

MIXED METHODS FOR TWO-PHASE DARCY–STOKES MIXTURES OF PARTIALLY MELTED MATERIALS WITH REGIONS OF ZERO POROSITY*

TODD ARBOGAST[†], MARC A. HESSE[‡], AND ABRAHAM L. TAICHER[§]

Abstract. The Earth’s mantle (or, e.g., a glacier) involves a deformable solid matrix phase within which a second phase, a fluid, may form due to melting processes. The system is modeled as a dual-continuum mixture, with at each point of space the solid matrix being governed by a Stokes flow and the fluid melt, if it exists, being governed by Darcy’s law. This system is mathematically degenerate when the porosity (volume fraction of fluid) vanishes. Assuming the porosity is given, we develop a mixed variational framework for the mechanics of the system by carefully scaling the Darcy variables by powers of the porosity. We prove that the variational problem is well-posed, even when there are regions of one and two phases. We then develop an accurate mixed finite element method for solving this Darcy–Stokes system and prove a convergence result. Numerical results are presented that illustrate and verify the convergence of the method.

Key words. degenerate elliptic, energy bounds, mixed finite element method, mantle dynamics, glaciers, midocean ridge

AMS subject classifications. 65N12, 65N30, 35J70, 76M10, 76S05, 76T99

DOI. 10.1137/16M1091095

1. Introduction. The equations of mantle dynamics introduced by McKenzie [28] have a wide range of applications in Earth physics [2, 27, 26, 25], such as in modeling midocean ridges, subduction zones, and hot-spot volcanism, as well as to glacier dynamics [22, 9, 37] and other two-phase flows in porous media [13, 18]. For example, at a midocean ridge, melt is believed to migrate upward until it reaches the lithospheric “tent” where it then moves toward the ridge within a high porosity band. Simulation of this phenomenon requires numerical methods that accurately handle highly heterogeneous porosity and the single-phase to two-phase transition.

The model assumes a dual-continuum mixture of solid matrix and fluid melt. The mixing parameter is the porosity ϕ , i.e., the volume fraction of fluid melt, which is assumed to be much smaller than one, but it may be zero in parts of the domain where there is no fluid melt.

We use subscripts f , s , and r to refer to a quantity associated with the fluid melt, the matrix solid, or the relative fluid minus solid, respectively. Fluid melt forms at the boundaries of rock crystals and so obeys Darcy’s law for fluid flow around solid

*Submitted to the journal’s Computational Methods in Science and Engineering section August 25, 2016; accepted for publication (in revised form) January 20, 2017; published electronically March 30, 2017.

<http://www.siam.org/journals/sisc/39-2/M109109.html>

Funding: This work was supported by the U.S. National Science Foundation under grants EAR-1025321 and DMS-1418752. The first two authors thank the Isaac Newton Institute for Mathematical Sciences, Cambridge (supported by EPSRC grant EP/K032208/1), for support during the programme Melt in the Mantle, 15 Feb. to 17 June 2016, when some work on this paper was undertaken.

[†]Department of Mathematics, University of Texas, 2515 Speedway, C1200, Austin, TX 78712-1202 and Institute for Computational Engineering and Sciences, University of Texas, 201 East 24th St., C0200, Austin, TX 78712-1229 (arbogast@ices.utexas.edu).

[‡]Jackson School of Geosciences, University of Texas, 2305 Speedway, C1160, Austin, TX 78712-1692 (mhesse@jsg.utexas.edu).

[§]Institute for Computational Engineering and Sciences, University of Texas, 201 East 24th St., C0200, Austin, TX 78712-1229 (ataicher@ices.utexas.edu).

matrix “grains,” which is

$$(1.1) \quad \mathbf{u} = \phi \mathbf{v}_r = \phi(\mathbf{v}_f - \mathbf{v}_s) = -\frac{k_0 \phi^{2+2\Theta}}{\mu_f} (\nabla p_f - \rho_f \mathbf{g}),$$

where \mathbf{u} is the Darcy velocity, \mathbf{v} and p are the velocity and pressure, μ is the viscosity, $k(\phi) = k_0 \phi^{2+2\Theta}$ is the porosity dependent permeability with Θ a constant between 0 and 1/2 (see, e.g., [13, 38]), ρ is the density, and \mathbf{g} is the downwards pointing gravitational vector.

As the two phases melt or solidify, total mass is conserved. After applying a *Boussinesq approximation* [35] (constant and equal densities for nonbuoyancy terms), this is expressed as

$$(1.2) \quad \nabla \cdot (\mathbf{u} + \mathbf{v}_s) = 0.$$

Conservation of momentum for the slowly creeping mixture obeys the Stokes equation. In terms of the deviatoric stress of the mixture

$$(1.3) \quad \hat{\boldsymbol{\sigma}} = \hat{\boldsymbol{\sigma}}(\mathbf{v}_s) = 2\mu_s(1 - \phi)(\mathcal{D}\mathbf{v}_s - \frac{1}{3}\nabla \cdot \mathbf{v}_s \mathbf{I}),$$

wherein $\mathcal{D}\mathbf{v}_s = \frac{1}{2}(\nabla \mathbf{v}_s + \nabla \mathbf{v}_s^T)$ is the symmetric gradient, we have that

$$(1.4) \quad -\nabla \bar{p} + \nabla \cdot \hat{\boldsymbol{\sigma}}(\mathbf{v}_s) = -(\rho_s + \phi \rho_r) \mathbf{g},$$

where $\bar{p} = \phi p_f + (1 - \phi)p_m = p_s + \phi p_r$ is the mixture pressure.

The mechanical system is closed by relating the solid and fluid pressures through a compaction relation [34]

$$(1.5) \quad p_s - p_f = -\frac{\mu_s}{\phi} \nabla \cdot \mathbf{v}_s,$$

where μ_s/ϕ is the solid matrix bulk viscosity.

When coupled with solute transport and thermal evolution, the model transitions dynamically in time from a nonporous single phase Stokes solid to a two-phase porous medium. Because the model is based on mixture theory, it has the advantage that the free boundary between the one- and two-phase regions need *not* be determined explicitly in the numerical approximation. Unfortunately, the disadvantage is that the Darcy part of the equations is mathematically degenerate in regions where the porosity is zero, since then there is only the one solid phase, even though the model equations continue to describe both phases over the entire domain Ω . In this paper we assume that $\phi(\mathbf{x})$ is given at some instant of time, and we discuss only the mechanics part of the full model.

A mixed finite element method (MFEM) is a good candidate for a computational approximation of the mechanics part of this model. MFEMs have an extensive theory for both Darcy and Stokes flows. Moreover, velocity fields computed using MFEM are continuous on each element and have a continuous normal component across element boundaries. This allows coupling with the transport equations of solute and thermal evolution, since the velocities unambiguously determine particle trajectories.

The Stokes part of the system is well-behaved, but the Darcy part has difficulties when ϕ vanishes. Later, we will see (2.26) and (2.27), which imply that

$$(1.6) \quad \|\phi^{-1-\Theta} \mathbf{u}\| + \|\phi^{-1/2} \nabla \cdot \mathbf{u}\| + \|\phi^{1/2} p_f\| \leq C$$

for some constant C , where $\|\cdot\|$ is the $L^2(\Omega)$ -norm. These estimates suggest that the fluid pressure p_f may be *unbounded* where porosity vanishes. Indeed, the fluid pressure is no longer a physical variable when there is no fluid. Moreover, any numerical method that does not take into account the degeneracy of ϕ , say by instead imposing a small nonzero porosity ϕ_0 everywhere, is sure to have a condition number that grows as $\phi_0 \rightarrow 0$. Our numerical results will show these issues.

Recently, two of the current authors [6, 5] developed an MFEM and cell-centered finite difference method for a single Darcy system with a similar degeneracy as appears in (2.4)–(2.5). The key is to follow the hint in the stability estimates and scale the fluid pressure and velocity to avoid problems with vanishing porosity. In this paper we apply this idea to the full set of mantle mechanics equations (1.1)–(1.5).

In the rest of the paper, we present in section 2 our scaled formulation that directly resolves the issue of degenerate porosity. We prove the existence and uniqueness of a solution to the scaled variational formulation. In section 3 we define our MFEM for the numerical approximation of the scaled variational formulation and prove its convergence. In section 4, we present a modification of the MFEM that is locally mass conservative. In section 5, we discuss implementation and give a mass lumping modification that results in a relatively simple solution procedure on rectangular meshes. Numerical results illustrating and evaluating the effects of degenerate porosity are given in sections 6–7. We include tests of a one-dimensional compacting column with various porosity functions, and a two-dimensional test example akin to a midocean ridge. We conclude the paper in section 8.

2. A scaled mixed variational formulation. Define the pressure potentials

$$(2.1) \quad q_f = p_f - \rho_f g z \quad \text{and} \quad q_s = p_s - \rho_f g z,$$

where z is depth and indeed q_s is defined using the *fluid* density ρ_f . Also let

$$(2.2) \quad q = \phi q_f + (1 - \phi) q_s = q_s + \phi(q_f - q_s)$$

be the mixture potential and note that

$$(2.3) \quad p_f - p_s = q_f - q_s = \frac{1}{1 - \phi}(q_f - q).$$

We find it convenient to remove q_s from (1.1)–(1.5). We obtain

$$(2.4) \quad \mathbf{u} + \frac{k_0 \phi^{2+2\Theta}}{\mu_f} \nabla q_f = 0,$$

$$(2.5) \quad \mu_s \nabla \cdot \mathbf{u} + \frac{\phi}{1 - \phi}(q_f - q) = 0,$$

$$(2.6) \quad \nabla q - \nabla \cdot \hat{\boldsymbol{\sigma}}(\mathbf{v}_s) = -(1 - \phi) \rho_r \mathbf{g},$$

$$(2.7) \quad \mu_s \nabla \cdot \mathbf{v}_s - \frac{\phi}{1 - \phi}(q_f - q) = 0,$$

where the deviatoric stress of the mixture is given in (1.3). For simplicity, the model parameters are assumed to be constant. Equation (2.4) represents Darcy's law for an incompressible fluid, (2.6)–(2.7), (1.3) is a Stokes system for a highly viscous, compressible material (matrix plus fluid), and (2.5) plus (2.7) enforces mass conservation.

We suppose that the spatial domain Ω is a bounded, simply connected, Lipschitz domain in \mathbb{R}^d , $d = 1, 2$, or 3 , with outward pointing unit normal vector ν . We impose

boundary conditions on the fluid and solid velocity of the form

$$(2.8) \quad \mathbf{u} \cdot \boldsymbol{\nu} = g_r \quad \text{and} \quad \mathbf{v}_s = \mathbf{g}_s \quad \text{on } \partial\Omega.$$

We need the compatibility condition

$$(2.9) \quad \int_{\partial\Omega} (g_r + \mathbf{g}_s \cdot \boldsymbol{\nu}) \, ds = 0.$$

2.1. Standard function spaces. The space $L^2(\Omega)$ consists of all square integrable, real-valued functions on Ω . It is equipped with the inner product $(u, v) = (u, v)_\Omega = \int_\Omega uv \, dx$ and associated norm $\|u\| = (u, u)^{1/2}$. Denote by $H^1(\Omega)$ all square integrable functions with square integrable weak derivatives. This space has the norm $\|u\|_1 = \{\|u\|^2 + \|\nabla u\|^2\}^{1/2}$. Let $H(\text{div}; \Omega)$ denote all square integrable vector-valued functions with square integrable weak divergence, and equip it with the norm $\|\mathbf{u}\|_{H(\text{div})} = \{\|\mathbf{u}\|^2 + \|\nabla \cdot \mathbf{u}\|^2\}^{1/2}$.

We can restrict functions in $H^1(\Omega)$ to the boundary $\partial\Omega$ using the trace lemma [1, 24]. The space of these restrictions is $H^{1/2}(\partial\Omega) \subset L^2(\partial\Omega)$, and we have the bound

$$(2.10) \quad \|u\|_{1/2, \partial\Omega} \leq C_\Omega \|u\|_1.$$

A similar lemma holds for functions in $H(\text{div}; \Omega)$ [17], and

$$(2.11) \quad \|\mathbf{u} \cdot \boldsymbol{\nu}\|_{-1/2, \partial\Omega} \leq C_\Omega \|\mathbf{u}\|_{H(\text{div}; \Omega)},$$

where $\|\cdot\|_{-1/2, \partial\Omega}$ is the norm of the dual space of $H^{1/2}(\partial\Omega)$.

The space $L^\infty(\Omega)$ consists of all essentially bounded functions on Ω equipped with the essential supremum norm $\|\cdot\|_{L^\infty(\Omega)}$. The space $W^{1,\infty}(\Omega)$ consists of the functions in $L^\infty(\Omega)$ that have weak derivatives also in $L^\infty(\Omega)$, and the norm is $\|\cdot\|_{W^{1,\infty}(\Omega)} = \|\cdot\|_{L^\infty(\Omega)} + \|\nabla(\cdot)\|_{(L^\infty(\Omega))^d}$.

2.2. The scaled formulation. Following [6], we define the *scaled relative velocity* and *scaled fluid potential*

$$(2.12) \quad \tilde{\mathbf{v}}_r = \phi^{-1-\Theta} \mathbf{u} \quad \text{and} \quad \tilde{q}_f = \phi^{1/2} q_f,$$

respectively, and we reformulate the problem (2.4)–(2.7) as

$$(2.13) \quad \tilde{\mathbf{v}}_r + \frac{k_0 \phi^{1+\Theta}}{\mu_f} \nabla(\phi^{-1/2} \tilde{q}_f) = 0,$$

$$(2.14) \quad \mu_s \phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} \tilde{\mathbf{v}}_r) + \frac{1}{1-\phi} (\tilde{q}_f - \phi^{1/2} q) = 0,$$

$$(2.15) \quad \nabla q - \nabla \cdot \hat{\boldsymbol{\sigma}}(\mathbf{v}_s) = -(1-\phi) \rho_r \mathbf{g},$$

$$(2.16) \quad \mu_s \nabla \cdot \mathbf{v}_s - \frac{\phi^{1/2}}{1-\phi} (\tilde{q}_f - \phi^{1/2} q) = 0,$$

wherein we have scaled the entire second equation by $\phi^{-1/2}$. The boundary condition on \mathbf{u} in (2.8) rescales to $\phi^{1+\Theta} \tilde{\mathbf{v}}_r \cdot \boldsymbol{\nu} = g_r$ on $\partial\Omega$.

The scaled equations make sense provided that the gradient and divergence terms are well-defined when $\phi = 0$. The divergence term in (2.14) expands to

$$(2.17) \quad \phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} \tilde{\mathbf{v}}_r) = \phi^{1/2+\Theta} \nabla \cdot \tilde{\mathbf{v}}_r + \phi^{\Theta-1/2} \nabla \phi \cdot \tilde{\mathbf{v}}_r,$$

and it is well-defined provided that, for example,

$$(2.18) \quad \phi^{\Theta-1/2} \nabla \phi \in (L^\infty(\Omega))^d.$$

The gradient terms in (2.13) make sense under the same condition. The porosity ϕ in the physical model satisfies the full set of equations, including solute and thermal transport equations. It is *not* clear if we should expect that this porosity satisfies our condition. Nevertheless, we will assume that the condition holds. Our numerical results suggest that it is not strictly necessary, and perhaps can be weakened (see also [6] for a discussion of the necessity of this condition).

