ell \mathbf{q} - \mathbf{q}), \mathbf{e}_q)_\Omega \\ &= - (a^{-1}(\mathbf{q}_h - \pi_\ell \mathbf{q}_{\nabla,h}), \mathbf{e}_q)_\Omega - (a^{-1}(Id - \pi_\ell)(\mathbf{q} - \mathbf{q}_{\nabla,h}), \mathbf{e}_q)_\Omega. \end{aligned}$$

Applying the Cauchy-Schwarz inequality,

$$C \|e_q\|_{0,\Omega} \leq \|\mathbf{q}_h - \pi_\ell \mathbf{q}_{\nabla,h}\|_{0,\Omega} + \|(Id - \pi_\ell)(\mathbf{q} - \mathbf{q}_{\nabla,h})\|_{0,\Omega}. \quad (4.14)$$

We now estimate each of the terms on the right hand side separately in the next two steps.

Step 5. *Estimating $\|\mathbf{q}_h - \pi_\ell \mathbf{q}_{\nabla,h}\|_{0,\Omega}$:* In order to estimate this term, we rewrite the equations (3.2b) and (3.2c) defining \mathbf{q}_h as

$$\begin{aligned} (\mathbf{q}_h - \pi_\ell \mathbf{q}_{\nabla,h}, \mathbf{v})_K &= 0, & \forall \mathbf{v} \in \mathcal{P}_{\ell-1}(K)^N, \\ \langle (\mathbf{q}_h - \pi_\ell \mathbf{q}_{\nabla,h}) \cdot \mathbf{n}, w \rangle_e &= \langle (\widehat{\mathbf{q}}_h - \mathbf{q}_{\nabla,h}) \cdot \mathbf{n}, w \rangle_e, & \forall w \in \mathcal{P}_\ell(e) \text{ and all faces } e \subset \partial K. \end{aligned}$$

Then, a standard scaling argument gives

$$\|\mathbf{q}_h - \pi_\ell \mathbf{q}_{\nabla,h}\|_{0,K}^2 \leq Ch_K \|(\widehat{\mathbf{q}}_h - \mathbf{q}_{\nabla,h}) \cdot \mathbf{n}\|_{0,\partial K}^2.$$

Summing over all mesh elements and using the definition of the numerical trace $\widehat{\mathbf{q}}_h$ given by equation (3.5),

$$\begin{aligned} \|\mathbf{q}_h - \pi_\ell \mathbf{q}_{\nabla,h}\|_{0,\Omega}^2 &\leq C \sum_{K \in \mathcal{T}_h} h_K \left(\left\| \left(\beta \cdot \mathbf{n} - \frac{1}{2} \right) [\mathbf{q}_{\nabla,h} \cdot \mathbf{n}] + \alpha \mathbf{J}_h \cdot \mathbf{n} \right\|_{0,\partial K \setminus \partial \Omega}^2 \right. \\ &\quad \left. + \|\alpha \mathbf{J}_h \cdot \mathbf{n}\|_{0,\partial K \cap \Gamma_D}^2 + \|\mathbf{q}_N - \mathbf{q}_{\nabla,h} \cdot \mathbf{n}\|_{0,\partial K \cap \Gamma_N}^2 \right) \\ &\leq Ch (T_1 + T_2), \end{aligned} \quad (4.15)$$

where

$$T_1 := \left\| [\mathbf{q}_{\nabla,h} \cdot \mathbf{n} - \pi_\ell \mathbf{q} \cdot \mathbf{n}] \right\|_{0,\mathcal{E}_h \setminus \Gamma_D}^2 \quad \text{and} \quad T_2 := \|\alpha \mathbf{J}_h \cdot \mathbf{n}\|_{0,\mathcal{E}_h \setminus \Gamma_N}^2.$$

The term T_1 can be easily estimated by an inverse inequality:

$$T_1 \leq Ch^{-1} \|\mathbf{q}_{\nabla,h} - \pi_\ell \mathbf{q}\|_{0,\Omega}^2 \leq Ch^{-1} (|u - U_h|_{1,\Omega}^2 + \|\mathbf{q} - \pi_\ell \mathbf{q}\|_{0,\Omega}^2). \quad (4.16)$$

The other term T_2 requires more work.

To estimate T_2 , we rewrite the definition of the jump \mathbf{J}_h , namely (3.8), as

$$\begin{aligned} \langle \alpha \mathbf{J}_h, \llbracket v \mathbf{n} \rrbracket \rangle_{\mathcal{E}_h \setminus \Gamma_N} &= \langle \llbracket \mathbf{q}_{\nabla, h} \cdot \mathbf{n} \rrbracket, \{\!\!\{ v \}\!\!\} - \beta \cdot \llbracket v \mathbf{n} \rrbracket \rangle_{\mathcal{E}_h \setminus \partial\Omega} \\ &\quad + \sum_{K \in \mathcal{T}_h} (f - \nabla \cdot \mathbf{q}_{\nabla, h}, v)_K \\ &= \langle \llbracket \mathbf{q}_{\nabla, h} \cdot \mathbf{n} \rrbracket, \{\!\!\{ v \}\!\!\} - \beta \cdot \llbracket v \mathbf{n} \rrbracket \rangle_{\mathcal{E}_h \setminus \partial\Omega} \\ &\quad + \sum_{K \in \mathcal{T}_h} (\nabla \cdot (\pi_k \mathbf{q} - \mathbf{q}_{\nabla, h}), v)_K, \end{aligned}$$

for all $v \in V_h$. Choosing $v = T_h(\mathbf{J}_h)$ and using the property (3.10a) of the operator T_h , we get

$$\begin{aligned} \langle \alpha \mathbf{J}_h, \mathbf{J}_h \rangle_{\mathcal{E}_h \setminus \Gamma_N} &= \langle \llbracket \mathbf{q}_{\nabla, h} \cdot \mathbf{n} \rrbracket, \{\!\!\{ T_h(\mathbf{J}_h) \}\!\!\} \rangle_{\mathcal{E}_h^\circ} - \langle \llbracket \mathbf{q}_{\nabla, h} \cdot \mathbf{n} \rrbracket, \beta \cdot \mathbf{J}_h \rangle_{\mathcal{E}_h^\circ} \\ &\quad + \sum_{K \in \mathcal{T}_h} (\nabla \cdot (\pi_k \mathbf{q} - \mathbf{q}_{\nabla, h}), T_h(\mathbf{J}_h))_K. \end{aligned}$$

