# A NESTED NEWTON-TYPE ALGORITHM FOR FINITE VOLUME METHODS SOLVING RICHARDS' EQUATION IN MIXED FORM[*]

VINCENZO CASULLI[†] AND PAOLA ZANOLLI[‡]

**Abstract.** A finite volume discretization of the mixed form of Richards' equation leads to a nonlinear numerical model which yields exact local and global mass conservation. The resulting nonlinear system requires sophisticated numerical strategies, especially in a variable saturated flow regime. In this paper a nested, Newton-type algorithm for the discretized Richards' equation is proposed and analyzed. With a judicious choice of the initial guess, the quadratic convergence rate is obtained for any time step size and for all flow regimes.

**Key words.** Richards' equation, variably saturated flow, finite volume, mildly nonlinear systems, Jordan decomposition, nested iterations

**AMS subject classifications.** 65M08, 65M12, 76M12, 76S05

**DOI.** 10.1137/100786320

**1. Introduction.** Flow through variably saturated porous media is governed by Richards' equation which is combined with the soil properties relating the moisture content and the hydraulic conductivity to the pressure head. Such an equation, in conservative form, can be written as

$$(1) \qquad \frac{\partial \theta(\psi)}{\partial t} = \nabla \cdot [K(\psi)\nabla(\psi + z)] + S,$$

where $\psi$ is the pressure head, $\theta(\psi)$ is the moisture content, $K(\psi)$ is the nonnegative hydraulic conductivity, $S$ is the source term, $t$ is the time, and $z$ is the vertical coordinate assumed positive upward. To complete the model, the moisture content and the hydraulic conductivity need to be specified by prescribing their constitutive relationships (see, e.g., [3, 24]).

Equation (1) is known as the *mixed form* of Richards' equation. This equation can also be expressed either in terms of $\theta$ or in terms of $\psi$. One of the advantages of the $\theta$-based formulation is that it can be solved perfectly by mass conservative methods. This form, however, degenerates under fully saturated conditions as heterogeneous material produces discontinuous $\theta$ profiles and a pressure-saturation relationship no longer exists. The $\psi$-based form allows for both unsaturated and saturated conditions. In highly nonlinear problems, however, numerical methods based on the $\psi$-based form can suffer from large mass balance errors (see, e.g., [7, 9, 11]).

In the mixed $\psi - \theta$-based form of Richards' equation (1), both variables, the moisture content and pressure head, are employed. Numerical techniques that employ both $\theta$ and $\psi$, with possibly sophisticated variable switching techniques, have been developed to minimize mass balance errors (see, e.g., [7, 9, 12, 15, 18, 19, 21, 22]).

Standard approximations typically applied to the spatial domain include finite differences [7, 9, 22], finite elements [2, 7, 12, 15], and finite volumes [12, 19, 21].

[†]Laboratory of Applied Mathematics, Department of Civil and Environmental Engineering, University of Trento, 38050 Trento, Italy (vincenzo.casulli@unitn.it).

[‡]Department of Mathematics, University of Trento, 38050 Trento, Italy (paola.zanolli@unitn.it).

Then, regardless of the approximation used, the nonlinear character of the constitutive relationships poses a major difficulty for solving the resulting large system of algebraic equations. This normally requires the use of iterative schemes, such as the Picard and Newton methods.

The moisture content $\theta(\psi)$ is, in general, a nonlinear function of the pressure head. Even when this function is sufficiently regular over the entire range of pressure heads, its behavior, especially near saturation, is such that the derivatives of $\theta(\psi)$ can exhibit sharp changes [3, 24]. It is this nonlinear dependency of the moisture content on the pressure head that makes the numerical solution of Richards' equation problematic and requires sophisticated numerical methods in order to overcome convergence problems and/or poor computational efficiency [2, 10, 12, 15, 22, 23].