We should not expect the scaled velocity $\tilde{\mathbf{v}}_r$ to lie in $H(\operatorname{div}; \Omega)$; rather, $\tilde{\mathbf{v}}_r$ should lie in the space

$$\tilde{\mathbb{V}}_r = H_\phi(\operatorname{div}; \Omega) = \{ \mathbf{v} \in (L^2(\Omega))^d : \phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} \mathbf{v}) \in L^2(\Omega) \}.$$

As discussed in [6], this is a Hilbert space with the inner product

$$(\tilde{\mathbf{u}}, \tilde{\mathbf{v}})_{\tilde{\mathbb{V}}_r} = (\tilde{\mathbf{u}}, \tilde{\mathbf{v}}) + (\phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} \tilde{\mathbf{u}}), \phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} \tilde{\mathbf{v}})).$$

Moreover, these vector functions have a well-defined normal trace on $\partial\Omega$, and, similarly to (2.11),

$$(2.19) \quad \|\phi^{1/2+\Theta} \tilde{\mathbf{v}} \cdot \nu\|_{-1/2, \partial\Omega} \leq C_\Omega \{ \|\tilde{\mathbf{v}}\| + \|\phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} \tilde{\mathbf{v}})\| \}.$$

We also have the space $H_\phi^{-1/2}(\partial\Omega)$, which is the image of this normal trace operator on $\tilde{\mathbb{V}}_r = H_\phi(\operatorname{div}; \Omega)$.

2.3. The scaled weak formulation. Define the function spaces

$$\begin{aligned} \tilde{\mathbb{V}}_{r,0} &= \{ \mathbf{v} \in H_\phi(\operatorname{div}; \Omega) : \phi^{1/2+\Theta} \mathbf{v} \cdot \nu = 0 \text{ on } \partial\Omega \}, \\ \mathbb{W}_f &= L^2(\Omega), \\ \mathbb{V}_{s,0} &= (H_0^1(\Omega))^d = \{ \mathbf{v} \in (H^1(\Omega))^d : \mathbf{v} = \mathbf{0} \text{ on } \partial\Omega \}, \\ \mathbb{W}_0 &= L^2(\Omega)/\mathbb{R} = \left\{ w \in L^2(\Omega) : \int_\Omega w \, dx = 0 \right\}, \end{aligned}$$

each with its natural norm.

To impose essential boundary conditions (2.8), we assume that $\mathbf{g}_s \in (H^{1/2}(\partial\Omega))^d$ and extend it continuously from the boundary into the domain, so that the extension $\mathbf{g}_s \in \mathbb{V}_s = (H^1(\Omega))^d$ and $\|\mathbf{g}_s\|_1 \leq C\|\mathbf{g}_s\|_{1/2, \partial\Omega}$. In a similar way, following [6], we assume that $\phi^{-1/2} g_r \in H_\phi^{-1/2}(\partial\Omega)$, the image of the scaled normal trace operator on $\tilde{\mathbb{V}}_r = H_\phi(\operatorname{div}; \Omega)$ which appears in (2.19). Then $\phi^{-1/2} g_r$ has a bounded extension $\mathbf{g}_r \in \tilde{\mathbb{V}}_r$ on Ω such that

$$(2.20) \quad \phi^{1/2+\Theta} \mathbf{g}_r \cdot \nu = \phi^{1/2+\Theta} \tilde{\mathbf{v}}_r \cdot \nu = \phi^{-1/2} g_r \quad \text{on } \partial\Omega.$$

We require the scaled compatibility condition

$$(2.21) \quad \int_{\partial\Omega} (\mathbf{g}_r + \mathbf{g}_s) \cdot ds = 0.$$

Scaled formulation. Find $\tilde{\mathbf{v}}_r \in \tilde{\mathbb{V}}_{r,0} + \mathbf{g}_r$, $\tilde{q}_f \in \mathbb{W}_f$, $\mathbf{v}_s \in \mathbb{V}_{s,0} + \mathbf{g}_s$, and $q \in \mathbb{W}_0$ such that

(2.22)

$$\left(\frac{\mu_f}{k_0} \tilde{\mathbf{v}}_r, \boldsymbol{\psi}_r \right) - (\tilde{q}_f, \phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} \boldsymbol{\psi}_r)) = 0 \quad \forall \boldsymbol{\psi}_r \in \tilde{\mathbb{V}}_{r,0},$$

$$(\phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} \tilde{\mathbf{v}}_r), w_f)$$

(2.23)

$$+ \left(\frac{1}{\mu_s(1-\phi)} (\tilde{q}_f - \phi^{1/2} q), w_f \right) = 0 \quad \forall w_f \in \mathbb{W}_f,$$

(2.24)

$$-(q, \nabla \cdot \boldsymbol{\psi}_s) + (\hat{\boldsymbol{\sigma}}(\mathbf{v}_s), \nabla \boldsymbol{\psi}_s) = -((1-\phi)\rho_r \mathbf{g}, \boldsymbol{\psi}_s) \quad \forall \boldsymbol{\psi}_s \in \mathbb{V}_{s,0},$$

(2.25)

$$(\nabla \cdot \mathbf{v}_s, w) - \left(\frac{\phi^{1/2}}{\mu_s(1-\phi)} (\tilde{q}_f - \phi^{1/2} q), w \right) = 0 \quad \forall w \in \mathbb{W}_0.$$

2.4. Existence and uniqueness of the solution. The following theorem shows that the scaled model is well-posed. (See also [4] for treatment of the Dirichlet condition on the Darcy system.)

THEOREM 1. *Assume that (2.18) holds on the porosity, $0 \leq \phi \leq \phi^* < 1$, and the extensions $\mathbf{g}_r \in \tilde{\mathbb{V}}_r$ and $\mathbf{g}_s \in \mathbb{V}_s$ satisfy (2.21). Then there exists a unique solution to the scaled formulation (2.22)–(2.25), (1.3), and it satisfies*

$$\begin{aligned} & \|\tilde{\mathbf{v}}_r\| + \|\phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} \tilde{\mathbf{v}}_r)\| + \|\tilde{q}_f\| + \|\mathbf{v}_s\|_1 + \|q\| \\ (2.26) \quad & \leq C\{|\rho_r| + \|\mathbf{g}_r\| + \|\phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} \mathbf{g}_r)\| + \|\mathbf{g}_s\|_1\}. \end{aligned}$$

The unscaled equations (2.4)–(2.7) are ill-posed where $\phi = 0$. If we restrict $\phi \geq \phi_* > 0$, the equations are well-posed, and we can unscale the variables in (2.26) to show the bound

$$(2.27) \quad \|\phi^{-1-\Theta} \mathbf{u}\| + \|\phi^{-1/2} \nabla \cdot \mathbf{u}\| + \|\phi^{1/2} q_f\| + \|\mathbf{v}_s\|_1 + \|q\| \leq C$$

(i.e., (1.6)). We conclude that the two velocities and the solid matrix pressure remain stable, i.e., they are bounded, as $\phi_* \rightarrow 0$, but the fluid potential *may* become unbounded. This potential loss of stability is a significant issue for numerical modeling. We remark that the correct scaling (2.12) is found by restricting $\phi \geq \phi_* > 0$ and showing directly the bound (2.27) (see also [6, 4]).

Before proving the theorem, we state a well-known result [10, 11, 30] that we need.

THEOREM 2 (Babuška–Lax–Milgram). *Let U and V be two real Hilbert spaces. Suppose that $a : U \times V \rightarrow \mathbb{R}$ is a continuous bilinear functional such that for some constant $\gamma > 0$ and all $u \in U$ and $v \in V$, $v \neq 0$,*

$$(2.28) \quad \sup_{\|v\|=1} |a(u, v)| \geq \gamma \|u\| \quad \text{and} \quad \sup_{\|u\|=1} |a(u, v)| > 0.$$

Then, for all $f \in V^$, there exists a unique solution $u \in U$ to*

$$a(u, v) = f(v) \quad \forall v \in V,$$

and

$$(2.29) \quad \|u\| \leq \frac{1}{\gamma} \|f\|.$$

Proof of Theorem 1. Let

$$\mathbb{X} = \tilde{\mathbb{V}}_{r,0} \times \mathbb{W}_f \times \mathbb{V}_{s,0} \times \mathbb{W}_0,$$

and take $U = V = \mathbb{X}$, which is indeed a real Hilbert space. The bilinear form is defined by (2.22)–(2.25) for any $\mathbf{U} = (\tilde{\mathbf{v}}_{r,0}, \tilde{q}_f, \mathbf{v}_{s,0}, q) \in \mathbb{X}$ and $\Psi = (\boldsymbol{\psi}_r, w_f, \boldsymbol{\psi}_s, w) \in \mathbb{X}$ as

$$\begin{aligned} a(\mathbf{U}, \Psi) &= \left(\frac{\mu_f}{k_0} \tilde{\mathbf{v}}_{r,0}, \boldsymbol{\psi}_r \right) - (\tilde{q}_f, \phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} \boldsymbol{\psi}_r)) + (\phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} \tilde{\mathbf{v}}_{r,0}), w_f) \\ &\quad + \left(\frac{1}{\mu_s(1-\phi)} (\tilde{q}_f - \phi^{1/2} q), w_f - \phi^{1/2} w \right) \\ &\quad - (q, \nabla \cdot \boldsymbol{\psi}_s) + (\hat{\boldsymbol{\sigma}}(\mathbf{v}_{s,0}), \nabla \boldsymbol{\psi}_s) + (\nabla \cdot \mathbf{v}_{s,0}, w). \end{aligned}$$

The linear functional is

$$\begin{aligned} f(\Psi) &= -((1-\phi)\rho_r \mathbf{g}, \boldsymbol{\psi}_s) - \left(\frac{\mu_f}{k_0} \mathbf{g}_r, \boldsymbol{\psi}_r \right) - (\phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} \mathbf{g}_r), w_f) \\ &\quad - (\hat{\boldsymbol{\sigma}}(\mathbf{g}_s), \nabla \boldsymbol{\psi}_s) - (\nabla \cdot \mathbf{g}_s, w). \end{aligned}$$

Clearly we have continuity (boundedness) of a on $\mathbb{X} \times \mathbb{X}$ and f on \mathbb{X} .

Our scaled formulation is written in the context of the Babuška–Lax–Milgram theorem as follows. We find $\mathbf{U} \in \mathbb{X}$ such that

$$(2.30) \quad a(\mathbf{U}, \Psi) = f(\Psi) \quad \forall \Psi \in \mathbb{X},$$

and then set $\tilde{\mathbf{v}}_r = \tilde{\mathbf{v}}_{r,0} + \mathbf{g}_r$ and $\mathbf{v}_s = \mathbf{v}_{s,0} + \mathbf{g}_s$.

We will need an estimate of the term $(\hat{\boldsymbol{\sigma}}(\mathbf{v}_s), \nabla \mathbf{v}_s)$. Using the definition (1.3),

$$\begin{aligned} (\hat{\boldsymbol{\sigma}}(\mathbf{v}_s), \nabla \mathbf{v}_s) &= (2\mu_s(1-\phi)(\mathcal{D}\mathbf{v}_s - \tfrac{1}{3}\nabla \cdot \mathbf{v}_s \mathbf{I}), \nabla \mathbf{v}_s) \\ &= 2\mu_s \left\{ ((1-\phi)\mathcal{D}\mathbf{v}_s, \mathcal{D}\mathbf{v}_s) - \tfrac{1}{3}((1-\phi)\nabla \cdot \mathbf{v}_s, \nabla \cdot \mathbf{v}_s) \right\}. \end{aligned}$$

We conclude that

$$(\hat{\boldsymbol{\sigma}}(\mathbf{v}_s), \nabla \mathbf{v}_s) \geq C \|\mathcal{D}\mathbf{v}_s\|^2$$

for some positive constant C . An application of Korn's inequality [23, 16] results in

$$(2.31) \quad (\hat{\boldsymbol{\sigma}}(\mathbf{v}_s), \nabla \mathbf{v}_s) \geq C \|\mathcal{D}\mathbf{v}_s\|^2 \geq C_1 \|\mathbf{v}_s\|_1^2.$$

We turn attention to the inf-sup condition, the first condition in (2.28). We recall the inf-sup condition for the Stokes problem [23, 17, 16, 15]. There exists $\gamma_S > 0$ such that for any $w \in \mathbb{W}_0 = L^2(\Omega)/\mathbb{R}$,

$$(2.32) \quad \sup_{\boldsymbol{\psi}_s \in \mathbb{V}_{s,0}} \frac{(w, \nabla \cdot \boldsymbol{\psi}_s)}{\|\boldsymbol{\psi}_s\|_1} \geq \gamma_S \|w\|.$$

We conclude that there is $\mathbf{v}_q \in \mathbb{V}_{s,0}$ normalized so that $\|\mathbf{v}_q\|_1 = \|q\|$ and satisfying

$$(2.33) \quad -(q, \nabla \cdot \mathbf{v}_q) \geq \frac{1}{2} \gamma_S \|q\|^2.$$

For any $\mathbf{U} = (\tilde{\mathbf{v}}_{r,0}, \tilde{q}_f, \mathbf{v}_{s,0}, q) \in \mathbb{X}$, we take the test function in (2.30) to be $\Psi = (\boldsymbol{\psi}_r, w_f, \boldsymbol{\psi}_s, w) \in \mathbb{X}$ defined by

$$(2.34) \quad \begin{aligned} \boldsymbol{\psi}_r &= \tilde{\mathbf{v}}_{r,0}, & w_f &= \tilde{q}_f + \delta_1 \phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} \tilde{\mathbf{v}}_{r,0}), \\ \boldsymbol{\psi}_s &= \mathbf{v}_{s,0} + \delta_2 \mathbf{v}_q, & \text{and } w &= q, \end{aligned}$$

where $\delta_1 > 0$ and $\delta_2 > 0$ will be determined below. After combining and canceling some terms,

$$\begin{aligned} a(\mathbf{U}, \Psi) &= \frac{\mu_f}{k_0} \|\tilde{\mathbf{v}}_{r,0}\|^2 + \delta_1 \|\phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} \tilde{\mathbf{v}}_{r,0})\|^2 + \frac{1}{\mu_s} \left\| \frac{1}{\sqrt{1-\phi}} (\tilde{q}_f - \phi^{1/2} q) \right\|^2 \\ &\quad + (\hat{\boldsymbol{\sigma}}(\mathbf{v}_{s,0}), \nabla \mathbf{v}_{s,0}) - \delta_2 (q, \nabla \cdot \mathbf{v}_q) \\ &\quad + \delta_1 \left(\frac{1}{\mu_s(1-\phi)} (\tilde{q}_f - \phi^{1/2} q), \phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} \tilde{\mathbf{v}}_{r,0}) \right) + \delta_2 (\hat{\boldsymbol{\sigma}}(\mathbf{v}_{s,0}), \nabla \mathbf{v}_q). \end{aligned}$$

There is some $C_2 > 0$ such that

$$(\hat{\boldsymbol{\sigma}}(\mathbf{v}_{s,0}), \nabla \mathbf{v}_q) \leq C_2 \|\mathbf{v}_{s,0}\|_1 \|\mathbf{v}_q\|_1 = C_2 \|\mathbf{v}_{s,0}\|_1 \|q\|,$$

so using (2.31) with its constant $C_1 > 0$ and (2.33), we see that

$$\begin{aligned} a(\mathbf{U}, \Psi) &\geq \frac{\mu_f}{k_0} \|\tilde{\mathbf{v}}_{r,0}\|^2 + \delta_1 \|\phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} \tilde{\mathbf{v}}_{r,0})\|^2 + \frac{1}{\mu_s} \|\tilde{q}_f - \phi^{1/2} q\|^2 \\ &\quad + C_1 \|\mathbf{v}_{s,0}\|_1^2 + \frac{1}{2} \delta_2 \gamma_S \|q\|^2 \\ &\quad + \delta_1 \left(\frac{1}{\mu_s(1-\phi)} (\tilde{q}_f - \phi^{1/2} q), \phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} \tilde{\mathbf{v}}_{r,0}) \right) + \delta_2 (\hat{\boldsymbol{\sigma}}(\mathbf{v}_{s,0}), \nabla \mathbf{v}_q) \\ &\geq \frac{\mu_f}{k_0} \|\tilde{\mathbf{v}}_{r,0}\|^2 + \frac{1}{2} \delta_1 \|\phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} \tilde{\mathbf{v}}_{r,0})\|^2 + \frac{1}{4} \delta_2 \gamma_S \|q\|^2 \\ &\quad + \frac{1}{\mu_s} \left(1 - \frac{\delta_1}{2\mu_s(1-\phi)^2} \right) \|\tilde{q}_f - \phi^{1/2} q\|^2 + \left(C_1 - \delta_2 \frac{C_2^2}{\gamma_S} \right) \|\mathbf{v}_{s,0}\|_1^2. \end{aligned}$$

Taking δ_1 and δ_2 positive but sufficiently small shows that for some $c > 0$,

$$\begin{aligned} a(\mathbf{U}, \Psi) &\geq c \{ \|\tilde{\mathbf{v}}_{r,0}\|^2 + \|\phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} \tilde{\mathbf{v}}_{r,0})\|^2 + \|\mathbf{v}_{s,0}\|_1^2 \\ &\quad + \|\tilde{q}_f - \phi^{1/2} q\|^2 + \|q\|^2 + \|\phi^{1/2+\Theta} \tilde{q}_f\|^2 \}. \end{aligned}$$

Moreover,

$$\|\tilde{q}_f\| \leq \|\tilde{q}_f - \phi^{1/2} q\| + \|q\|,$$

and we have shown the first condition in (2.28). The second follows by symmetry.