Using (4.1) of Lemma 4.1 after applying suitable inverse inequalities, we obtain

$$\langle \alpha \mathbf{J}_h, \mathbf{J}_h \rangle_{\mathcal{E}_h \setminus \Gamma_N} \leq C \left(\|\llbracket \mathbf{q}_{\nabla, h} \cdot \mathbf{n} \rrbracket\|_{0, \mathcal{E}_h^\circ} + h^{1/2} \|\nabla \cdot (\pi_k \mathbf{q} - \mathbf{q}_{\nabla, h})\|_{0, \Omega} \right) \|\mathbf{J}_h\|_{0, \mathcal{E}_h \setminus \Gamma_N}.$$

This implies that

$$T_2 \leq C \frac{\alpha^*}{\alpha_*} \left(\|\llbracket \mathbf{q}_{\nabla, h} \cdot \mathbf{n} \rrbracket\|_{0, \mathcal{E}_h^\circ}^2 + h \|\nabla \cdot (\pi_k \mathbf{q} - \mathbf{q}_{\nabla, h})\|_{0, \Omega}^2 \right).$$

Treating the first term inside the parentheses above as in the proof of (4.16), and applying an inverse inequality to the second, we get

$$T_2 \leq C h^{-1} (|u - U_h|_{1, \Omega}^2 + \|\mathbf{q} - \pi_\ell \mathbf{q}\|_{0, \Omega}^2 + \|\mathbf{q} - \pi_k \mathbf{q}\|_{0, \Omega}^2).$$

Using the estimates for T_1 and T_2 in (4.15), we conclude that

$$C \|\mathbf{q}_h - \pi_\ell \mathbf{q}_{\nabla, h}\|_{0, \Omega}^2 \leq |u - U_h|_{1, \Omega}^2 + \|\mathbf{q} - \pi_\ell \mathbf{q}\|_{0, \Omega}^2 + \|\mathbf{q} - \pi_k \mathbf{q}\|_{0, \Omega}^2. \quad (4.17)$$

Step 6. *Estimating* $\|(Id - \pi_\ell)(\mathbf{q} - \mathbf{q}_{\nabla, h})\|_{0, \Omega}$: On an element K , using the well known approximation property of π_ℓ [6, 29]

$$\begin{aligned} \|(Id - \pi_\ell)(\mathbf{q} - \mathbf{q}_{\nabla, h})\|_{0, K} &\leq Ch_K |\mathbf{q} - \mathbf{q}_{\nabla, h}|_{1, K} \\ &\leq Ch_K \|a\|_{W_\infty^1(K)} \|\nabla u - \nabla U_h\|_{1, K}. \end{aligned}$$

Now using any projector Π_K well defined on $L^2(K)^N$ with standard approximation properties, e.g. the one constructed in [30], we have

$$\begin{aligned} \|(Id - \pi_\ell)(\mathbf{q} - \mathbf{q}_{\nabla, h})\|_{0, K} &\leq Ch_K (\|\nabla u - \Pi_K \nabla u\|_{1, K} + \|\Pi_K(\nabla u - \nabla U_h)\|_{1, K}) \\ &\leq C(h_K \|\nabla u - \Pi_K \nabla u\|_{1, K} + \|\Pi_K(\nabla u - \nabla U_h)\|_{0, K}) \\ &\leq C(h_K \|\nabla u - \Pi_K \nabla u\|_{1, K} + \|\nabla u - \nabla U_h\|_{0, K}). \end{aligned}$$

Thus, we obtain

$$\|(Id - \pi_\ell)(\mathbf{q} - \mathbf{q}_{\nabla, h})\|_{0, \Omega} \leq Ch^{\min(s, k)} |u|_{s+1, \Omega}. \quad (4.18)$$

Step 7. *Completing the proof of Theorem 3.3:* Now we use the results of the previous two steps, namely (4.17) and (4.18), in the splitting (4.14) of the error term. Then we obtain

$$C\|\mathbf{e}_q\|_{0,\Omega} \leq |u - U_h|_{1,\Omega} + \|\mathbf{q} - \pi_\ell \mathbf{q}\|_{0,\Omega} + \|\mathbf{q} - \pi_k \mathbf{q}\|_{0,\Omega} + h^{\min(s,k)}|u|_{s+1,\Omega}.$$

and the estimate of \mathbf{q}_h immediately follows from the standard approximation results

$$\begin{aligned} |u - U_h|_{1,\Omega} &\leq Ch^{\min\{k,s\}}|u|_{s+1,\Omega}, \\ \|\mathbf{q} - \pi_m \mathbf{q}\|_{0,\Omega} &\leq Ch^{\min\{m+1,s\}}|\mathbf{q}|_{s,\Omega}, \end{aligned}$$

This concludes the proof. \square

5. Numerical results. In this section, we carry out some numerical experiments to verify the theoretical results when the exact solution is smooth (Test 1) and to test the performance of the method when the exact solution has a singularity (Test 2). For the sake of simplicity, we use uniform meshes and pick $\beta = \mathbf{0}$ in the definition of the numerical trace $\widehat{\mathbf{q}}_h$, (3.5).