In linearization schemes such as the Picard and the Newton, the number of iterations needed to converge is a determining factor for the simulation efficiency. To this purpose, convergence rate is often enhanced by providing the solver with an initial estimate that is closer to the final solution for the current time step. This can be obtained by taking the initial guess from the previous time step and by choosing a sufficiently small time step size [2]. Thus, numerical algorithms often include an empirical time step adaptation criterion (see, e.g., [10, 12, 15, 23]).

In this paper the moisture capacity $c(\psi)$ within each control volume, which is a nonnegative function with bounded variations, is first expressed as a difference of two nonnegative, nondecreasing, and bounded functions $c(\psi) = p(\psi) - q(\psi)$. Accordingly, the moisture content is expressed as a difference of volumes $\theta(\psi) = \theta_1(\psi) - \theta_2(\psi)$, where $\theta_1(\psi)$ and $\theta_2(\psi)$ are integrals of $p(\psi)$ and $q(\psi)$, respectively. Then a nested Newton-type algorithm for a finite volume discretization of Richards' equation is derived by linearizing, in order, $\theta_1(\psi)$ and $\theta_2(\psi)$ in the inner and in the outer cycle, respectively. Convergence of the iterations and exact mass conservation are ensured for any time step size, for all flow regimes, and for a wide class of constitutive relationships.

The remainder of this paper is organized as follows. In section 2 the Picard linearization for a finite volume discretization of the mixed form of Richards' equation is introduced. A nested Newton-type iterative method for solving the resulting mildly nonlinear system is given in section 3. The convergence analysis of the proposed algorithm is given in section 4. Hints and remarks are summarized in section 5. The applicability of the proposed algorithm to the most commonly employed constitutive relationships is illustrated in section 6. Finally, a few severe numerical tests are illustrated in section 7 to emphasize efficiency and robustness of the proposed method.

**2. Finite volume discretization.** Richards' equation (1) is based on mass conservation. Accordingly, the numerical method to be formulated is required to be exactly mass conservative.

**2.1. Time and spatial approximation.** To solve (1) numerically, the flow region $\Omega$ is partitioned by an *unstructured orthogonal grid* (see, e.g., [5, 6]) consisting of a set of nonoverlapping convex volumes $\Omega_i$, $i = 1, 2, \ldots, N_{\mathrm{V}}$, separated by $M$ internal faces $\Gamma_j$, $j = 1, 2, \ldots, M$. Let $\mathcal{A}_j$ denote the nonzero $j$th face area. Within each control volume a *center* must be identified in such a way that the segment joining the centers of two adjacent volumes and the face shared by the two volumes have a nonempty intersection and are *orthogonal* to each other.

Each control volume $\Omega_i$ may have an arbitrary number of faces. Let $\mathcal{F}_i$ denote the nonempty set of faces of the $i$th volume, with the exclusion of boundary faces. Moreover, let $\wp(i, j)$ be the neighbor of volume $i$ that shares face $j$ with the $i$th

control volume so that $1 \leq \wp(i,j) \leq N_V$ for all $j \in \mathcal{F}_i$. The nonzero distance between the centers of two adjacent volumes which share the $j$th internal face is denoted with $\delta_j$. The discrete variables $\theta_i$ and $\psi_i$ are located at the center of the $i$th control volume. Let $N \leq N_V$ be the number of cells where the pressure head is unknown, and assume that these cells are ordered first for $i = 1, 2, \ldots, N$; whereas, the control volumes where the pressure head is specified as a boundary condition are listed next for $i = N + 1, N + 2, \ldots, N_V$.

Using a fully implicit formulation, at every time step n, the pressure $\psi_i^n$ is assumed to be constant over each control volume $\Omega_i$. Thus, for all $i = 1, 2, \ldots, N$, a finite volume form of (1) is taken to be

$$(2) \qquad \theta_i(\psi_i^n) = \theta_i(\psi_i^{n-1}) + \Delta t \left[ \sum_{j \in \mathcal{F}_i} \mathcal{K}_j^n \frac{\psi_{\wp(i,j)}^n - \psi_i^n}{\delta_j} + \sum_{j \in \mathcal{F}_i} \mathcal{K}_j^n \mathsf{n}_{i,j}^z + S_i^n \right],$$