We have thus met the conditions of the Babuška–Lax–Milgram theorem, and we conclude that the problem (2.22)–(2.25), (1.3) has a unique solution. Moreover, the bound (2.29) is what is written in Theorem 1. \square

3. The mixed finite element method. Assume Ω is a polygonal domain in one, two, or three dimensions. Let \mathcal{T}_h be a conforming finite element mesh of simplices or rectangular parallelepipeds covering Ω with maximal spacing h , and let \mathcal{E}_h denote the set of element endpoints, edges, or faces.

To continue the exposition, we will restrict ourselves to two dimensions. Extension to one and three dimensions should be clear. Let \mathbb{P}_n denote the space of polynomials of degree n and \mathbb{P}_{n_1, n_2} denote the polynomials of degree n_1 in x and n_2 in z (taking the second coordinate to be the depth z).

3.1. Finite element spaces. For the Darcy part of the system, we choose the lowest order Raviart–Thomas (RT_0) finite element space $\mathbb{V}_{\text{RT}} \times \mathbb{W}_{\text{RT}}$ [31, 17, 33]. On an element $E \in \mathcal{T}_h$, $\mathbb{V}_{\text{RT}}(E) = (\mathbb{P}_0 \times \mathbb{P}_0) \oplus \begin{pmatrix} x \\ z \end{pmatrix} \mathbb{P}_0$ if E is a triangle and $\mathbb{P}_{1,0} \times \mathbb{P}_{0,1}$ if E is a rectangle, and $\mathbb{W}_{\text{RT}}(E) = \mathbb{P}_0$. The degrees of freedom are the normal fluxes on the edges for \mathbb{V}_{RT} , and the average values over the elements for \mathbb{W}_{RT} , i.e.,

$$(3.1) \quad \mathbb{V}_{\text{RT}} = \text{span} \left\{ \mathbf{v}_e : \int_f \mathbf{v}_e \cdot \boldsymbol{\nu}_f ds = \delta_{e,f} \quad \forall e, f \in \mathcal{E}_h \right\},$$

$$(3.2) \quad \mathbb{W}_{\text{RT}} = \text{span} \{ w_E : w_E|_F = \delta_{E,F} \quad \forall E, F \in \mathcal{T}_h \},$$

where $\delta_{i,j}$ is the Kronecker delta function for indices i and j . RT_0 is first order accurate in $H(\text{div}; \Omega) \times L^2(\Omega) \supset \mathbb{V}_{\text{RT}} \times \mathbb{W}_{\text{RT}}$. We could use quadrilateral elements as well, as long as we substitute the Arbogast–Correa (AC_0) space [3] for RT_0 .

For the Stokes part of the system, we could choose any inf-sup stable finite element space $\mathbb{V}_S \times \mathbb{W}_S \subset (H_0^1(\Omega))^2 \times (L^2(\Omega)/\mathbb{R})$. A good choice on rectangular meshes is the Bernardi–Raugel (BR) space $\mathbb{V}_{\text{BR}} \times \mathbb{W}_{\text{BR}}$ [14, 7]. On a rectangular element $E \in \mathcal{T}_h$, $\mathbb{V}_{\text{BR}}(E) = \mathbb{P}_{1,2} \times \mathbb{P}_{2,1}$ and $\mathbb{W}_{\text{BR}}(E) = \mathbb{W}_{\text{RT}}(E)$. BR is first order accurate in $(H^1(\Omega))^2 \times (L^2(\Omega)/\mathbb{R})$. The space was first introduced to solve the Stokes equation, and it has been used to solve Darcy problems with continuous velocities [7]. It is a natural choice for our coupled Darcy–Stokes system, since the convergence rates of the two spaces match.

We could also use standard Taylor–Hood (TH) elements [23, 17, 21]. If $E \in \mathcal{T}_h$ is triangular, $\mathbb{V}_{\text{TH}}(E) = \mathbb{P}_2 \times \mathbb{P}_2$ and $\mathbb{W}_{\text{TH}}(E) = \mathbb{P}_1$ and, if E is rectangular, $\mathbb{V}_{\text{TH}}(E) = \mathbb{P}_{2,2} \times \mathbb{P}_{2,2}$ and $\mathbb{W}_{\text{TH}}(E) = \mathbb{P}_{1,1}$. On rectangular meshes, TH is more accurate than BR, but TH has more degrees of freedom. However, we would not gain any additional overall convergence within the coupled system because of the Darcy part.

3.2. The scaled mixed finite element method. To impose the essential boundary conditions, the extensions \mathbf{g}_r and \mathbf{g}_s are projected into the finite element spaces as $\hat{\mathbf{g}}_r \in \mathbb{V}_{\text{RT}}$ and $\hat{\mathbf{g}}_s \in \mathbb{V}_S$ (\mathbb{V}_{BR} or \mathbb{V}_{TH}) in such a way that the following two compatibility conditions hold:

$$(3.3) \quad \int_{\partial\Omega} (\hat{\mathbf{g}}_s - \mathbf{g}_s) \cdot \boldsymbol{\nu} ds = 0 \quad \text{and} \quad \int_{\partial\Omega} (\phi^{1+\Theta} \hat{\mathbf{g}}_r + \hat{\mathbf{g}}_s) \cdot \boldsymbol{\nu} ds = 0.$$

We also need to define

$$\mathbb{V}_{\text{RT},0} = \{ \mathbf{v} \in \mathbb{V}_{\text{RT}} : \mathbf{v} \cdot \boldsymbol{\nu} = 0 \text{ on } \partial\Omega \},$$

$$\mathbb{V}_{S,0} = \{ \mathbf{v} \in \mathbb{V}_S : \mathbf{v} = \mathbf{0} \text{ on } \partial\Omega \},$$

$$\mathbb{W}_{S,0} = \left\{ w \in \mathbb{W}_S : \int_{\Omega} w dx = 0 \right\}.$$

Scaled mixed finite element method. Find $\tilde{\mathbf{v}}_{r,h} \in \mathbb{V}_{\text{RT},0} + \hat{\mathbf{g}}_r$, $\tilde{q}_{f,h} \in \mathbb{W}_{\text{RT}}$, $\mathbf{v}_{s,h} \in \mathbb{V}_{\text{S},0} + \hat{\mathbf{g}}_s$, and $q_h \in \mathbb{W}_{\text{S},0}$ such that

$$(3.4) \quad \left(\frac{\mu_f}{k_0} \tilde{\mathbf{v}}_{r,h}, \boldsymbol{\psi}_r \right) - (\tilde{q}_{f,h}, \phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} \boldsymbol{\psi}_r)) = 0 \quad \forall \boldsymbol{\psi}_r \in \mathbb{V}_{\text{RT},0},$$

$$(\phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} \tilde{\mathbf{v}}_{r,h}), w_f) + \left(\frac{1}{\mu_s(1-\phi)} (\tilde{q}_{f,h} - \phi^{1/2} q_h), w_f \right) = 0 \quad \forall w_f \in \mathbb{W}_{\text{RT}},$$

$$(3.5) \quad - (q_h, \nabla \cdot \boldsymbol{\psi}_s) + (\hat{\boldsymbol{\sigma}}(\mathbf{v}_{s,h}), \nabla \boldsymbol{\psi}_s) = -((1-\phi)\rho_r \mathbf{g}, \boldsymbol{\psi}_s) \quad \forall \boldsymbol{\psi}_s \in \mathbb{V}_{\text{S},0},$$

$$(3.6) \quad - (q_h, \nabla \cdot \boldsymbol{\psi}_s) + (\hat{\boldsymbol{\sigma}}(\mathbf{v}_{s,h}), \nabla \boldsymbol{\psi}_s) = -((1-\phi)\rho_r \mathbf{g}, \boldsymbol{\psi}_s) \quad \forall \boldsymbol{\psi}_s \in \mathbb{V}_{\text{S},0},$$

$$(3.7) \quad (\nabla \cdot \mathbf{v}_{s,h}, w) - \left(\frac{\phi^{1/2}}{\mu_s(1-\phi)} (\tilde{q}_{f,h} - \phi^{1/2} q_h), w \right) = 0 \quad \forall w \in \mathbb{W}_{\text{S},0},$$

where the term $\hat{\boldsymbol{\sigma}}$ is defined by (1.3). While the scaled finite element method is well-defined when ϕ vanishes due to the condition (2.18), it is important to avoid division by zero in the implementation. One must evaluate the two divergence terms containing ϕ to a negative power in (3.4)–(3.5) at quadrature points. Because the divergence terms scale with ϕ to the overall power $1/2 + \Theta > 0$, these terms should be set to zero when ϕ vanishes. That is, at a quadrature point where $\phi = 0$, take the value of the entire term to be zero at that point.

LEMMA 3. *If (2.18) holds, then there exists a unique solution to the scaled mixed finite element method (3.4)–(3.7).*

Proof. The scaled method gives rise to a square linear system when restricted to bases for the finite element spaces, so existence of a solution is equivalent to uniqueness. To show uniqueness, set to zero the quantities $\hat{\mathbf{g}}_r$, $\hat{\mathbf{g}}_s$, and \mathbf{g} . The test functions

$$\boldsymbol{\psi}_r = \tilde{\mathbf{v}}_{r,h} \in \mathbb{V}_{\text{RT},0}, \quad w_f = \tilde{q}_{f,h} \in \mathbb{W}_{\text{RT}}, \quad \boldsymbol{\psi}_s = \mathbf{v}_{s,h} \in \mathbb{V}_{\text{S},0}, \quad \text{and} \quad w = q_h \in \mathbb{W}_{\text{S},0},$$

when substituted into (3.4)–(3.7) and after the equations are added, imply that

$$\frac{\mu_f}{k_0} \|\tilde{\mathbf{v}}_{r,h}\|^2 + \left(\frac{1}{\mu_s(1-\phi)} (\tilde{q}_{f,h} - \phi^{1/2} q_h), \tilde{q}_{f,h} - \phi^{1/2} q_h \right) + (\hat{\boldsymbol{\sigma}}(\mathbf{v}_{s,h}), \nabla \mathbf{v}_{s,h}) = 0.$$

Thus $\tilde{\mathbf{v}}_{r,h} = 0$, the estimate (2.31) shows $\mathbf{v}_{s,h} = 0$, and $\tilde{q}_{f,h} = \phi^{1/2} q_h$.

The discrete version of the inf-sup condition (2.32) holds for BR and TH Stokes elements with a possibly smaller constant $0 < \gamma_S^* \leq \gamma_S$ independent of h . Therefore there is some $\mathbf{v}_{q,h} \in \mathbb{V}_{\text{S},0}$ such that $\|\mathbf{v}_{q,h}\|_1 = \|q_h\|$ and

$$(3.8) \quad - (q_h, \nabla \cdot \mathbf{v}_{q,h}) \geq \frac{1}{2} \gamma_S^* \|q_h\|^2.$$

The choice $\boldsymbol{\psi}_s = \mathbf{v}_{q,h}$ in (3.6) shows that $q_h = 0$, and so also $\tilde{q}_{h,f} = \phi^{1/2} q_h = 0$. \square

3.3. Convergence of the scaled method. To derive a bound for the error, we first take the difference of (2.22)–(2.25) and (3.4)–(3.7) and add the resulting

equations to see that

$$\begin{aligned}
 & \left(\frac{\mu_f}{k_0} (\tilde{\mathbf{v}}_r - \tilde{\mathbf{v}}_{r,h}), \boldsymbol{\psi}_r \right) - (\tilde{q}_f - \tilde{q}_{f,h}, \phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} \boldsymbol{\psi}_r)) \\
 & + (\phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} (\tilde{\mathbf{v}}_r - \tilde{\mathbf{v}}_{r,h})), w_f) \\
 & + \left(\frac{1}{\mu_s(1-\phi)} (\tilde{q}_f - \tilde{q}_{f,h} - \phi^{1/2}(q - q_h)), w_f - \phi^{1/2}w \right) \\
 (3.9) \quad & - (q - q_h, \nabla \cdot \boldsymbol{\psi}_s) + (\hat{\boldsymbol{\sigma}}(\mathbf{v}_s - \mathbf{v}_{s,h}), \nabla \boldsymbol{\psi}_s) + (\nabla \cdot (\mathbf{v}_s - \mathbf{v}_{s,h}), w) = 0
 \end{aligned}$$

for any $\boldsymbol{\psi}_r \in \mathbb{V}_{\text{RT},0}$, $w_f \in \mathbb{W}_{\text{RT}}$, $\boldsymbol{\psi}_s \in \mathbb{V}_{\text{S},0}$, and $w \in \mathbb{W}_{\text{S},0}$.

Before defining our choice of test functions, we need the usual projection operators associated with RT_0 (or AC_0). Let $\mathcal{P}_{\mathbb{W}_{\text{RT}}} : L^2(\Omega) \rightarrow \mathbb{W}_{\text{RT}}$ denote the $L^2(\Omega)$ -projection operator mapping onto the space of piecewise constant functions \mathbb{W}_{RT} . Let $\pi_{\text{RT}} : H(\text{div}; \Omega) \cap L^{2+\epsilon}(\Omega) \rightarrow \mathbb{V}_{\text{RT}}$ (any $\epsilon > 0$) denote the standard RT or Fortin operator that preserves element average divergence and average edge normal fluxes [31, 17, 33, 3]. We also need the usual $H^1(\Omega)$ -projection $\pi_{\text{S}} : H^1(\Omega) \rightarrow \mathbb{V}_{\text{S}}$ and the $L^2(\Omega)$ -projection $\mathcal{P}_{\mathbb{W}_{\text{S}}} : L^2(\Omega) \rightarrow \mathbb{W}_{\text{S}}$.

Let the function $\mathbf{v}_{q,h} \in \mathbb{V}_{\text{S},0}$ arise from the discrete version of the inf-sup condition for Stokes (2.33) (as in (3.8)), normalized so that $\|\mathbf{v}_{q,h}\|_1 = \|\mathcal{P}_{\mathbb{W}_{\text{S}}}q - q_h\|$ and satisfying

$$(3.10) \quad -(\mathcal{P}_{\mathbb{W}_{\text{S}}}q - q_h, \nabla \cdot \mathbf{v}_{q,h}) \geq \frac{1}{2} \gamma_{\text{S}}^* \|\mathcal{P}_{\mathbb{W}_{\text{S}}}q - q_h\|^2.$$

Similarly to the test functions taken in (2.34), we take

$$\begin{aligned}
 \boldsymbol{\psi}_r &= (\tilde{\mathbf{v}}_r - \tilde{\mathbf{v}}_{r,h}) - (\tilde{\mathbf{v}}_r - \pi_{\text{RT}} \tilde{\mathbf{v}}_r) - (\pi_{\text{RT}} \mathbf{g}_r - \hat{\mathbf{g}}_r) && \in \mathbb{V}_{\text{RT},0}, \\
 w_f &= \mathcal{P}_{\mathbb{W}_{\text{RT}}} \tilde{q}_f - \tilde{q}_{f,h} + \delta_1 \mathcal{P}_{\mathbb{W}_{\text{RT}}} [\phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} (\tilde{\mathbf{v}}_r - \tilde{\mathbf{v}}_{r,h}))] && \in \mathbb{W}_{\text{RT}}, \\
 \boldsymbol{\psi}_s &= (\mathbf{v}_s - \mathbf{v}_{s,h}) - (\mathbf{v}_s - \pi_{\text{S}} \mathbf{v}_s) - (\pi_{\text{S}} \mathbf{g}_s - \hat{\mathbf{g}}_s) + \delta_2 \mathbf{v}_{q,h} && \in \mathbb{V}_{\text{S},0}, \\
 w &= \mathcal{P}_{\mathbb{W}_{\text{S}}}q - q_h && \in \mathbb{W}_{\text{S},0},
 \end{aligned}$$

where $\delta_1 > 0$ and $\delta_2 > 0$ will be determined below. We remark that the term multiplying δ_1 must be projected back into the discrete space, and so our derivation is not completely straightforward.