In what follows, by the approximation given by the “RT $_\ell$ method” we mean the pair (\mathbf{q}_h, U_h) obtained as follows. The function $(\mathbf{q}_h, u_h, \lambda_h)$ is the solution of the hybridized Raviart-Thomas method whose Lagrange multipliers are piecewise polynomials of degree ℓ . The function U_h is obtained from (u_h, λ_h) by using the local postprocessing described in [1]. The resulting pair (\mathbf{q}_h, U_h) is then compared to the solution of our RT $_\ell$ -postprocessed CG $_k$ method, for which \mathbf{q}_h is the RT $_\ell$ -postprocessed CG $_k$ flux and U_h is the solution of the CG method with piecewise polynomials of degree k .

Test 1. We take

$$a = \begin{pmatrix} x+2 & x+y \\ x+y & y+2 \end{pmatrix}$$

and then g and f so that the exact solution is

$$u(x, y) = \sin(\pi x) \sin(\pi y).$$

The history of convergence of the approximations given by the RT $_{k-1}$ and the RT $_\ell$ -postprocessed CG $_k$ methods, for $\ell \in \{k-1, k\}$, are displayed in Tables 5.1, 5.2 and 5.3 for $k = 1, 2$, and 3, respectively. Plots of these results are also displayed in Fig. 5.1 for an easier comparison.

We see that the approximation given by the RT $_\ell$ -postprocessed CG $_k$ method converges with the orders predicted by Theorem 3.3. Observe that the errors of the divergence between the RT $_{k-1}$ and the RT $_{k-1}$ -postprocessed CG $_k$ methods are exactly the same, as predicted by the theory. Moreover, as can be clearly seen from Fig. 5.1, the approximations of the RT $_{k-1}$ and the RT $_{k-1}$ -postprocessed CG $_k$ methods are comparable in accuracy. We also see that the approximate flux provided by the RT $_{k-1}$ -postprocessed CG $_k$ is better than the approximation $-a\nabla U_h$ provided by the CG $_k$ method. Finally, note that if we increase ℓ by one more degree than $k-1$ in RT $_\ell$ -postprocessing, there is no improvement – in fact, the approximate flux given by the RT $_k$ -postprocessed CG $_k$ method produces an approximate flux that is worse than that provided by the RT $_{k-1}$ -postprocessed CG $_k$ method.

TABLE 5.1

Comparison of the history of convergence of the RT_0 and the postprocessed CG_1 methods.

grid level	$\ e_{U_h}\ _0$		$\ e_{\text{div} \mathbf{q}_h}\ _0$		$\ e_{\mathbf{q}_h}\ _0$		$\ e_{a \nabla U_h}\ _0$	
	error	order	error	order	error	order	error	order
RT ₀ method								
1	.42e+0	–	.18e+2	–	.35e+1	–	.44e+1	–
2	.12e+0	1.77	.11e+2	0.69	.27e+1	0.39	.29e+1	0.59
3	.35e-1	1.81	.58e+1	0.95	.14e+1	0.96	.15e+1	0.95
4	.90e-2	1.95	.30e+1	0.99	.69e+0	0.99	.76e+0	0.99
5	.23e-2	1.99	.15e+1	1.00	.35e+0	1.00	.38e+0	1.00
6	.57e-3	2.00	.74e+0	1.00	.17e+0	1.00	.19e+0	1.00
7	.14e-3	2.00	.37e+0	1.00	.87e-1	1.00	.95e-1	1.00
RT ₀ -postprocessed CG ₁ method								
1	.50e+0	–	.18e+2	–	.38e+1	–	.61e+1	–
2	.27e+0	0.89	.11e+2	0.69	.30e+1	0.33	.42e+1	0.54
3	.94e-1	1.52	.58e+1	0.95	.16e+1	0.96	.25e+1	0.76
4	.26e-1	1.83	.30e+1	0.99	.74e+0	1.09	.13e+0	0.92
5	.69e-2	1.95	.15e+1	1.00	.36e+0	1.05	.66e+0	0.98
6	.17e-2	1.98	.74e+0	1.00	.18e+0	1.02	.33e+0	0.99
7	.43e-3	2.00	.37e+0	1.00	.87e-1	1.01	.17e+0	1.00
RT ₁ -postprocessed CG ₁ method								
1	.50e+0	–	.90e+1	–	.83e+1	–	.61e+1	–
2	.27e+0	0.89	.28e+1	1.70	.60e+1	0.47	.42e+1	0.54
3	.94e-1	1.52	.73e+0	1.93	.37e+1	0.67	.25e+1	0.76
4	.26e-1	1.83	.19e+0	1.98	.20e+1	0.91	.13e+1	0.92
5	.69e-2	1.95	.46e-1	2.00	.10e+1	0.98	.66e+0	0.98
6	.17e-2	1.98	.12e-1	2.00	.51e+0	1.00	.33e+0	0.99
7	.43e-3	2.00	.29e-2	2.00	.25e+0	1.00	.17e+0	1.00

Test 2. Now we work on a problem in which the solution has singularities produced by drastic changes in the permeability a , see Fig. 1.1 in the introduction. We compare the streamlines of the approximate flux obtained by the RT_{k-1} -postprocessed CG_k method and that of the RT_{k-1} method around the left upper corner of the rock in Figs. 5.2, 5.3 and 5.4. We see that the presence of the singularity at the corner induces small distortions in the streamlines. However, even in this hard case, the flux produced by the RT_{k-1} method and the solution given by the RT_{k-1} -postprocessed CG_k method are remarkably similar.

6. Concluding remarks. We have shown that a new postprocessing of the CG_k solution gives rise to an $H(\text{div})$ -conforming approximation to the flux which renders the CG method locally conservative. The postprocessing belongs to the Raviart-Thomas space of degree $k-1$ and displays convergence properties similar to the approximation given by the Raviart-Thomas method of degree $k-1$ itself. By counting the degrees of freedom we have established that the computational effort needed to obtain the new postprocessed flux is less than that of the Raviart-Thomas method.