where $\Delta t$ is the time step size; $\mathcal{K}_j^n = \mathcal{A}_j \max[K_i(\psi_i^n), K_{\wp(i,j)}(\psi_{\wp(i,j)}^n)]$; $\mathsf{n}_{i,j}^z$ is the $j$th projection of the outward normal vector; $S_i = \int_{\Omega_i} S \, d\Omega$ is the $i$th total source, including possibly nonzero boundary fluxes; whereas $\theta_i(\psi)$ is the $i$th fluid volume given by

$$(3) \qquad \theta_i(\psi) = \int_{\Omega_i} \theta(\psi) d\Omega.$$

For given initial conditions $\psi_i^0$, at every time step n $= 1, 2, \ldots$, (2) constitutes a fully nonlinear system of equations to be solved for $\psi_i^n$. To solve (2), one sets $\psi_i^{n,0} = \psi_i^{n-1}$. Then the Picard iterations are taken to be

$$(4) \qquad \theta_i(\psi_i^{n,m}) - \Delta t \sum_{j \in \mathcal{F}_i} \mathcal{K}_j^{n,m-1} \frac{\psi_{\wp(i,j)}^{n,m} - \psi_i^{n,m}}{\delta_j} = b_i^{n,m-1},$$

where

$$b_i^{n,m-1} = \theta_i(\psi_i^{n-1}) + \Delta t \left[ \sum_{j \in \mathcal{F}_i} \mathcal{K}_j^{n,m-1} \mathsf{n}_{i,j}^z + S_i^{n,m-1} \right].$$

At each iteration m $= 1, 2, \ldots$, (4) constitutes a *mildly nonlinear* system [14] for $\psi_i^{n,m}$, with the diagonal nonlinearity being represented by the volumes $\theta_i(\psi_i^{n,m})$. This system of equations represents a consistent and conservative discretization of (1). Hence, regardless of the chosen spatial and temporal accuracy, each Picard iterate $\psi_i^{n,m}$ is a *conservative* approximation for the new pressure. In general, however, an inexact solution of (4) will not be conservative.

In the present study, local and global mass conservation will be enforced at *each* Picard iteration by solving (4) to the best possible accuracy so that the resulting mass balance error will be negligible. Consequently, convergence of the Picard iterations is not essential, but a few steps can be allowed with the only purpose to update the hydraulic conductivity to the nth time level.

By omitting the time index n and the Picard iteration index m, system (4), at every time step and for each Picard iteration, can be written in matrix form as

$$(5) \qquad \boldsymbol{\theta}(\boldsymbol{\psi}) + T\boldsymbol{\psi} = \boldsymbol{b},$$

where $\boldsymbol{\psi} = (\psi_i)$ is the unknown vector, $\boldsymbol{\theta}(\boldsymbol{\psi}) = (\theta_i(\psi_i))$ is a nonnegative vectorial function representing the discrete fluid volumes, $T$ is the diffusive flux matrix, and $\boldsymbol{b}$ is a known vector whose elements are the right-hand side of (4), properly augmented by the known Dirichlet boundary conditions.

**2.2. Moisture capacity.** By denoting with $\theta_r$ the residual moisture content and by $c(\psi) = \frac{\partial \theta}{\partial \psi}$ the (nonnegative) specific moisture capacity, the moisture content can be expressed in terms of $c(\psi)$ as

$$\theta(\psi) = \theta_r + \int_{-\infty}^{\psi} c(\xi)d\xi$$

so that $\theta_s = \theta_r + \int_{-\infty}^{+\infty} c(\xi)d\xi$ is the soil porosity and $\theta(\psi) \leq \theta_s$ for all $\psi \in \mathbb{R}$.