Introducing $\mathcal{P}_{\mathbb{W}_{\text{RT}}}$ thrice into (3.9) yields

$$\begin{aligned}
 & \left(\frac{\mu_f}{k_0} (\tilde{\mathbf{v}}_r - \tilde{\mathbf{v}}_{r,h}), \boldsymbol{\psi}_r \right) - (\mathcal{P}_{\mathbb{W}_{\text{RT}}} \tilde{q}_f - \tilde{q}_{f,h}, \mathcal{P}_{\mathbb{W}_{\text{RT}}} [\phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} \boldsymbol{\psi}_r)]) \\
 & + (\mathcal{P}_{\mathbb{W}_{\text{RT}}} [\phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} (\tilde{\mathbf{v}}_r - \tilde{\mathbf{v}}_{r,h}))], w_f) \\
 & - (\tilde{q}_f - \mathcal{P}_{\mathbb{W}_{\text{RT}}} \tilde{q}_f, \phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} \boldsymbol{\psi}_r)) \\
 & + \left(\frac{1}{\mu_s(1-\phi)} (\tilde{q}_f - \tilde{q}_{f,h} - \phi^{1/2}(q - q_h)), w_f - \phi^{1/2}w \right) \\
 & - (q - q_h, \nabla \cdot \boldsymbol{\psi}_s) + (\hat{\boldsymbol{\sigma}}(\mathbf{v}_s - \mathbf{v}_{s,h}), \nabla \boldsymbol{\psi}_s) + (\nabla \cdot (\mathbf{v}_s - \mathbf{v}_{s,h}), w) \\
 (3.11) \quad & = T_1 + \cdots + T_8 \text{ (respectively)} = 0.
 \end{aligned}$$

For the first term in (3.11), we deduce that for some generic constant $C > 0$,

$$\begin{aligned} T_1 &= \left(\frac{\mu_f}{k_0} (\tilde{\mathbf{v}}_r - \tilde{\mathbf{v}}_{r,h}), \boldsymbol{\psi}_r \right) \\ &= \frac{\mu_f}{k_0} \|\tilde{\mathbf{v}}_r - \tilde{\mathbf{v}}_{r,h}\|^2 - \left(\frac{\mu_f}{k_0} (\tilde{\mathbf{v}}_r - \tilde{\mathbf{v}}_{r,h}), \tilde{\mathbf{v}}_r - \pi_{\text{RT}} \tilde{\mathbf{v}}_r + \pi_{\text{RT}} \mathbf{g}_r - \hat{\mathbf{g}}_r \right) \\ &\geq \frac{\mu_f}{2k_0} \|\tilde{\mathbf{v}}_r - \tilde{\mathbf{v}}_{r,h}\|^2 - C \{ \|\tilde{\mathbf{v}}_r - \pi_{\text{RT}} \tilde{\mathbf{v}}_r\|^2 + \|\pi_{\text{RT}} \mathbf{g}_r - \hat{\mathbf{g}}_r\|^2 \}. \end{aligned}$$

For the next two terms, for any $\epsilon > 0$,

$$\begin{aligned} T_2 + T_3 &= -(\mathcal{P}_{\text{WRT}} \tilde{q}_f - \tilde{q}_{f,h}, \mathcal{P}_{\text{WRT}} [\phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} \boldsymbol{\psi}_r)]) \\ &\quad + (\mathcal{P}_{\text{WRT}} [\phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} (\tilde{\mathbf{v}}_r - \tilde{\mathbf{v}}_{r,h}))], w_f) \\ &= (\mathcal{P}_{\text{WRT}} \tilde{q}_f - \tilde{q}_{f,h}, \mathcal{P}_{\text{WRT}} [\phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} (\tilde{\mathbf{v}}_r - \pi_{\text{RT}} \tilde{\mathbf{v}}_r + \pi_{\text{RT}} \mathbf{g}_r - \hat{\mathbf{g}}_r)]) \\ &\quad + \delta_1 \|\mathcal{P}_{\text{WRT}} [\phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} (\tilde{\mathbf{v}}_r - \tilde{\mathbf{v}}_{r,h}))]\|^2 \\ &\geq \delta_1 \|\mathcal{P}_{\text{WRT}} [\phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} (\tilde{\mathbf{v}}_r - \tilde{\mathbf{v}}_{r,h}))]\|^2 - \epsilon \|\mathcal{P}_{\text{WRT}} \tilde{q}_f - \tilde{q}_{f,h}\|^2 \\ &\quad - C \{ \|\mathcal{P}_{\text{WRT}} [\phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} (\tilde{\mathbf{v}}_r - \pi_{\text{RT}} \tilde{\mathbf{v}}_r))]\|^2 \\ &\quad + \|\mathcal{P}_{\text{WRT}} [\phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} (\pi_{\text{RT}} \mathbf{g}_r - \hat{\mathbf{g}}_r))]\|^2 \}. \end{aligned}$$

Skipping T_4 for the moment, the next term is

$$\begin{aligned} T_5 &= \left(\frac{1}{\mu_s(1-\phi)} (\tilde{q}_f - \tilde{q}_{f,h} - \phi^{1/2}(q - q_h)), w_f - \phi^{1/2}w \right) \\ &= \left(\frac{1}{\mu_s(1-\phi)} (\tilde{q}_f - \tilde{q}_{f,h} - \phi^{1/2}(q - q_h)), \tilde{q}_f - \tilde{q}_{f,h} - \phi^{1/2}(q - q_h) \right) \\ &\quad - \left(\frac{1}{\mu_s(1-\phi)} (\tilde{q}_f - \tilde{q}_{f,h} - \phi^{1/2}(q - q_h)), \tilde{q}_f - \mathcal{P}_{\text{WRT}} \tilde{q}_f - \phi^{1/2}(q - \mathcal{P}_{\text{WS}} q) \right) \\ &\quad + \delta_1 \left(\frac{1}{\mu_s(1-\phi)} (\tilde{q}_f - \tilde{q}_{f,h} - \phi^{1/2}(q - q_h)), \mathcal{P}_{\text{WRT}} [\phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} (\tilde{\mathbf{v}}_r - \tilde{\mathbf{v}}_{r,h}))] \right) \\ &\geq \frac{1}{2\mu_s(1-\phi^*)} \|\tilde{q}_f - \tilde{q}_{f,h} - \phi^{1/2}(q - q_h)\|^2 \\ &\quad - C \{ \|\tilde{q}_f - \mathcal{P}_{\text{WRT}} \tilde{q}_f\|^2 + \|q - \mathcal{P}_{\text{WS}} q\|^2 \\ &\quad + \delta_1^2 \|\mathcal{P}_{\text{WRT}} [\phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} (\tilde{\mathbf{v}}_r - \tilde{\mathbf{v}}_{r,h}))]\|^2 \}. \end{aligned}$$

Noting that $w = q - q_h - (q - \mathcal{P}_{\text{WS}} q)$, the sixth and eighth terms satisfy

$$\begin{aligned} T_6 + T_8 &= -(q - q_h, \nabla \cdot \boldsymbol{\psi}_s) + (\nabla \cdot (\mathbf{v}_s - \mathbf{v}_{s,h}), w) \\ &= (q - q_h, \nabla \cdot (\mathbf{v}_s - \pi_{\text{S}} \mathbf{v}_s + \pi_{\text{S}} \mathbf{g}_s - \hat{\mathbf{g}}_s)) - (\nabla \cdot (\mathbf{v}_s - \mathbf{v}_{s,h}), q - \mathcal{P}_{\text{WS}} q) \\ &\quad - \delta_2 [(\mathcal{P}_{\text{WS}} q - q_h, \nabla \cdot \mathbf{v}_{q,h}) + (q - \mathcal{P}_{\text{WS}} q, \nabla \cdot \mathbf{v}_{q,h})]. \end{aligned}$$

Recalling (3.10) and that $\|\mathbf{v}_{q,h}\|_1 = \|\mathcal{P}_{\text{WS}} q - q_h\|$, we have

$$\begin{aligned} T_6 + T_8 &\geq \frac{1}{4} \delta_2 \gamma_{\text{S}}^* \|\mathcal{P}_{\text{WS}} q - q_h\|^2 - \frac{1}{8} \delta_2 \gamma_{\text{S}}^* \|q - q_h\|^2 - \epsilon \|\nabla \cdot (\mathbf{v}_s - \mathbf{v}_{s,h})\|^2 \\ &\quad - C \{ \|q - \mathcal{P}_{\text{WS}} q\|^2 + \|\nabla \cdot (\mathbf{v}_s - \pi_{\text{S}} \mathbf{v}_s)\|^2 + \|\nabla \cdot (\pi_{\text{S}} \mathbf{g}_s - \hat{\mathbf{g}}_s)\|^2 \} \\ &\geq \frac{1}{8} \delta_2 \gamma_{\text{S}}^* \|q - q_h\|^2 - \epsilon \|\nabla \cdot (\mathbf{v}_s - \mathbf{v}_{s,h})\|^2 \\ &\quad - C \{ \|q - \mathcal{P}_{\text{WS}} q\|^2 + \|\nabla \cdot (\mathbf{v}_s - \pi_{\text{S}} \mathbf{v}_s)\|^2 + \|\nabla \cdot (\pi_{\text{S}} \mathbf{g}_s - \hat{\mathbf{g}}_s)\|^2 \}. \end{aligned}$$

Finally, for the next to last term, note that $\boldsymbol{\psi}_s = \pi_S \mathbf{v}_s - \mathbf{v}_{s,h} - \pi_S \mathbf{g}_s + \hat{\mathbf{g}}_s + \delta_2 \mathbf{v}_{q,h}$, so we have from (2.31) that

$$\begin{aligned} T_7 &= (\hat{\boldsymbol{\sigma}}(\mathbf{v}_s - \mathbf{v}_{s,h}), \nabla \boldsymbol{\psi}_s) \\ &= (\hat{\boldsymbol{\sigma}}(\pi_S \mathbf{v}_s - \mathbf{v}_{s,h} - \pi_S \mathbf{g}_s + \hat{\mathbf{g}}_s), \nabla \boldsymbol{\psi}_s) + (\hat{\boldsymbol{\sigma}}(\mathbf{v}_s - \pi_S \mathbf{v}_s + \pi_S \mathbf{g}_s - \hat{\mathbf{g}}_s), \nabla \boldsymbol{\psi}_s) \\ &\geq C_1 \|\pi_S \mathbf{v}_s - \mathbf{v}_{s,h} - \pi_S \mathbf{g}_s + \hat{\mathbf{g}}_s\|_1^2 - \frac{1}{2} C_1 \|\pi_S \mathbf{v}_s - \mathbf{v}_{s,h}\|_1^2 \\ &\quad - C \{ \|\mathbf{v}_s - \pi_S \mathbf{v}_s\|_1^2 + \|\pi_S \mathbf{g}_s - \hat{\mathbf{g}}_s\|_1^2 + \delta_2^2 \|\mathbf{v}_{q,h}\|_1^2 \} \\ &\geq \frac{1}{2} C_1 \|\mathbf{v}_s - \mathbf{v}_{s,h}\|_1^2 - C \{ \|\mathbf{v}_s - \pi_S \mathbf{v}_s\|_1^2 + \|\pi_S \mathbf{g}_s - \hat{\mathbf{g}}_s\|_1^2 + \delta_2^2 \|\mathcal{P}_{\mathbb{W}_S} q - q_h\|^2 \}. \end{aligned}$$

We turn now to the fourth term T_4 in (3.11), which we estimate similarly to a term in [6, section 7] for the degenerate Darcy system. That is, we introduce the projection $I - \mathcal{P}_{\mathbb{W}_{\text{RT}}}$ and compute as follows:

$$\begin{aligned} -T_4 &= (\tilde{q}_f - \mathcal{P}_{\mathbb{W}_{\text{RT}}} \tilde{q}_f, \phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} \boldsymbol{\psi}_r)) \\ &= (\tilde{q}_f - \mathcal{P}_{\mathbb{W}_{\text{RT}}} \tilde{q}_f, (I - \mathcal{P}_{\mathbb{W}_{\text{RT}}}) \phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} \boldsymbol{\psi}_r)) \\ &= (\tilde{q}_f - \mathcal{P}_{\mathbb{W}_{\text{RT}}} \tilde{q}_f, (I - \mathcal{P}_{\mathbb{W}_{\text{RT}}}) \phi^{1/2+\Theta} \nabla \cdot \boldsymbol{\psi}_r + (I - \mathcal{P}_{\mathbb{W}_{\text{RT}}}) (1 + \Theta) \phi^{\Theta-1/2} \nabla \phi \cdot \boldsymbol{\psi}_r) \\ &\leq C \|\tilde{q}_f - \mathcal{P}_{\mathbb{W}_{\text{RT}}} \tilde{q}_f\| \{ \|(I - \mathcal{P}_{\mathbb{W}_{\text{RT}}}) \phi^{1/2+\Theta} \nabla \cdot \boldsymbol{\psi}_r\| + \|\boldsymbol{\psi}_r\| \}, \end{aligned}$$

since we have assumed the bound (2.18) on the term $\phi^{\Theta-1/2} \nabla \phi$. Because $\nabla \cdot \boldsymbol{\psi}_r$ is piecewise constant, we have that

$$\begin{aligned} \|(I - \mathcal{P}_{\mathbb{W}_{\text{RT}}}) \phi^{1/2+\Theta} \nabla \cdot \boldsymbol{\psi}_r\| &\leq \|(I - \mathcal{P}_{\mathbb{W}_{\text{RT}}}) \phi^{1/2+\Theta}\|_{L^\infty(\Omega)} \|\nabla \cdot \boldsymbol{\psi}_r\| \\ &\leq Ch \|\phi^{1/2+\Theta}\|_{W^{1,\infty}(\Omega)} \|\nabla \cdot \boldsymbol{\psi}_r\| \\ &\leq Ch \|\nabla \cdot \boldsymbol{\psi}_r\|, \end{aligned}$$

using [20] for the approximation of the L^2 -projection in L^∞ and (2.18) again. If we assume that the mesh is quasi-uniform, then we can remove the divergence operator in the final expression at the expense of a power of the mesh spacing h . Thus we have

$$\begin{aligned} -T_4 &\leq C \|\tilde{q}_f - \mathcal{P}_{\mathbb{W}_{\text{RT}}} \tilde{q}_f\| \|\boldsymbol{\psi}_r\| \\ &\leq \epsilon \|\tilde{\mathbf{v}}_r - \tilde{\mathbf{v}}_{r,h}\|^2 + C \{ \|\tilde{q}_f - \mathcal{P}_{\mathbb{W}_{\text{RT}}} \tilde{q}_f\|^2 + \|\tilde{\mathbf{v}}_r - \pi_{\text{RT}} \tilde{\mathbf{v}}_r\|^2 + \|\pi_{\text{RT}} \mathbf{g}_r - \hat{\mathbf{g}}_r\|^2 \}. \end{aligned}$$