TABLE 5.2

Comparison of the history of convergence of the RT_1 and the postprocessed CG_2 methods.

grid	$\ e_{U_h}\ _0$		$\ e_{\text{div} \mathbf{q}_h}\ _0$		$\ e_{\mathbf{q}_h}\ _0$		$\ e_{a \nabla U_h}\ _0$	
level	error	order	error	order	error	order	error	order
RT ₁ -method								
1	.16e+0	–	.90e+1	–	0.23e+1	–	.42e+1	–
2	.28e-1	2.49	.28e+1	1.70	0.50e+0	2.17	.12e+1	1.87
3	.37e-2	2.94	.73e+0	1.93	0.13e+0	1.94	.30e+0	1.95
4	.47e-3	2.96	.19e+0	1.98	0.33e-1	1.98	.76e-1	1.99
5	.60e-4	2.98	.46e-1	2.00	0.83e-2	1.99	.19e-1	2.00
6	.75e-5	2.99	.12e-1	2.00	0.21e-2	2.00	.47e-2	2.00
7	.94e-6	3.00	.29e-2	2.00	0.52e-3	2.00	.12e-2	2.00
RT ₁ -postprocessed CG_2 method								
1	.24e+0	–	.90e+1	–	.25e+1	–	.40e+1	–
2	.35e-1	2.77	.28e+1	1.70	.76e+0	1.73	.13e+1	1.64
3	.46e-2	2.93	.73e+0	1.93	.16e+0	2.25	.37e+0	1.79
4	.56e-3	3.04	.19e+0	1.98	.35e-1	2.17	.97e-1	1.93
5	.69e-4	3.02	.46e-1	2.00	.85e-2	2.07	.25e-1	1.98
6	.86e-5	3.01	.12e-1	2.00	.21e-2	2.02	.62e-2	1.99
7	.11e-5	3.00	.29e-2	2.00	.52e-3	2.00	.15e-2	2.00
RT ₂ -postprocessed CG_2 method								
1	.24e+0	–	.29e+1	–	.66e+1	–	.40e+1	–
2	.35e-1	2.77	.54e+0	2.40	.21e+1	1.65	.13e+1	1.64
3	.46e-2	2.93	.72e-1	2.92	.67e+0	1.64	.37e+0	1.79
4	.56e-3	3.04	.91e-2	2.98	.18e+0	1.89	.97e-1	1.93
5	.69e-4	3.02	.11e-2	2.99	.46e-1	1.97	.25e-1	1.98
6	.86e-5	3.01	.14e-3	3.00	.12e-1	1.99	.62e-2	1.99
7	.11e-5	3.00	.18e-4	3.00	.29e-2	2.00	.15e-2	2.00

We have also shown how to hybridize the CG method making it easier to treat variable degree approximation spaces and hanging nodes.

The study of the effect of the numerical trace parameter β on the quality of the approximation and the extension of this approach to linear elasticity are subjects of ongoing research.

Acknowledgments. The authors would like to thank the reviewers for bringing to their attention the papers [7] and [26]. They would also like to thank Clint Dawson and Graham F. Carey for bringing to their attention the papers [27, 32] and [9, 28, 10, 14, 11, 8], respectively.

REFERENCES

- [1] D. N. ARNOLD AND F. BREZZI, *Mixed and nonconforming finite element methods: implementation, postprocessing and error estimates*, RAIRO Modél. Math. Anal. Numér., 19 (1985), pp. 7–32.

TABLE 5.3

Comparison of the history of convergence of the RT_2 and the postprocessed CG_3 methods.

grid	$\ e_{U_h}\ _0$		$\ e_{\text{div} \mathbf{q}_h}\ _0$		$\ e_{\mathbf{q}_h}\ _0$		$\ e_{a \nabla U_h}\ _0$	
level	error	order	error	order	error	order	error	order
RT ₂ -method								
1	.12e+0	–	.29e+1	–	0.47e+0	–	.35e+1	–
2	.85e-2	3.84	.54e+0	2.40	0.85e-1	2.46	.50e+0	2.78
3	.60e-3	3.84	.72e-1	2.92	0.11e-1	2.92	.72e-1	2.81
4	.39e-4	3.94	.91e-2	2.98	0.14e-2	2.98	.94e-2	2.94
5	.25e-5	3.98	.11e-2	2.99	0.18e-3	2.99	.12e-2	2.98
6	.16e-6	3.99	.14e-3	3.00	0.22e-4	3.00	.15e-3	2.99
7	.97e-8	4.00	.18e-4	3.00	0.28e-5	3.00	.19e-4	3.00
RT ₂ -postprocessed CG ₃ method								
1	.96e-1	–	.29e+1	–	.14e+1	–	.20e+1	–
2	.63e-2	3.93	.54e+0	2.40	.20e+0	2.82	.27e+0	2.90
3	.35e-3	4.15	.72e-1	2.92	.21e-1	3.23	.35e-1	2.94
4	.20e-4	4.12	.91e-2	2.98	.22e-2	3.29	.43e-2	3.00
5	.12e-5	4.06	.11e-2	2.99	.23e-3	3.23	.54e-3	3.01
6	.75e-7	4.02	.14e-3	3.00	.26e-4	3.15	.67e-4	3.01
7	.47e-8	4.01	.18e-4	3.00	.31e-5	3.09	.83e-5	3.00
RT ₃ -postprocessed CG ₃ method								
1	.96e-1	–	.16e+1	–	.31e+1	–	.20e+1	–
2	.63e-2	3.93	.90e-1	4.12	.48e+0	2.71	.27e+0	2.90
3	.35e-3	4.15	.59e-2	3.93	.63e-1	2.93	.35e-1	2.94
4	.20e-4	4.12	.37e-3	3.98	.79e-2	3.00	.43e-2	3.00
5	.12e-5	4.06	.24e-4	4.99	.98e-3	3.01	.54e-3	3.01
6	.75e-7	4.02	.15e-5	4.00	.12e-3	3.01	.67e-4	3.01
7	.47e-8	4.01	.93e-7	3.98	.15e-4	3.00	.83e-5	3.00