Combining these estimates results in

$$\begin{aligned} &\frac{\mu_f}{2k_0} \|\tilde{\mathbf{v}}_r - \tilde{\mathbf{v}}_{r,h}\|^2 + \delta_1 \|\mathcal{P}_{\mathbb{W}_{\text{RT}}} [\phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} (\tilde{\mathbf{v}}_r - \tilde{\mathbf{v}}_{r,h}))]\|^2 + \frac{1}{2} C_1 \|\mathbf{v}_s - \mathbf{v}_{s,h}\|_1^2 \\ &\quad + \frac{1}{2\mu_s(1-\phi^*)} \|\tilde{q}_f - \tilde{q}_{f,h} - \phi^{1/2}(q - q_h)\|^2 + \frac{1}{4} \delta_2 \gamma_S^* \|q - q_h\|^2 \\ &\leq \epsilon \{ \|\mathcal{P}_{\mathbb{W}_{\text{RT}}} \tilde{q}_f - \tilde{q}_{f,h}\|^2 + \|\nabla \cdot (\mathbf{v}_s - \mathbf{v}_{s,h})\|^2 + \|\tilde{\mathbf{v}}_r - \tilde{\mathbf{v}}_{r,h}\|^2 \} \\ &\quad + C \{ \|\tilde{\mathbf{v}}_r - \pi_{\text{RT}} \tilde{\mathbf{v}}_r\|^2 + \|\mathcal{P}_{\mathbb{W}_{\text{RT}}} [\phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} (\tilde{\mathbf{v}}_r - \pi_{\text{RT}} \tilde{\mathbf{v}}_r))]\|^2 \\ &\quad + \|\pi_{\text{RT}} \mathbf{g}_r - \hat{\mathbf{g}}_r\|^2 + \|\mathcal{P}_{\mathbb{W}_{\text{RT}}} [\phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} (\pi_{\text{RT}} \mathbf{g}_r - \hat{\mathbf{g}}_r))]\|^2 \\ &\quad + \|\tilde{q}_f - \mathcal{P}_{\mathbb{W}_{\text{RT}}} \tilde{q}_f\|^2 + \|q - \mathcal{P}_{\mathbb{W}_S} q\|^2 \\ &\quad + \delta_1^2 \|\mathcal{P}_{\mathbb{W}_{\text{RT}}} [\phi^{-1/2} \nabla \cdot (\phi^{1+\Theta} (\tilde{\mathbf{v}}_r - \tilde{\mathbf{v}}_{r,h}))]\|^2 \\ &\quad + \|\mathbf{v}_s - \pi_S \mathbf{v}_s\|_1^2 + \|\pi_S \mathbf{g}_s - \hat{\mathbf{g}}_s\|_1^2 + \delta_2^2 \|\mathcal{P}_{\mathbb{W}_S} q - q_h\|^2 \}. \end{aligned} \tag{3.12}$$

Note that

$$\|\tilde{q}_f - \tilde{q}_{f,h}\|^2 \leq \|\tilde{q}_f - \tilde{q}_{f,h} - \phi^{1/2}(q - q_h)\|^2 + \|q - q_h\|^2.$$

Therefore, if we take ϵ , δ_1 , and δ_2 small enough, we have proven the following theorem.

THEOREM 4. *Assume that (2.18) holds on the porosity, $0 \leq \phi \leq \phi^* < 1$, the mesh is quasi-uniform, and the extensions $\mathbf{g}_r \in \tilde{\mathbb{V}}_r$ and $\mathbf{g}_s \in \mathbb{V}_s$ satisfy (2.21) and their approximations satisfy (3.3). Then the difference of the solution to the scaled formulation (2.22)–(2.25), (1.3), and its finite element approximation satisfy*

$$\begin{aligned} & \|\tilde{\mathbf{v}}_r - \tilde{\mathbf{v}}_{r,h}\| + \|\mathcal{P}_{\mathbb{W}_{\text{RT}}}[\phi^{-1/2}\nabla \cdot (\phi^{1+\Theta}(\tilde{\mathbf{v}}_r - \tilde{\mathbf{v}}_{r,h}))]\| + \|\mathbf{v}_s - \mathbf{v}_{s,h}\|_1 \\ & + \|\tilde{q}_f - \tilde{q}_{f,h}\| + \|q - q_h\| \\ & \leq C\{\|\tilde{\mathbf{v}}_r - \pi_{\text{RT}}\tilde{\mathbf{v}}_r\| + \|\mathcal{P}_{\mathbb{W}_{\text{RT}}}[\phi^{-1/2}\nabla \cdot (\phi^{1+\Theta}(\tilde{\mathbf{v}}_r - \pi_{\text{RT}}\tilde{\mathbf{v}}_r))]\| \\ & + \|\pi_{\text{RT}}\mathbf{g}_r - \hat{\mathbf{g}}_r\| + \|\mathcal{P}_{\mathbb{W}_{\text{RT}}}[\phi^{-1/2}\nabla \cdot (\phi^{1+\Theta}(\pi_{\text{RT}}\mathbf{g}_r - \hat{\mathbf{g}}_r))]\| \\ (3.13) \quad & + \|\tilde{q}_f - \mathcal{P}_{\mathbb{W}_{\text{RT}}}\tilde{q}_f\| + \|q - \mathcal{P}_{\mathbb{W}_S}q\| + \|\mathbf{v}_s - \pi_S\mathbf{v}_s\|_1 + \|\pi_S\mathbf{g}_s - \hat{\mathbf{g}}_s\|_1\}. \end{aligned}$$

If the solution is sufficiently smooth, this bound implies first order convergence. It also implies stability of the scheme even when the solution is not very smooth.

4. A modification for local mass conservation. As in [5], we define a locally mass conservative implementation of the scaled method by using the quantity $\hat{\phi} = \mathcal{P}_{\mathbb{W}_{\text{RT}}}\phi \in \mathbb{W}_{\text{RT}}$ given by taking the average over each element $E \in \mathcal{T}_h$, i.e.,

$$(4.1) \quad \text{for } \mathbf{x} \in E, \quad \hat{\phi}(\mathbf{x}) = \hat{\phi}_E = \frac{1}{|E|} \int_E \phi \, dx,$$

where $|E|$ is the area of E . When $\hat{\phi}|_E = \hat{\phi}_E = 0$ vanishes on an element E , ϕ is identically zero on E . We modify the two divergence terms in the scaled MFEM (3.4)–(3.7) by replacing

$$\phi^{-1/2}\nabla \cdot (\phi^{1+\Theta}\mathbf{v})|_E \quad \text{by} \quad \begin{cases} \hat{\phi}_E^{-1/2}\nabla \cdot (\phi^{1+\Theta}\mathbf{v}) & \text{if } \hat{\phi}_E \neq 0, \\ 0 & \text{if } \hat{\phi}_E = 0. \end{cases}$$

We also modify the two terms in (3.5) and (3.7) involving the pressure potentials. These changes make the method locally mass conservative, as we show later.

Locally conservative scaled MFEM. Find $\tilde{\mathbf{v}}_{r,h} \in \mathbb{V}_{\text{RT},0} + \hat{\mathbf{g}}_r$, $\tilde{q}_{f,h} \in \mathbb{W}_{\text{RT}}$, $\mathbf{v}_{s,h} \in \mathbb{V}_{S,0} + \hat{\mathbf{g}}_s$, and $q_h \in \mathbb{W}_{S,0}$ such that

$$(4.2) \quad \left(\frac{\mu_f}{k_0} \tilde{\mathbf{v}}_{r,h}, \boldsymbol{\psi}_r \right) - (\tilde{q}_{f,h}, \hat{\phi}^{-1/2}\nabla \cdot (\phi^{1+\Theta}\boldsymbol{\psi}_r)) = 0 \quad \forall \boldsymbol{\psi}_r \in \mathbb{V}_{\text{RT},0},$$

$$(\hat{\phi}^{-1/2}\nabla \cdot (\phi^{1+\Theta}\tilde{\mathbf{v}}_{r,h}), w_f) + \left(\frac{\phi\hat{\phi}^{-1}}{\mu_s(1-\phi)} (\tilde{q}_{f,h} - \hat{\phi}^{1/2}q_h), w_f \right) = 0 \quad \forall w_f \in \mathbb{W}_{\text{RT}},$$

$$(4.3) \quad - (q_h, \nabla \cdot \boldsymbol{\psi}_s) + (\hat{\boldsymbol{\sigma}}(\mathbf{v}_{s,h}), \nabla \boldsymbol{\psi}_s) = -((1-\phi)\rho_r \mathbf{g}, \boldsymbol{\psi}_s) \quad \forall \boldsymbol{\psi}_s \in \mathbb{V}_{S,0},$$

$$(4.4) \quad (\nabla \cdot \mathbf{v}_{s,h}, w) - \left(\frac{\phi\hat{\phi}^{-1/2}}{\mu_s(1-\phi)} (\tilde{q}_{f,h} - \hat{\phi}^{1/2}q_h), w \right) = 0 \quad \forall w \in \mathbb{W}_{S,0},$$

where the term $\hat{\sigma}$ is defined by (1.3). On an element $E \in \mathcal{T}_h$, where $\hat{\phi}_E = 0$, the three terms in (4.2)–(4.3), (4.5) involving $\hat{\phi}^{-1/2}$ are set to zero, and in the second term in (4.3), we interpret $\hat{\phi}\hat{\phi}^{-1}$ as one. Furthermore, we define the discrete Darcy velocity $\mathbf{u}_h \in \mathbb{V}_{\text{RT}}$ and fluid potential $q_{f,h} \in \mathbb{W}_{\text{RT}}$ by their degrees of freedom:

$$(4.6) \quad \mathbf{u}_h \cdot \nu|_e = \frac{1}{|e|} \int_e \phi^{1+\Theta} ds \, \tilde{\mathbf{v}}_{r,h} \cdot \nu|_e \quad \forall e \in \mathcal{E}_h,$$

$$(4.7) \quad q_{f,h}|_E = \hat{\phi}_E^{-1/2} \tilde{q}_{f,h}|_E \quad \forall E \in \mathcal{T}_h,$$

wherein we arbitrarily set $q_{f,h}|_E = 0$ if $\hat{\phi}_E = 0$.

The existence of a unique solution can be shown in a way completely analogous to that for the nonconservative scaled MFEM in section 3. Moreover, one can show that the locally conservative MFEM is stable. We have no proof of convergence of the locally conservative method at this time, but the numerical results show optimal convergence and even superconvergence.

To see local mass conservation of the fluid, let $E \in \mathcal{T}_h$ be any element. With w_E defined in (3.2), the test function $w_f = \hat{\phi}_E^{1/2} w_E \in \mathbb{W}_{\text{RT}}$ in (4.3) gives

$$\mu_s \int_E \nabla \cdot (\phi^{1+\Theta} \tilde{\mathbf{v}}_{r,h}) dx + \int_E \frac{\phi}{1-\phi} (\hat{\phi}^{-1/2} \tilde{q}_{f,h} - q_h) dx = 0.$$

Since $\tilde{\mathbf{v}}_r \cdot \nu$ and $\mathbf{u}_h \cdot \nu$ are constant on each edge $e \subset \partial E$, we see from (4.6) that

$$\int_E \nabla \cdot (\phi^{1+\Theta} \tilde{\mathbf{v}}_{r,h}) dx = \int_{\partial E} \phi^{1+\Theta} \tilde{\mathbf{v}}_{r,h} \cdot \nu ds = \int_{\partial E} \mathbf{u}_h \cdot \nu ds = \int_E \nabla \cdot \mathbf{u}_h dx.$$

The definition of $q_{f,h}$ (4.7) gives

$$(4.8) \quad \mu_s \int_E \nabla \cdot \mathbf{u}_h dx + \int_E \frac{\phi}{1-\phi} (q_{f,h} - q_h) dx = 0,$$

which is local mass conservation, i.e., (2.5) holds locally.

We obtain local mass conservation of the solid matrix if we use BR spaces. In that case, we can take the test function $w = w_E \in \mathbb{W}_{\text{BR}}$ in (4.5) to see

$$\mu_s \int_E \nabla \cdot \mathbf{v}_{s,h} dx - \int_E \frac{\phi}{1-\phi} (q_{f,h} - q_{s,h}) dx = 0,$$

which is (2.7) holding locally.

5. Implementation of the methods on rectangular meshes. The linear system corresponding to either of the methods (3.4)–(3.7) or (4.2)–(4.5) has the form

$$(5.1) \quad \begin{pmatrix} A & -B_\phi & 0 & 0 \\ B_\phi^T & C_{f,\phi} & 0 & -C_{f,s,\phi} \\ 0 & 0 & D_\phi & -B \\ 0 & -C_{f,s,\phi}^T & B^T & C_{s,\phi} \end{pmatrix} \begin{pmatrix} \tilde{v}_r \\ \tilde{q}_f \\ v_s \\ q \end{pmatrix} = \begin{pmatrix} a_\phi \\ 0 \\ b_\phi \\ 0 \end{pmatrix},$$

wherein the solution represents the degrees of freedom of $\tilde{\mathbf{v}}_{r,h}$, $\tilde{q}_{f,h}$, $\mathbf{v}_{s,h}$, and q_h with respect to the bases of the finite element spaces. We remark only on the evaluation of B_ϕ and D_ϕ . To avoid approximating derivatives of ϕ , the matrix B_ϕ should be

computed using the divergence theorem. For the locally conservative method, for any element $E \in \mathcal{T}_h$ and edge $e \in \mathcal{E}_h$,

$$\begin{aligned} B_{\phi,e,E} &= (\hat{\phi}^{-1/2} \nabla \cdot (\phi^{1+\Theta} \mathbf{v}_e), w_E) \\ &= \begin{cases} \hat{\phi}_E^{-1/2} \int_e \phi^{1+\Theta} ds \mathbf{v}_e \cdot \nu_E & \text{if } e \subset \partial E \text{ and } \hat{\phi}_E \neq 0, \\ 0 & \text{if } e \not\subset \partial E \text{ or } \hat{\phi}_E = 0. \end{cases} \end{aligned}$$

A similar expression is used for the nonconservative scaled MFEM of section 3. The matrix D_ϕ is symmetric, and the (k, ℓ) entry is computed using (1.3) as

$$D_{\phi,k,\ell} = (\hat{\sigma}(\mathbf{v}_{s,k}), \nabla \psi_{s,\ell}) = 2\mu_s [((1-\phi) \mathcal{D} \mathbf{v}_{s,k}, \mathcal{D} \psi_{s,\ell}) - \frac{1}{3} ((1-\phi) \nabla \cdot \mathbf{v}_{s,k}, \nabla \cdot \psi_{s,\ell})].$$

We can simplify the implementation when Ω is a union of rectangular subdomains in one, two, or three dimensions, and \mathcal{T}_h is a rectangular finite element mesh. We modify either method by approximating the first integral in (3.4) or (4.2) using what is known as mass lumping. The integral is approximated by a trapezoidal quadrature rule $(\cdot, \cdot)_Q$, so that for any two edges $e, f \in \mathcal{E}_h$,

$$(5.2) \quad A_{e,f} = \left(\frac{\mu_f}{k_0} \mathbf{v}_e, \mathbf{v}_f \right)_Q = \frac{\mu_f}{2k_0} |E_e| \delta_{e,f},$$

where E_e is the one element or union of two elements that have e as an edge. This approximation diagonalizes A and enables us to eliminate the scaled relative velocity using the Schur complement from the first row of (5.1),

$$\tilde{v}_r = A^{-1} (B_\phi \tilde{q}_f + a_\phi).$$

What remains is a Stokes-like system with *two* pressure potentials. One can further eliminate $v_s = D_\phi^{-1} (Bq + b_\phi)$ to obtain

$$(5.3) \quad \begin{pmatrix} B_\phi^T A^{-1} B_\phi + C_{f,\phi} & -C_{f,s,\phi} \\ -C_{f,s,\phi}^T & B^T D_\phi^{-1} B + C_{s,\phi} \end{pmatrix} \begin{pmatrix} \tilde{q}_f \\ q \end{pmatrix} = \begin{pmatrix} -B_\phi^T A^{-1} a_\phi \\ -B^T D_\phi^{-1} b_\phi \end{pmatrix},$$

but the matrix $B^T D_\phi^{-1} B$ is not easily formed. Nevertheless, one can apply this matrix and therefore solve a Schur complement system for the two pressure potentials. The system can be preconditioned by a diagonal preconditioner, using any good preconditioners for the two diagonal blocks, and solved by conjugate gradients; see, e.g., the block preconditioner defined in [32].

6. Numerical results in one dimension. In this section we simulate a compacting column in one dimension [29]. The column extends over $z \in [-L, L]$ and has no flow through the top and bottom boundaries, i.e.,

$$(6.1) \quad v_s(-L) = v_s(L) = u(-L) = u(L) = 0;$$

moreover, the fluid potential scale is set so that

$$(6.2) \quad q_f(0) = 0.$$

We nondimensionalize using the compaction length scale [27]

$$(6.3) \quad L_c = \left(\frac{k_0 \mu_s}{\mu_f} \right)^{1/2}$$

by defining the dimensionless variables

$$x = L_c \tilde{x}, \quad q_f = |\rho_r| g L_c \tilde{q}_f, \quad q_s = |\rho_r| g L_c \tilde{q}_s, \quad u = \frac{k_0 |\rho_r| g}{\mu_f} \tilde{u}, \quad v_s = \frac{k_0 |\rho_r| g}{\mu_f} \tilde{v}_s.$$

After dropping the check accent marks, (1.1)–(1.5) become

$$(6.4) \quad u + \phi^{2+2\Theta} q'_f = 0,$$

$$(6.5) \quad u' + \phi(q_f - q_s) = 0,$$

$$(6.6) \quad [q_s - \frac{1}{3}(1 - 4\phi)v'_s]' = 1 - \phi,$$

$$(6.7) \quad v'_s - \phi(q_f - q_s) = 0.$$

Where $\phi > 0$, we can reduce the system to a single equation in terms of u as follows. First, (6.5) and (6.7) imply that $v'_s = -u'$ and $q_s = u'/\phi + q_f$. Equation (6.4) gives $q'_f = -\phi^{-2-2\Theta}u$. Finally, (6.6) reduces to

$$(6.8) \quad \phi^{2+2\Theta} \left(\frac{3 + \phi - 4\phi^2}{3\phi} u' \right)' - u = \phi^{2+2\Theta}(1 - \phi).$$

On an open interval where $\phi = 0$, the equations reduce to $u = 0$, $q'_s = 1$, and $v'_s = 0$.

6.1. Closed form solutions. We consider the three porosity functions

$$(6.9) \quad \phi_0(z) = \phi_0,$$

$$(6.10) \quad \phi_J(z) = \begin{cases} \phi_- & \text{if } z \leq 0, \\ \phi_+ & \text{if } z > 0, \end{cases}$$

$$(6.11) \quad \phi_2(z) = \begin{cases} 0 & \text{if } z \leq 0, \\ \phi_+ z^2 & \text{if } z > 0, \end{cases}$$

where $\phi_0 > 0$ and $\phi_- \neq \phi_+$ gives a discontinuous jump in ϕ_J . Note that ϕ_0 and ϕ_2 satisfy the condition (2.18), since in fact $\phi_2^{\Theta-1/2} \nabla \phi_2 = 2\phi_+ z^{2\Theta}$ for $z > 0$ is indeed in $L^\infty(\Omega)$. However, ϕ_J does not satisfy this condition.

6.1.1. Constant porosity. Taking $\phi(z) = \phi_0 > 0$, (6.8) reduces to

$$R^{-2}u'' - u = \phi_0^{2+2\Theta}(1 - \phi_0), \quad R = R(\phi_0) = \left(\frac{3 + \phi_0 - 4\phi_0^2}{3} \phi_0^{1+2\Theta} \right)^{-1/2}.$$

Solving the differential equation with the potential scale condition (6.2) gives the full solution in terms of the constants a and b as

$$(6.12) \quad u = -\phi_0^{2+2\Theta}(1 - \phi_0)[1 + a \cosh(Rz) + b \sinh(Rz)],$$

$$(6.13) \quad q_f = (1 - \phi_0) \left\{ z - \frac{b}{R} + \frac{1}{R} [a \sinh(Rz) + b \cosh(Rz)] \right\},$$

$$(6.14) \quad q_s = (1 - \phi_0) \left\{ z - \frac{b}{R} + \frac{1 - 4\phi_0}{3 + \phi_0 - 4\phi_0^2} \frac{\phi_0}{R} [a \sinh(Rz) + b \cosh(Rz)] \right\}.$$

The boundary conditions (6.1) imply that

$$(6.15) \quad v_s = -u, \quad a = -\frac{1}{\cosh(RL)}, \quad \text{and} \quad b = 0.$$

6.1.2. Discontinuous porosity. For $\phi = \phi_J$ in (6.10), we can solve (6.8) on each subdomain where ϕ is constant. If both ϕ_+ and ϕ_- are positive, the result is (6.12)–(6.14), i.e.,

$$(6.16) \quad u_{\pm} = \phi_{\pm}^{2+2\Theta}(1 - \phi_{\pm})[1 + a_{\pm} \cosh(R_{\pm}z) + b_{\pm} \sinh(R_{\pm}z)].$$

For (6.4)–(6.7) to make sense at the interface $z = 0$, the functions that are differentiated must be continuous. The scale condition (6.2) enforces continuity of q_f . We must impose continuity at $z = 0$ on u (and thereby on v_s) and on

$$(6.17) \quad q_s(0) - \frac{1}{3}(1 - 4\phi)v'_s(0) = \left(\frac{1}{\phi} + \frac{1}{3}(1 - 4\phi)\right)u'(0) = R^{-2}\phi^{-2-2\Theta}u'(0).$$

With the boundary condition (6.1), i.e., $u_{\pm}(\pm L) = 0$, we have four conditions that determine a_{\pm} and b_{\pm} . Letting $\mathcal{F}_{\pm} = \phi_{\pm}^{2+2\Theta}(1 - \phi_{\pm})$, the coefficients are determined by solving the relatively simple linear system

$$(6.18) \quad \begin{bmatrix} \cosh(R_+L) & 0 & \sinh(R_+L) & 0 \\ 0 & \cosh(R_-L) & 0 & -\sinh(R_-L) \\ \mathcal{F}_+ & -\mathcal{F}_- & 0 & 0 \\ 0 & 0 & R_-(1 - \phi_+) & -R_+(1 - \phi_-) \end{bmatrix} \begin{bmatrix} a_+ \\ a_- \\ b_+ \\ b_- \end{bmatrix} = \begin{bmatrix} -1 \\ -1 \\ \mathcal{F}_- - \mathcal{F}_+ \\ 0 \end{bmatrix}.$$

In the case that $\phi_- = 0$ but $\phi_+ > 0$, the solution to (6.8) is (6.12)–(6.14) for $z > 0$, but for $z < 0$, $u = v_s = 0$ (i.e., $a_- = b_- = 0$) and $q_s = z + c_-$. The interface conditions imply that

$$(6.19) \quad a_+ = -1, \quad b_+ = \frac{\cosh(R_+L) - 1}{\sinh(R_+L)}, \quad \text{and} \quad c_- = -b_+(1 - \phi_+)/R_+.$$

Finally, if $\phi_- > 0$ and $\phi_+ = 0$, then

$$(6.20) \quad a_- = -1, \quad b_- = \frac{1 - \cosh(R_-L)}{\sinh(R_-L)}, \quad \text{and} \quad c_+ = -b_-(1 - \phi_-)/R_-,$$

where (6.12)–(6.14) gives the solution for $z < 0$ and for $z > 0$, $u = v_s = 0$ (i.e., $a_+ = b_+ = 0$) and $q_s = z + c_+$.

6.1.3. Quadratic porosity approximation. The final closed form solution is an approximation to the system. Set $\Theta = 0$ and take $\phi(z) = \phi_2(z)$ from (6.11). Working on $z > 0$, the differential equation (6.8) reduces to

$$\phi_+^2 z^4 \left(\frac{3 + \phi_+ z^2 - 4\phi_+^2 z^4}{3\phi_+ z^2} u' \right)' - u = \phi_+^2 z^4 (1 - \phi_+ z^2).$$

Assuming that $\phi = \phi_+ z^2 \ll 1$, we retain only the lowest order terms, i.e.,

$$\phi_+ z^4 (z^{-2} u')' - u = \phi_+^2 z^4,$$

which reduces to the Euler equation $\phi_+ z^2 u'' - 2\phi_+ z u' - u = \phi_+^2 z^4$. The Euler exponents are

$$r_1 = \frac{3 + \sqrt{9 + 4/\phi_+}}{2} > 3 \quad \text{and} \quad r_2 = \frac{3 - \sqrt{9 + 4/\phi_+}}{2} < 0,$$

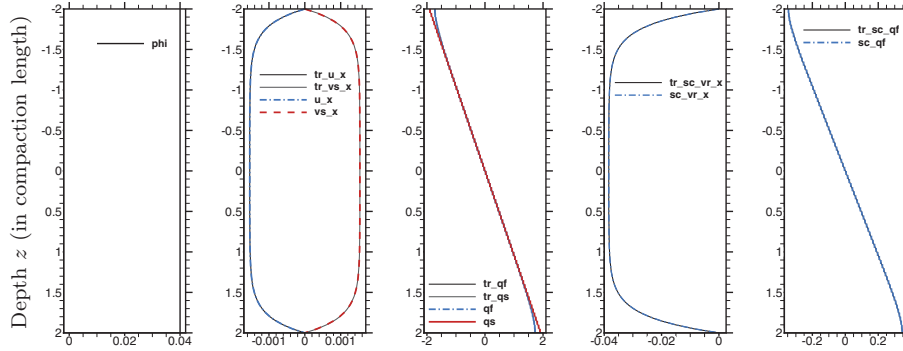


FIG. 1. Constant porosity (6.9) with $\phi_0 = 0.04$. The computed solution is drawn as thick dashed lines, and the closed form solution as thin, solid dark lines. Shown are the porosity ϕ (ϕ), u , $v_s = -u$, q_f , and q_s , as well as the scaled \tilde{v}_r and \tilde{q}_f .

and, if $\phi_+ \neq 1/4$, the solution for $z > 0$ is

$$(6.21) \quad u = -v_s = \frac{\phi_+^2}{1 - 4\phi_+} (L^{4-r_1} z^{r_1} - z^4),$$

$$(6.22) \quad q_f = \frac{1}{1 - 4\phi_+} \left(z - \frac{L^{4-r_1} z^{r_1-3}}{r_1 - 3} \right),$$

$$(6.23) \quad q_s = z.$$

For $z < 0$ where $\phi = 0$, the solution is $u = v_s = 0$ and $q_s = z$.

6.2. Verification of the scaled method. We now present numerical results for the locally conservative scaled mixed method (4.2)–(4.5) and its mass lumped approximation using (5.2). We use the BR spaces for the Stokes part of the system. In one dimension, the velocity part of the RT and BR spaces reduces to piecewise continuous linear functions, and the pressure part is the set of piecewise discontinuous constant functions. The theoretical bound in Theorem 4 would guarantee a convergence rate of $\mathcal{O}(h)$ for the potentials and velocities, provided ϕ satisfies (2.18).

In all tests, we take $L = 2$, so that the domain extends four compaction lengths, and we fix $\Theta = 0$. Each problem is solved on a uniform mesh of n cells. Our computer code is based on the deal.II software library [12].

We chose above to fix the pressure scale of \tilde{q}_f by imposing (6.2) in the interior of the domain. This works well for the closed form solution. However, it does not set the scale properly in our numerical implementation of the problems. Instead, we set the pressure scale of q at the point where it achieves its maximum value. Moreover, in the two cases where the porosity degenerates, we also set the scale of q at the point where it achieves its minimum value.

6.2.1. Constant porosity tests. For the first set of tests, $\phi(z) = \phi_0 = 0.04$. This problem tests the overall performance of the code when there is no degeneracy in the porosity. The computed and closed form solutions using $n = 80$ are shown in Figure 1, although the former is so accurate that it obscures the latter.

In Table 1 we give the relative errors of the potentials as measured in the L^2 -norm for both the scaled mixed method and its mass lumped approximation. The optimal rates of convergence $\mathcal{O}(h)$ are observed for \tilde{q}_f , q_f , and q . We also measured the errors in the discrete L^2 norm, which is the usual L^2 -norm but evaluated using the midpoint

TABLE 1

Constant porosity potential errors. Relative L^2 errors and convergence rates for the potentials. We show results for the scaled mixed method, the mass lumped approximation, and the scaled mixed method but using the discrete L^2 -norm given by using the midpoint rule.

n	\tilde{q}_f		q_f		q	
	L^2 error	rate	L^2 error	rate	L^2 error	rate
Scaled mixed method						
20	1.428e-02	1.00	3.237e-02	1.00	3.435e-02	1.00
40	7.139e-03	1.00	1.618e-02	1.00	1.717e-02	1.00
80	3.569e-03	1.00	8.090e-03	1.00	8.581e-03	1.00
160	1.784e-03	1.00	4.044e-03	1.00	4.290e-03	1.00
Mass lumped method						
20	1.427e-02	1.00	3.236e-02	1.00	3.434e-02	1.00
40	7.139e-03	1.00	1.618e-02	1.00	1.717e-02	1.00
80	3.569e-03	1.00	8.090e-03	1.00	8.581e-03	1.00
160	1.784e-03	1.00	4.044e-03	1.00	4.290e-03	1.00
Scaled mixed method, discrete norm (midpoint rule)						
20	8.597e-04	1.46	1.949e-03	1.46	1.411e-03	0.91
40	2.794e-04	1.62	6.334e-04	1.62	5.422e-04	1.38
80	8.263e-05	1.76	1.873e-04	1.76	1.708e-04	1.67
160	2.271e-05	1.86	5.149e-05	1.86	4.813e-05	1.83
320	5.972e-06	1.93	1.354e-05	1.93	1.279e-05	1.91

TABLE 2

Constant porosity velocity errors. Relative L^2 errors and convergence rates for the velocities using the scaled mixed method and the mass lumped approximation.

n	\tilde{v}_f		u		v_s	
	L^2 error	rate	L^2 error	rate	L^2 error	rate
Scaled mixed method						
20	1.147e-03	1.82	4.897e-05	1.82	4.897e-05	1.82
40	2.972e-04	1.95	1.269e-05	1.95	1.269e-05	1.95
80	7.500e-05	1.99	3.203e-06	1.99	3.203e-06	1.99
160	1.879e-05	2.00	8.027e-07	2.00	8.027e-07	2.00
Mass lumped method						
20	1.650e-03	1.72	7.047e-05	1.72	7.047e-05	1.72
40	4.381e-04	1.91	1.871e-05	1.91	1.871e-05	1.91
80	1.113e-04	1.98	4.753e-06	1.98	4.753e-06	1.98
160	2.794e-05	1.99	1.193e-06	1.99	1.193e-06	1.99

quadrature rule. This is a norm for which one might expect to see superconvergence. Indeed, we see superconvergence for all three potentials. On coarser meshes we see $\mathcal{O}(h^{3/2})$ for the fluid potentials and $\mathcal{O}(h)$ for the mixture, but on fine meshes the rates rise to $\mathcal{O}(h^2)$ for all three variables. Similar superconvergence results hold for the mass lumped approximation.

In Table 2 we give the relative errors of the velocities in the L^2 -norm for both the scaled mixed method and its mass lumped approximation. The velocities are approximated by piecewise linears, so the optimal rates of convergence would be $\mathcal{O}(h^2)$. This is precisely what is observed for each of the velocities \tilde{v}_f , u , and v_s .

6.2.2. Discontinuous porosity tests. For the next set of tests, we use the discontinuous porosity function (6.10) with $\phi_- = 0$ and $\phi_+ = 0.04$. Not only is there a jump in porosity, but it is also degenerate for $z < 0$. The discontinuity will land on a mesh point if n is even, and it will land in the center of a cell if n is odd.