- [2] D. N. ARNOLD, F. BREZZI, B. COCKBURN, AND L. D. MARINI, *Unified analysis of discontinuous Galerkin methods for elliptic problems*, SIAM J. Numer. Anal., 39 (2002), pp. 1749–1779.
- [3] J. BARRETT AND C. ELLIOTT, *Total flux estimates for a finite-element approximation of elliptic equations*, IMA J. Numer. Anal., 7 (1987), pp. 129–148.
- [4] P. BASTIAN AND B. RIVIÈRE, *Superconvergence and $H(\text{div})$ projection for discontinuous Galerkin methods*, Internat. J. Numer. Methods Fluids, 42 (2003), pp. 1043–1057.
- [5] J. H. BRAMBLE, J. E. PASCIAK, AND A. H. SCHATZ, *The construction of preconditioners for elliptic problems by substructuring. I*, Math. Comp., 47 (1986).
- [6] F. BREZZI AND M. FORTIN, *Mixed and Hybrid finite element methods*, Springer Verlag, 1991.
- [7] F. BREZZI, T. J. R. HUGHES, AND E. SÜLI, *Variational approximation of flux in conforming finite element methods for elliptic partial differential equations: a model problem*, Atti Accad. Naz. Lincei Cl. Sci. Fis. Mat. Natur. Rend. Lincei (9) Mat. Appl., 12 (2001), pp. 159–166 (2002).
- [8] G. CAREY, *Derivative calculation from finite element solutions*, Comput. Methods Appl. Mech. Engrg., 35 (1982), pp. 1–14.
- [9] ———, *Some further properties of the superconvergent flux projection*, Comput. Methods Appl. Mech. Engrg., 18 (2002), pp. 241–250.
- [10] G. F. CAREY, G. BICKEN, V. CAREY, C. BERGER, AND J. SANCHEZ, *Locally constrained projections on grids*, Internat. J. Numer. Methods Engrg., 50 (2001), pp. 549–577.

- [11] G. F. CAREY, S.-S. CHOW, AND M. K. SEAGER, *Approximate boundary-flux calculations*, Comput. Methods Appl. Mech. Engrg., 50 (1985), pp. 107–120.
- [12] F. CELIKER AND B. COCKBURN, *Superconvergence of the numerical traces of discontinuous Galerkin and hybridized mixed methods for convection-diffusion problems in one space dimension*, Math. Comp., 67 (2007), pp. 67–96.
- [13] P. G. CIARLET, *The Finite Element Method for Elliptic Problems*, North-Holland Publishing Company, Amsterdam, 1978.
- [14] S.-S. CHOW, G. F. CAREY, AND R. D. LAZAROV, *Natural and postprocessed superconvergence in semilinear problems*, Numer. Methods Partial Differential Equations, 7 (1991), pp. 245–259.
- [15] B. COCKBURN AND J. GOPALAKRISHNAN, *A characterization of hybridized mixed methods for second order elliptic problems*, SIAM J. Numer. Anal., 42 (2004), pp. 283–301.
- [16] ———, *Error analysis of variable degree mixed methods for elliptic problems via hybridization*, Math. Comp., 74 (2005), pp. 1653–1677 (electronic).
- [17] ———, *Incompressible finite elements via hybridization. Part I: The Stokes system in two space dimensions*, SIAM J. Numer. Anal., 43 (2005), pp. 1627–1650.
- [18] ———, *Incompressible finite elements via hybridization. Part II: The Stokes system in three space dimensions*, SIAM J. Numer. Anal., 43 (2005), pp. 1651–1672.
- [19] B. COCKBURN, G. KANSCHAT, AND D. SCHÖTZAU, *A locally conservative LDG method for the incompressible Navier-Stokes equations*, Math. Comp., 74 (2005), pp. 1067–1095.
- [20] R. COURANT, *Variational methods for the solution of problems of equilibrium and vibrations*, Bull. Amer. Math. Soc., 49 (1943), pp. 1–23.
- [21] L. DEMKOWICZ, *2D hp-adaptive finite element package (2Dhp90). version 2.0*, Technical Report 02–06, Texas Institute for Computational and Applied Mathematics, 2002.
- [22] L. DEMKOWICZ AND A. BUFFA, H^1 , $H(\text{curl})$ and $H(\text{div})$ -conforming projection-based interpolation in three dimensions. Quasi-optimal p -interpolation estimates, Comput. Methods Appl. Mech. Engrg., 194 (2005), pp. 267–296.
- [23] J. DOUGLAS, JR., T. DUPONT, AND M. F. WHEELER, *A Galerkin procedure for approximating the flux on the boundary for elliptic and parabolic boundary value problems*, RAIRO Modél. Math. Anal. Numér., (1974), pp. 47–59.
- [24] J. GOPALAKRISHNAN, *A Schwarz preconditioner for a hybridized mixed method*, Comput. Methods Appl. Math., 3 (2003), pp. 116–134.
- [25] J. T. R. HUGHES, G. ENGEL, L. MAZZEI, AND M. LARSON, *The continuous Galerkin method is locally conservative*, J. Comput. Phys., 163 (2000), pp. 467–488.
- [26] J. T. R. HUGHES AND G. N. WELLS, *Conservation properties for the Galerkin and stabilised forms of the advection-diffusion and incompressible Navier-Stokes equations*, Comput. Methods Appl. Mech. Engrg., 194 (2005), pp. 1141–1159.
- [27] M. LARSON AND A. NIKLASSON, *A conservative flux for the continuous Galerkin method based on discontinuous enrichment*, Calcolo, 41 (2004), pp. 65–76.
- [28] A. I. PEHLIVANOV, R. D. LAZAROV, G. F. CAREY, AND S.-S. CHOW, *Superconvergence analysis of approximate boundary-flux calculations*, Numer. Math., 63 (1992), pp. 483–501.
- [29] P. A. RAVIART AND J. M. THOMAS, *A mixed finite element method for second order elliptic problems*, in Mathematical Aspects of Finite Element Method, I. Galligani and E. Magenes, eds., Lecture Notes in Math. 606, Springer-Verlag, New York, 1977, pp. 292–315.
- [30] L. R. SCOTT AND S. ZHANG, *Finite element interpolation of nonsmooth functions satisfying boundary conditions*, Math. Comp., 54 (1990), pp. 483–493.
- [31] E. STEIN AND M. RÜTTER, *Finite Element Methods for Elasticity with Error-Controlled Discretization and Model Adaptivity*, in Encyclopedia of Computational Mechanics, R. d. B. E. Stein and T. Hughes, eds., vol. 2, John Wiley & Sons, Ltd., England, 2004, pp. 5–58.
- [32] S. SUN AND M. WHEELER, *Projections of velocity data for the compatibility with transport*, Comput. Methods Appl. Mech. Engrg., 195 (2006), pp. 653–673.
- [33] J. A. WHEELER, *Simulation of heat transfer from a warm pipe buried in permafrost*, 74th. National Meeting of the American Institute of Chemical Engineers, New Orleans, (March, 1973).
- [34] M. F. WHEELER, *A Galerkin procedure for estimating the flux for two-point boundary value problems*, SIAM J. Numer. Anal., 11 (1974), pp. 764–768.