The computed and closed form solutions using $n = 80$ are shown in Figure 2. The discontinuity in q_s is clearly evident. Note that the computation of \tilde{v}_r has some difficulty near $z = 0$ where the porosity is discontinuous. This difficulty is not seen

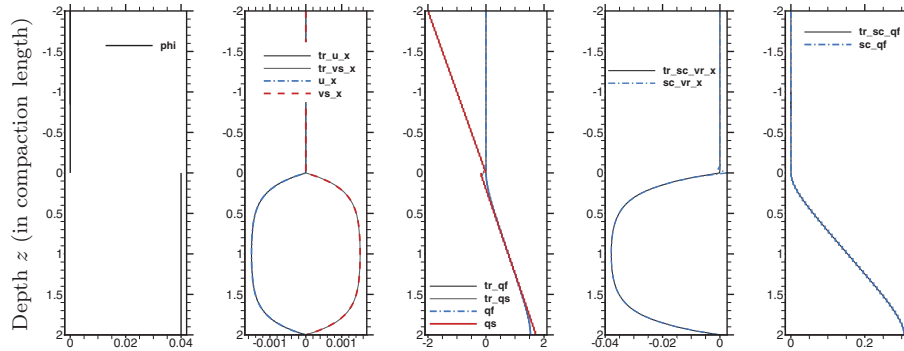


FIG. 2. Discontinuous porosity (6.10) with $\phi_- = 0$ and $\phi_+ = 0.04$. The computed solution is drawn as thick dashed lines, and the closed form solution as thin, solid dark lines. Shown are the porosity ϕ (ϕ), u , $v_s = -u$, q_f , and q_s , as well as the scaled \tilde{v}_r and \tilde{q}_f .

TABLE 3

Discontinuous porosity potential errors. Relative L^2 errors and convergence rates for the potentials. We show results for the scaled mixed method, including one case using the discrete L^2 -norm given by using the midpoint rule, and one case with the L^2 -norm restricted to the interior (i.e., away from the discontinuity). When n is even, the discontinuity is at a computational mesh point, but not when n is odd.

n	\tilde{q}_f		q_f		q	
	L^2 error	rate	L^2 error	rate	L^2 error	rate
Scaled mixed method, n even						
20	1.040e-02	1.00	2.852e-02	1.00	3.622e-02	1.00
40	5.202e-03	1.00	1.426e-02	1.00	1.811e-02	1.00
80	2.601e-03	1.00	7.133e-03	1.00	9.055e-03	1.00
160	1.301e-03	1.00	3.567e-03	1.00	4.527e-03	1.00
Scaled mixed method, n odd						
21	9.961e-03	0.98	2.744e-02	0.98	3.955e-02	0.87
41	5.140e-03	0.99	1.416e-02	0.99	2.321e-02	0.80
81	2.611e-03	0.99	7.184e-03	1.00	1.428e-02	0.71
161	1.316e-03	1.00	3.615e-03	1.00	9.218e-03	0.64
Scaled mixed method, n even, discrete norm (midpoint rule)						
20	9.536e-04	1.64	2.615e-03	1.64	1.231e-04	1.46
40	2.529e-04	1.91	6.935e-04	1.91	3.860e-05	1.67
80	6.225e-05	2.02	1.707e-04	2.02	1.102e-05	1.81
160	1.505e-05	2.05	4.128e-05	2.05	2.964e-06	1.89
Scaled mixed method, n odd, interior						
21	9.126e-03	0.68	2.502e-02	0.68	3.016e-02	0.74
41	4.983e-03	0.90	1.367e-02	0.90	1.658e-02	0.89
81	2.572e-03	0.97	7.052e-03	0.97	8.670e-03	0.95
161	1.304e-03	0.99	3.576e-03	0.99	4.431e-03	0.98

in u and v_s , since compared to \tilde{v}_r , these velocities are multiplied by ϕ . Overall, the computed solution is an excellent match to the closed form one.

In Table 3 we give convergence results for the potentials using the scaled mixed method. The mass lumped approximation has nearly identical results. Even though ϕ does not satisfy the condition (2.18) and is in fact discontinuous, we see good convergence results. When n is even and the grid resolves the discontinuity in ϕ , we see optimal convergence rates $\mathcal{O}(h)$ for all three potentials and superconvergence $\mathcal{O}(h^2)$ when using the discrete norm.

When n is odd and the discontinuity in ϕ is not resolved, we see some degradation in the convergence rate for q . Not shown is that superconvergence is not seen in the

TABLE 4

Discontinuous porosity velocity errors. Relative L^2 errors and convergence rates for the velocities. We show results for the scaled mixed method and the mass lumped approximation using (5.2). When n is even, the discontinuity is at a computational mesh point, but not when n is odd.

n	\tilde{v}_f		u		v_s	
	L^2 error	rate	L^2 error	rate	L^2 error	rate
Scaled mixed method, n even						
20	2.929e-03	1.33	4.714e-05	1.85	4.714e-05	1.85
40	1.116e-03	1.39	1.213e-05	1.96	1.213e-05	1.96
80	4.120e-04	1.44	3.090e-06	1.97	3.090e-06	1.97
160	1.491e-04	1.47	7.850e-07	1.98	7.850e-07	1.98
Mass lumped method, n even						
20	1.695e-03	1.73	7.076e-05	1.73	7.076e-05	1.73
40	4.499e-04	1.91	1.878e-05	1.91	1.878e-05	1.91
80	1.143e-04	1.98	4.770e-06	1.98	4.770e-06	1.98
160	2.869e-05	1.99	1.197e-06	1.99	1.197e-06	1.99
Scaled mixed method, n odd						
21	2.027e-03	1.20	9.004e-05	1.25	9.004e-05	1.25
41	1.055e-03	0.98	4.524e-05	1.03	4.524e-05	1.03
81	5.614e-04	0.93	2.368e-05	0.95	2.368e-05	0.95
161	2.925e-04	0.95	1.227e-05	0.96	1.227e-05	0.96

discrete norm. To test whether the error near the discontinuity pollutes the solution, we computed the *interior* errors, given by computing the error in all cells of the mesh except the five near the discontinuity. That is, we restrict the domain of integration of the L^2 -norm to be interior to where ϕ is smooth by removing the center cell and its two neighbors on each side. This mesh dependent norm shows good $\mathcal{O}(h)$ convergence, and so indeed the error is localized to the region of the discontinuity. We do not, however, observe superconvergence in the discrete interior norm when n is odd.

The errors in the velocities are given in Table 4. We see good rates of convergence when n is even, being $\mathcal{O}(h^2)$ for all cases except the scaled method's \tilde{v}_f , which is still $\mathcal{O}(h^{3/2})$. When n is odd, we observe $\mathcal{O}(h)$ convergence (we show only the scaled method, but the mass lumped approximation is similar).

6.2.3. Quadratic porosity tests. For the final set of tests, we use the quadratic porosity (6.11) with $\phi_+ = 0.001$, i.e., $\phi_2(z) = 0.001 z^2$ for $z > 0$ and $\phi_2(z) = 0$ for $z \leq 0$. The maximal value of ϕ is $\phi(2) = 0.004$, so the analytic solution (6.21)–(6.23) should approximate the true solution reasonably well, at least if n is not too large.

The computed and closed form solutions agree quite well, as shown in Figure 3 using $n = 80$, even though there is a boundary layer near $z = 2$ in the velocities that is difficult to resolve. Convergence results for the potentials and velocities are given in Tables 5 and 6 using the scaled method (the mass lumped approximation gives nearly identical results). We expect convergence only if the approximate true solution is adequate. Indeed, we see some degradation of the results as n becomes large, because the numerical solution does not converge to the closed form solution (6.21)–(6.23). When the approximation is adequate, we see that the potentials converge to $\mathcal{O}(h)$. (The discrete norm does not display superconvergence for this test problem, but the errors are much smaller.) The velocities converge to at least $\mathcal{O}(h)$, and may approach $\mathcal{O}(h^2)$ before the grid becomes too fine.

6.3. Condition number as positive ϕ tends to zero. We now turn our attention to the nondegenerate problem, so that we can solve the system of equations using other mixed formulations. We compare our method to a relatively standard

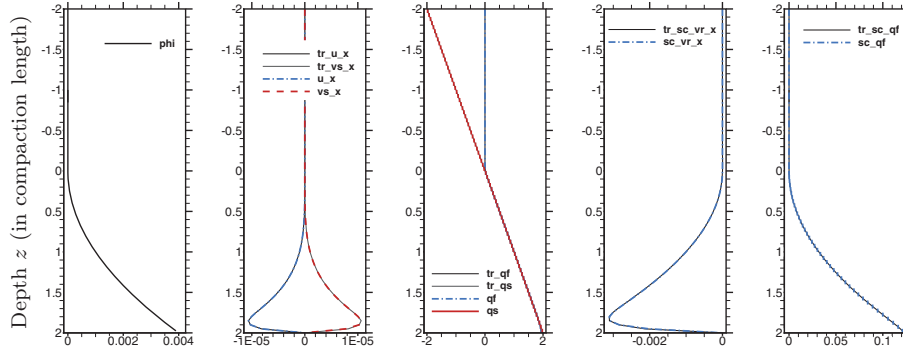


FIG. 3. Quadratic porosity (6.11) with $\phi_+ = 0.001$. The computed solution is drawn as thick dashed lines, and the closed form solution as thin, solid dark lines. Shown are the porosity ϕ (ϕ), u , $v_s = -u$, q_f , and q_s , as well as the scaled \tilde{v}_r and \tilde{q}_f .

TABLE 5

Quadratic porosity potential errors. Relative L^2 errors and convergence rates for the potentials for the scaled mixed method. When n is even, the transition to positive porosity is at a computational mesh point, but not when n is odd.

n	\tilde{q}_f		q_f		q	
	L^2 error	rate	L^2 error	rate	L^2 error	rate
Scaled mixed method, n even						
20	5.326e-03	1.01	3.037e-02	1.01	3.490e-02	1.00
40	2.656e-03	1.00	1.520e-02	1.00	1.747e-02	1.00
80	1.329e-03	1.00	7.667e-03	0.99	8.768e-03	0.99
160	6.675e-04	0.99	3.963e-03	0.95	4.457e-03	0.98
Scaled mixed method, n odd						
21	5.070e-03	1.01	2.893e-02	1.01	3.324e-02	1.00
41	2.592e-03	1.00	1.484e-02	1.00	1.704e-02	1.00
81	1.313e-03	1.00	7.575e-03	0.99	8.660e-03	0.99
161	6.634e-04	0.99	3.940e-03	0.95	4.430e-03	0.98

TABLE 6

Quadratic porosity velocity errors. Relative L^2 errors and convergence rates for the velocities for the scaled mixed method. When n is even, the transition to positive porosity is at a computational mesh point, but not when n is odd.

n	\tilde{v}_f		u		v_s	
	L^2 error	rate	L^2 error	rate	L^2 error	rate
Scaled mixed method, n even						
20	3.663e-04	1.23	1.546e-06	1.14	1.546e-06	1.14
40	1.166e-04	1.65	5.104e-07	1.60	5.104e-07	1.60
80	3.252e-05	1.84	1.429e-07	1.84	1.429e-07	1.84
160	1.079e-05	1.59	4.342e-08	1.72	4.342e-08	1.72
Scaled mixed method, n odd						
21	3.408e-04	1.27	1.444e-06	1.19	1.444e-06	1.19
41	1.115e-04	1.67	4.886e-07	1.62	4.886e-07	1.62
81	3.179e-05	1.84	1.397e-07	1.84	1.397e-07	1.84
161	1.071e-05	1.58	4.304e-08	1.71	4.304e-08	1.71

MFEM and to the expanded MFEM [8]. We also compare our method to a symmetry preserving formulation in which we balance the degeneracy by using a square-root scaling of the coefficient $\phi^{2+2\Theta}$ in (2.4) and modify (2.5), as was done in, e.g., [19]. However, we would still need to divide by ϕ in a standard approach, so we modify

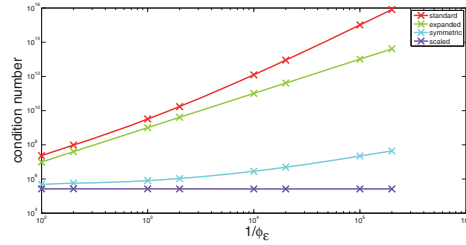


FIG. 4. *Compacting column. Condition numbers for the standard, expanded, symmetric, and scaled formulations as $\phi_\epsilon \rightarrow 0^+$ for the porosities defined in (6.9)–(6.11) plus ϕ_ϵ .*

the test function as was done in our scaled method and in [6]. (See [4] for explicit definitions of the three methods.) We consider the two nonconstant porosities (6.10)–(6.11) ($\phi_- = 0$), but add a small positive constant $\phi_\epsilon > 0$ to each. We take $\phi_\epsilon \rightarrow 0^+$ and observe the condition number of the linear system that is solved by each method.

In Figure 4 we see that the condition number increases rapidly as $\phi_\epsilon \rightarrow 0$ except for the scaled method, which remains stable. Indeed, as we saw already, the scaled method works well even when the porosity is identically zero in parts of the domain.

7. Numerical results in two dimensions. Finally, we present numerical results for the mass lumped approximation of the scaled method (4.2)–(4.5), (5.2) using the TH Stokes spaces and the deal.II software library [12]. Our two-dimensional problem is related to simulation of the mantle near a midocean ridge (MOR).

If porosity is constant, $\phi = \phi_0$, and one sets $\nabla \cdot \mathbf{u} = \nabla \cdot \mathbf{v}_s = 0$, then (1.1)–(1.5) can be solved in an infinite quarter-plane $\{x > 0, z > 0\}$ [36]. This problem describes viscous corner flow if, at the top of the mantle $\{z = 0\}$, one sets the MOR spreading rate as a boundary condition $\mathbf{v}_s \cdot \boldsymbol{\tau} = \mathbf{v}_s \cdot \hat{\mathbf{x}} = U_0$, and on the ridge axis $\{x = 0\}$, one sets the symmetry condition $\mathbf{v}_f \cdot \boldsymbol{\nu} = \mathbf{v}_f \cdot \hat{\mathbf{z}} = 0$ and $\partial \mathbf{v}_s / \partial x = 0$. The solution is

$$(7.1) \quad q = q_s = q_f = (1 - \phi_0) \left(\frac{4\mu_s U_0}{\pi(x^2 + z^2)} + |\rho_r|g \right) z,$$

$$(7.2) \quad \mathbf{v}_s = \frac{2U_0}{\pi(x^2 + z^2)} \begin{pmatrix} \tan^{-1}(x/z)(x^2 + z^2) - xz \\ -z^2 \end{pmatrix},$$

$$(7.3) \quad \mathbf{u} = \frac{k_0(1 - \phi_0)\phi_0^{2+2\Theta}}{\mu_f} \left\{ \frac{4\mu_s U_0}{\pi(x^2 + z^2)^2} \begin{pmatrix} 2xz \\ z^2 - x^2 \end{pmatrix} + \rho_r \mathbf{g} \right\}.$$

We solve the full system of equations on the rectangular domain $-160 \text{ km} < x < 160 \text{ km}$, $0 < z < 160 \text{ km}$. The MOR is at $(0, 0)$. We take $\mu_s = 10^{19} \text{ Pa}\cdot\text{s}$, $\mu_f = 1 \text{ Pa}\cdot\text{s}$, $\rho_s = 3300 \text{ kg/m}^3$, $\rho_f = 2800 \text{ kg/m}^3$, $k_0 = 10^{-8} \text{ m}^2$, $\Theta = 0$, and $U_0 = 10^{-9} \text{ m/s} = 3.1536 \text{ cm/yr}$. The boundary conditions are defined by imposing (7.2) on the Stokes velocity \mathbf{v}_s and (7.1) on the potential $q = q_f$. However, to avoid the singularity at the corner, we translate x to $x - \ell$ when $x < 0$ and $x + \ell$ when $x > 0$ before evaluating (7.1)–(7.2), where we arbitrarily set $\ell = 20 \text{ m}$. We use a mesh of 160×80 elements.

In Figure 5 we show the Stokes solution using $\phi = 0$. Note that the mantle flows up to the MOR and outward from there. There is no fluid melt in this computation, although our code solves for the Darcy system as well as the Stokes system. Rather than normalizing the average of q to zero, we set a single point to zero.

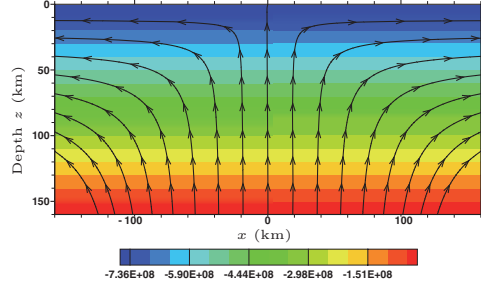


FIG. 5. MOR-like example with zero porosity. We show the solid matrix potential q_s as a contour in $\text{Kg}/(\text{m}\cdot\text{s}^2)$ and the velocity v_s as streamlines.

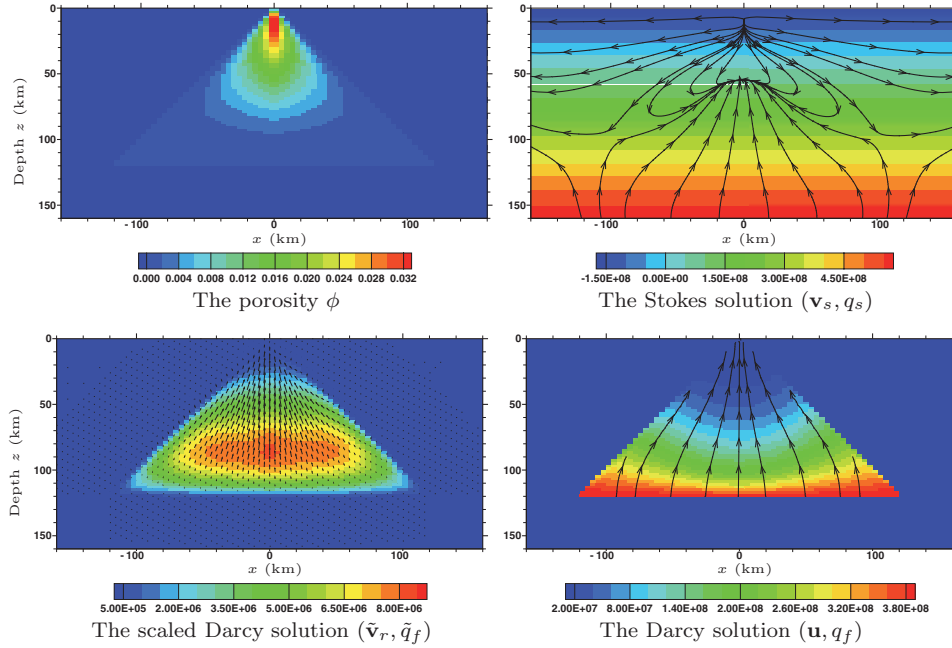


FIG. 6. MOR-like example. We show the porosity and the solutions (\mathbf{v}_s, q_s) , $(\tilde{\mathbf{v}}_r, \tilde{q}_f)$, and (\mathbf{u}, q_f) . The porosity and potentials in $\text{Kg}/(\text{m}\cdot\text{s}^2)$ are shown as contours, the relative velocity $\tilde{\mathbf{v}}_r$ as arrows, and the other velocities as streamlines.

We then set a not identically vanishing porosity by the formula

$$(7.4) \quad \phi(x, z) = \begin{cases} 0.05 \left(\frac{120 \text{ km} - z}{120 \text{ km}} \right)^2 \left(1 - \frac{|x|}{z + \ell} \right) & \text{if } z \leq 120 \text{ km and } |x| \leq z + \ell, \\ 0 & \text{otherwise,} \end{cases}$$

wherein we could offset by any value of ℓ , but we simply chose to use the previously chosen value $\ell = 20 \text{ m}$. In Figure 6, we show the porosity and the solutions (\mathbf{v}_s, q_s) to the Stokes system, $(\tilde{\mathbf{v}}_r, \tilde{q}_f)$ to the scaled Darcy system, and (\mathbf{u}, q_f) to the Darcy system. The form of the solution is dictated by our (arbitrary) choice of ϕ . The scaled and unscaled Darcy solutions show fluid melt rising and focusing into the MOR, and some melt leaving the domain to form new crust. The Stokes solution varies significantly from the case in Figure 5 where there is no melt. The solid matrix rises at the bottom, but it falls at the top near the MOR to compensate for the rise of

fluid melt to the surface. The scaled potential \tilde{q}_f is much smoother than q_f , and so the former is much easier to approximate. The porosity vanishes in a significant portion of the domain; nevertheless, there is no difficulty solving the system accurately.

8. Conclusions. We developed a mathematically well-posed, mixed variational framework for McKenzie's equations governing the mechanics of a mixture of molten and solid materials [28], assuming that the porosity ϕ is given and satisfies the hypothesis (2.18). Our formulation handles the regions where there are two phases (i.e., the mixture variable $\phi > 0$) as well as the mathematically degenerate regions where there is only the single solid matrix phase (i.e., $\phi = 0$). The formulation is based on a careful scaling of the Darcy variables by powers of the porosity [6].

We defined an MFEM based on our scaled variational formulation and proved its optimal order convergence. We also presented two modifications, one that is locally mass conservative, and the other involving mass lumping (5.2) to simplify and increase solver efficiency on rectangular meshes.

Numerical results of a one-dimensional compacting column with various porosity functions (6.9)–(6.11) showed an excellent match to the closed form solutions for ϕ_0 and ϕ_J , as well as a good match to the approximate solution for ϕ_2 . Degeneracies in the porosity posed no difficulties for the simulations; in fact, the condition number of the linear system is nearly insensitive to degeneracies in ϕ . The results showed that the method indeed achieves optimal convergence and that the mass lumping approximation does not degrade the results in any way.

The nondegenerate constant porosity example showed $\mathcal{O}(h)$ convergence for the potentials and superconvergence of order $\mathcal{O}(h^2)$ when measured in the discrete midpoint rule norm. The velocity achieved the optimal $\mathcal{O}(h^2)$ convergence for this one-dimensional problem. The degenerate, quadratic porosity example also showed optimal $\mathcal{O}(h)$ convergence of the potentials and perhaps $\mathcal{O}(h^2)$ convergence of the velocities, regardless of whether the computational mesh resolved the point where ϕ transitions from zero to positive.

The degenerate and discontinuous porosity example had an interesting set of results. Even though the porosity does not satisfy (2.18), the MFEM achieved good, but not necessarily optimal, convergence in all cases. When the computational mesh resolved the transition point of ϕ , we saw $\mathcal{O}(h)$ convergence for the potentials and superconvergence of order $\mathcal{O}(h^2)$ when measured in the discrete midpoint rule norm. We also saw $\mathcal{O}(h^{3/2})$ convergence for \tilde{v}_r and $\mathcal{O}(h^2)$ convergence for u and v_s . The mass lumped approximation actually improved the convergence to $\mathcal{O}(h^2)$ for all three velocities. However, when the computational mesh did not resolve the transition point in ϕ , we saw $\mathcal{O}(h)$ convergence for the potentials \tilde{q}_f and q_f , but only $\mathcal{O}(h^{1/2})$ for q . The discrete norm did not help, but we did verify that the main errors were localized to a region near the transition point, since removing the error there led to $\mathcal{O}(h)$ convergence for all three potentials. The velocities converged to order $\mathcal{O}(h)$. This example suggests that the condition (2.18) may not be strictly necessary.

In the full model of mantle dynamics, the porosity evolves and so must be approximated. In a finite element or discontinuous Galerkin method, one would naturally approximate ϕ by continuous or discontinuous polynomials on each element of the computational mesh. Any jumps in the porosity will then naturally lie on the boundaries of the elements, and so we would expect our method to perform well.

A two-dimensional test example akin to an MOR showed the strong effect that melt can have on the velocity field. Even though the porosity vanished in much of the domain, our locally conservative scaled finite element method showed good results.

Using the mass lumped approximation, the method easily reduces to a single Stokes system with two potentials, and the efficiency of the linear solver is fairly insensitive to the absence of melt. We believe that our method is highly suited to realistic problems of the mechanics of mantle dynamics, and that it can be used effectively as a component of the full mantle dynamics problem.

REFERENCES

- [1] R. A. ADAMS, *Sobolev Spaces*, Academic Press, New York, 1975.
- [2] E. AHARONOV, J. A. WHITEHEAD, P. B. KELEMEN, AND M. SPIEGELMAN, *Channeling instability of upwelling melt in the mantle*, J. Geophys. Res., 100 (1995), pp. 20,433–20,450.
- [3] T. ARBOGAST AND M. R. CORREA, *Two families of $H(\text{div})$ mixed finite elements on quadrilaterals of minimal dimension*, SIAM J. Numer. Anal., 54 (2016), pp. 3332–3356, <https://doi.org/10.1137/15M1013705>.
- [4] T. ARBOGAST, M. A. HESSE, AND A. L. TAICHER, *Mixed Methods for Two-Phase Darcy-Stokes Mixtures of Partially Melted Materials with Regions of Zero Porosity*, ICES Technical report 16–18, Institute for Computational Engineering and Sciences, University of Texas at Austin, Austin, TX 2016.
- [5] T. ARBOGAST AND A. L. TAICHER, *A cell-centered finite difference method for a degenerate elliptic equation arising from two-phase mixtures*, Comput. Geosci., submitted.
- [6] T. ARBOGAST AND A. L. TAICHER, *A linear degenerate elliptic equation arising from two-phase mixtures*, SIAM J. Numer. Anal., 54 (2016), pp. 3105–3122, <https://doi.org/10.1137/16M1067846>.
- [7] T. ARBOGAST AND M. F. WHEELER, *A family of rectangular mixed elements with a continuous flux for second order elliptic problems*, SIAM J. Numer. Anal., 42 (2005), pp. 1914–1931.
- [8] T. ARBOGAST, M. F. WHEELER, AND I. YOTOV, *Mixed finite elements for elliptic problems with tensor coefficients as cell-centered finite differences*, SIAM J. Numer. Anal., 34 (1997), pp. 828–852.
- [9] A. ASCHWANDEN, E. BUELER, C. KHROULEV, AND H. BLATTER, *An enthalpy formulation for glaciers and ice sheets*, J. Glaciol., 58 (2012), pp. 441–457.
- [10] I. BABUŠKA, *Error-bounds for finite element method*, Numer. Math., 16 (1971), pp. 322–333.
- [11] I. BABUŠKA, *The finite element method with Lagrangian multipliers*, Numer. Math., 20 (1973), pp. 179–192.
- [12] W. BANGERTH, T. HEISTER, L. HELTAI, G. KANSCHAT, M. KRONBICHLER, M. MAIER, AND B. TURCKIN, *The deal.II library, version 8.3*, Arch. Numer. Software, 4 (2016).
- [13] J. BEAR AND A. H.-D. CHENG, *Modeling Groundwater Flow and Contaminant Transport*, Springer, New York, 2010.
- [14] C. BERNARDI AND G. RAUGEL, *Analysis of some finite elements for the Stokes problem*, Math. Comp., 44 (1985), pp. 71–79.
- [15] J. H. BRAMBLE, *A proof of the inf-sup condition for the Stokes equations on Lipschitz domains*, Math. Models Methods Appl. Sci., 13 (2003), pp. 361–371.
- [16] S. C. BRENNER AND L. R. SCOTT, *The Mathematical Theory of Finite Element Methods*, Springer, New York, 1994.
- [17] F. BREZZI AND M. FORTIN, *Mixed and Hybrid Finite Element Methods*, Springer, New York, 1991.
- [18] Z. CHEN, G. HUAN, AND Y. MA, *Computational Methods for Multiphase Flows in Porous Media*, Comput. Sci. Eng. 2, SIAM, Philadelphia, 2006.
- [19] B. COCKBURN AND C.-W. SHU, *Runge-Kutta discontinuous Galerkin methods for convection-dominated problems*, J. Sci. Comput., 16 (2001), pp. 173–261.
- [20] J. DOUGLAS, JR., T. DUPONT, AND L. WAJLBIN, *The stability in L^q of the L^2 -projection into finite element function spaces*, Numer. Math., 23 (1975), pp. 193–197.
- [21] A. ERN AND J.-L. GUERMOND, *Theory and Practice of Finite Elements*, Appl. Math. Sci. 159, Springer, New York, 2004.
- [22] A. C. FOWLER, *On the transport of moisture in polythermal glaciers*, Geophys. Astrophys. Fluid Dyn., 28 (1984), pp. 99–140.
- [23] V. GIRAULT AND P. A. RAVIART, *Finite Element Methods for Navier-Stokes Equations: Theory and Algorithms*, Springer, Berlin, 1986.
- [24] P. GRISVARD, *Elliptic Problems in Nonsmooth Domains*, Pitman, Boston, 1985.
- [25] M. A. HESSE, A. R. SCHIEMENZ, Y. LIANG, AND E. M. PARMENTIER, *Compaction-dissolution waves in an upwelling mantle column*, Geophys. J. Int., 187 (2011), pp. 1057–1075.

- [26] I. J. HEWITT AND A. C. FOWLER, *Partial melting in an upwelling mantle column*, R. Soc. Lond. Proc. Ser. A Math. Phys. Sci. Eng., 464 (2008), pp. 2467–2491.
- [27] R. F. KATZ, *Magma dynamics with the enthalpy method: Benchmark solutions and magmatic focusing at mid-ocean ridges*, J. Petrol., 49 (2008), pp. 2099–2121.
- [28] D. MCKENZIE, *The generation and compaction of partially molten rock*, J. Petrol., 25 (1984), pp. 713–765.
- [29] D. MCKENZIE, *Th-U disequilibrium and the melting processes beneath ridge axes*, Earth Planetary Sci. Lett., 72 (1985), pp. 149–157.
- [30] J. T. ODEN AND L. F. DEMKOWICZ, *Applied Functional Analysis*, CRC Press, Boca Raton, FL, 1996.
- [31] R. A. RAVIART AND J. M. THOMAS, *A mixed finite element method for 2nd order elliptic problems*, in Mathematical Aspects of Finite Element Methods, I. Galligani and E. Magenes, eds., Lecture Notes in Math. 606, Springer, New York, 1977, pp. 292–315.
- [32] S. RHEBERGEN, G. N. WELLS, R. F. KATZ, AND A. J. WATHEN, *Analysis of block preconditioners for models of coupled magma/mantle dynamics*, SIAM J. Sci. Comput., 36 (2014), pp. A1960–A1977.
- [33] J. E. ROBERTS AND J.-M. THOMAS, *Mixed and hybrid methods*, in Handbook of Numerical Analysis, Vol. 2, P. G. Ciarlet and J. L. Lions, eds., North-Holland, Amsterdam, 1991, pp. 523–639.
- [34] N. SLEEP, *Tapping of melt by veins and dikes*, J. Geophys. Res., 93 (1988), pp. 255–272.
- [35] E. A. SPIEGEL AND G. VERONIS, *On the Boussinesq approximation for a compressible fluid*, Astrophys. J., 131 (1960), pp. 442–447.
- [36] M. SPIEGELMAN AND D. MCKENZIE, *Simple 2-D models for melt extraction at mid-ocean ridges and island arcs*, Earth Planetary Sci. Lett., 83 (1987), pp. 137–152.
- [37] C. J. VAN DER VEEN, *Fundamentals of Glacier Dynamics*, 2nd ed., CRC Press, Boca Raton, FL, 2013.
- [38] W. ZHU, G. GAETANI, F. FUSSEIS, L. MONTESI, AND F. D. CARLO, *Microtomography of partially molten rocks: Three dimensional melt distribution in mantle peridotite*, Science, 332 (2011), pp. 88–91.