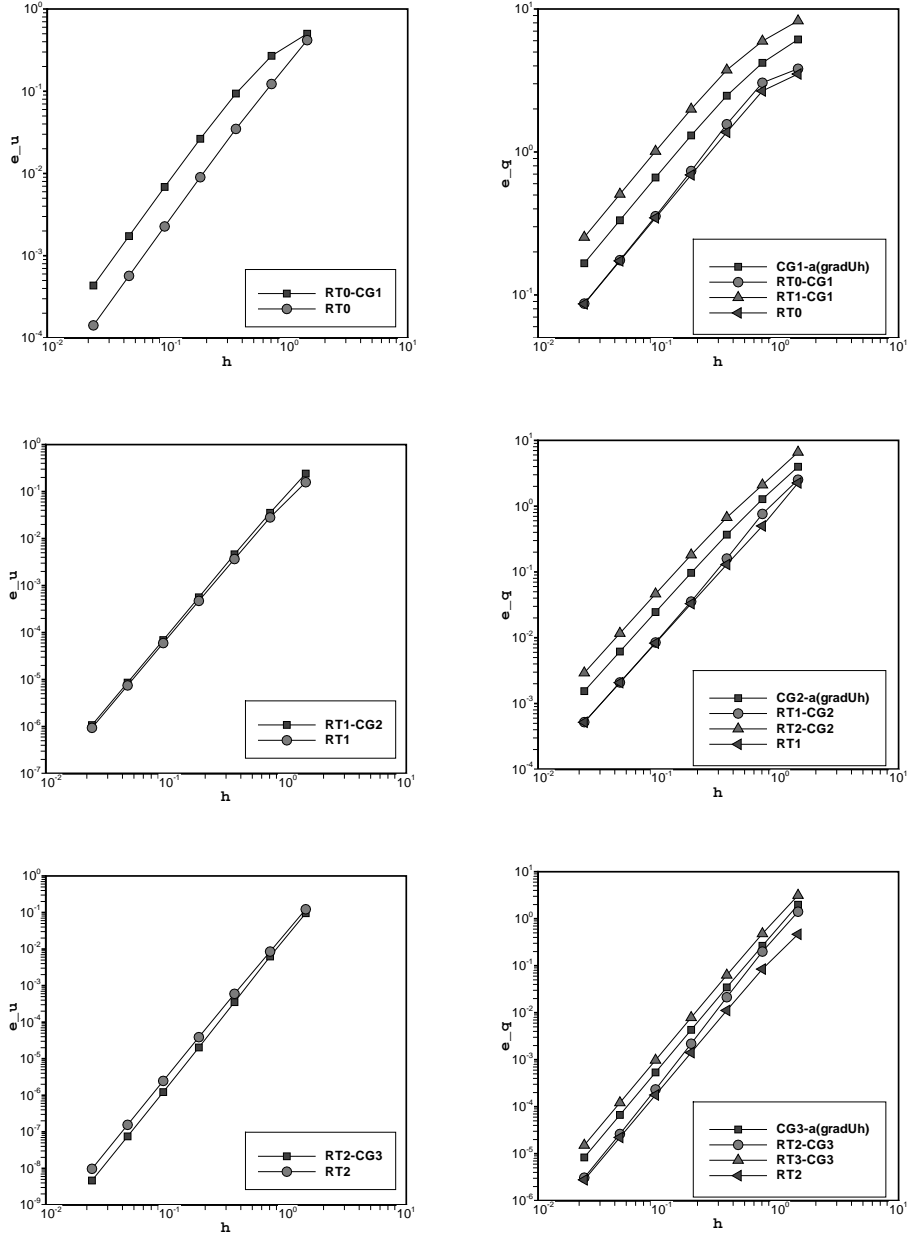


FIG. 5.1. History of convergence for u_h (left) and q_h (right), for $k = 1$ (first row), $k = 2$ (second row) and $k = 3$ (third row). Note: CGk-a(gradUh) is $a\nabla U_h$ from CG_k , $k = 1, 2$ or 3 .

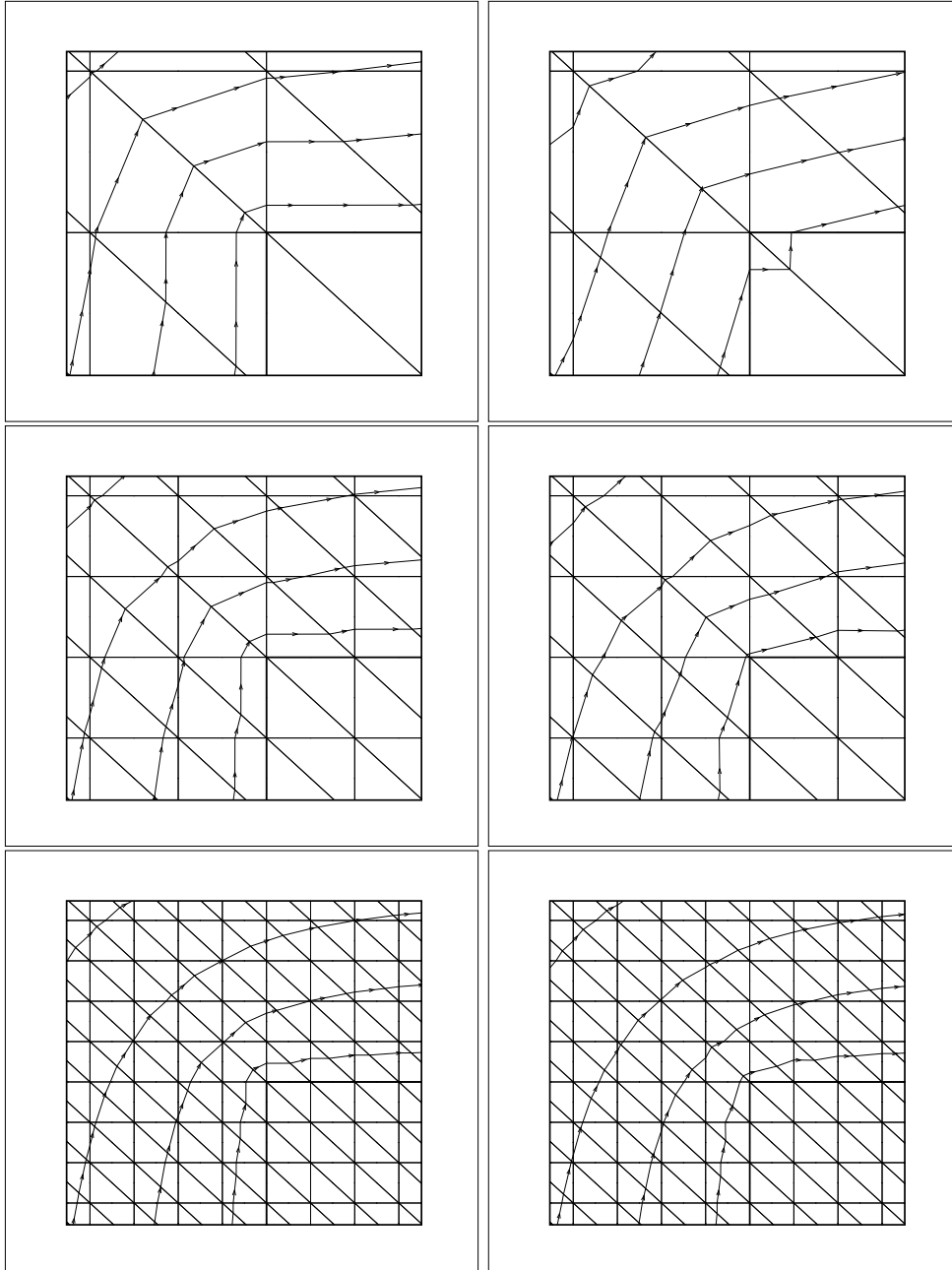


FIG. 5.2. Streamlines in the left upper corner with $k = 1$. On the left column is the solution given by the RT_k method and on the right column that of the $RT_{(k-1)}$ -postprocessed CG_k method. From top to bottom, mesh size $h = \frac{1}{8}, \frac{1}{16}, \frac{1}{32}$.

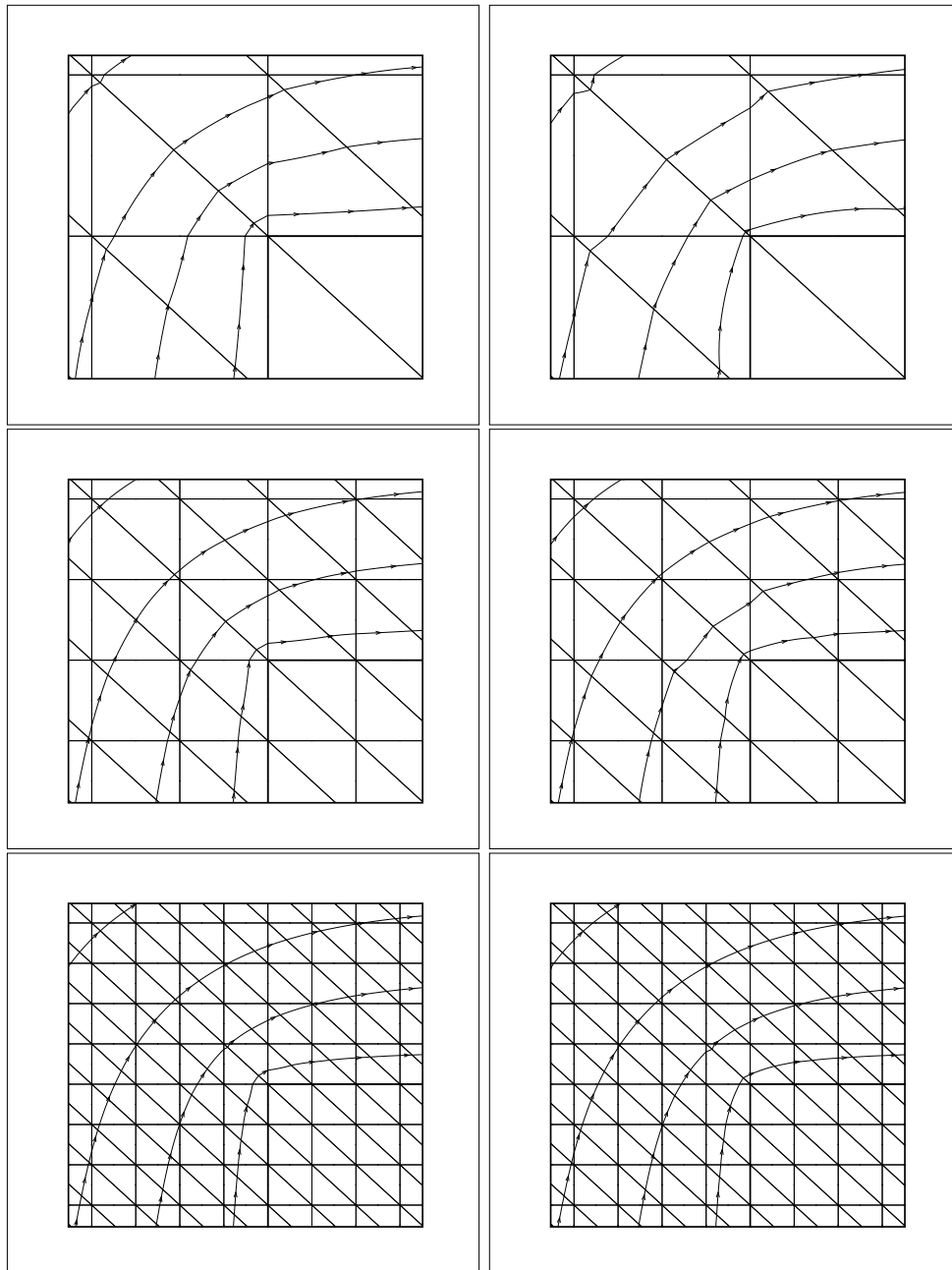


FIG. 5.3. Streamlines in the left upper corner with $k = 2$. On the left column is the solution given by the RT_k method and on the right column that of the $RT_{(k-1)}$ -postprocessed CG_k method. From top to bottom, mesh size $h = \frac{1}{8}, \frac{1}{16}, \frac{1}{32}$.

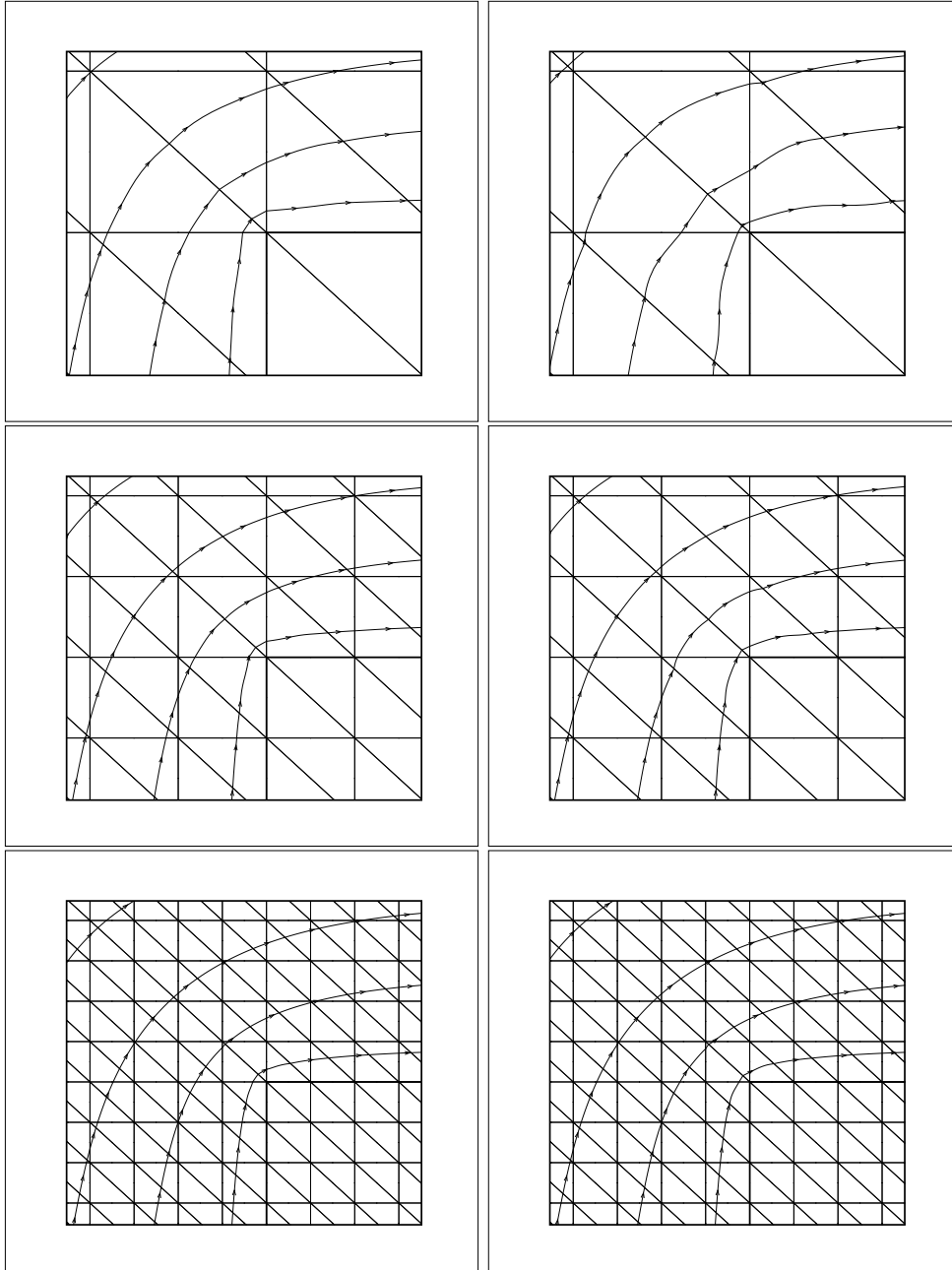


FIG. 5.4. Streamlines in the left upper corner with $k = 3$. On the left column is the solution given by the RT_k method and on the right column that of the $RT_{(k-1)}$ -postprocessed CG_k method. From top to bottom, mesh size $h = \frac{1}{8}, \frac{1}{16}, \frac{1}{32}$.