On each control volume, by denoting with $\theta_{r,i}$ the $i$th residual $\theta_{r,i} = \int_{\Omega_i} \theta_r d\Omega$ and with $c_i(\psi)$ the $i$th moisture capacity $c_i(\psi) = \int_{\Omega_i} c(\psi)d\Omega$, the $i$th porous volume is $\theta_{s,i} = \theta_{r,i} + \int_{-\infty}^{+\infty} c_i(\xi)d\xi$, and the $i$th fluid volume can be written as

$$(6) \qquad \theta_i(\psi) = \theta_{r,i} + \int_{-\infty}^{\psi} c_i(\xi)d\xi.$$

For all $i$, the following assumptions are made on the cell moisture capacity $c_i(\psi)$.

*Assumption* C1: $c_i(\psi)$ is defined for all $\psi \in \mathbb{R}$ and is a nonnegative function with bounded variations.

*Assumption* C2: There exist $\psi_i^* \in \mathbb{R}$ such that $c_i(\psi)$ is strictly positive and nondecreasing in $(-\infty, \psi_i^*)$ and nonincreasing in $(\psi_i^*, +\infty)$.

Thus, $c_i(\psi) = \frac{d\theta_i(\psi)}{d\psi} \geq 0$ and $\theta_{r,i} < \theta_i(\psi) \leq \theta_{s,i}$ for all $\psi \in \mathbb{R}$.

As will be seen in section 6, the most commonly used constitutive equations, relating the moisture content to the pressure head, satisfy Assumptions C1–C2.

Since $c_i(\psi)$ are nonnegative functions with bounded variations, they are *almost everywhere differentiable*, admit only *discontinuities of the first kind*, and can be expressed as the difference of two nonnegative, nondecreasing, and bounded functions (the *Jordan decomposition* [8]), say $p_i(\psi)$ and $q_i(\psi)$, so that $c_i(\psi) = p_i(\psi) - q_i(\psi) \geq 0$ and $0 \leq q_i(\psi) \leq p_i(\psi)$ for all $\psi \in \mathbb{R}$. When $c_i(\psi)$ satisfies Assumptions C1–C2, the corresponding Jordan decomposition is given by

$$(7) \qquad \begin{cases} p_i(\psi) = c_i(\psi), \quad q_i(\psi) = 0 & \text{if } \psi \leq \psi_i^*, \\[2mm] p_i(\psi) = c_i(\psi_i^*), \quad q_i(\psi) = p_i(\psi) - c_i(\psi) & \text{if } \psi > \psi_i^*. \end{cases}$$

Furthermore, the fluid volumes $\boldsymbol{\theta}(\boldsymbol{\psi})$ can be written as $\boldsymbol{\theta}(\boldsymbol{\psi}) = \boldsymbol{\theta}_1(\boldsymbol{\psi}) - \boldsymbol{\theta}_2(\boldsymbol{\psi})$, where each component of $\boldsymbol{\theta}_1(\boldsymbol{\psi})$ and $\boldsymbol{\theta}_2(\boldsymbol{\psi})$, respectively, is given by

$$(8) \qquad \theta_{1,i}(\psi) = \theta_{r,i} + \int_{-\infty}^{\psi} p_i(\xi)d\xi \qquad \text{and} \qquad \theta_{2,i}(\psi) = \int_{-\infty}^{\psi} q_i(\xi)d\xi$$

or, equivalently,

$$(9) \qquad \begin{cases} \theta_{1,i}(\psi) = \theta_i(\psi), \qquad\qquad\qquad \theta_{2,i}(\psi) = 0 & \text{if } \psi \leq \psi_i^*, \\[2mm] \theta_{1,i}(\psi) = \theta_i(\psi_i^*) + c_i(\psi_i^*)(\psi - \psi_i^*), \quad \theta_{2,i}(\psi) = \theta_{1,i}(\psi) - \theta_i(\psi) & \text{if } \psi > \psi_i^* \end{cases}$$

so that $\theta_i(\psi) = \theta_{1,i}(\psi) - \theta_{2,i}(\psi)$, $p_i(\psi) = \frac{d\theta_{1,i}(\psi)}{d\psi}$ and $q_i(\psi) = \frac{d\theta_{2,i}(\psi)}{d\psi}$, $i = 1, 2, \ldots, N$.

Let $C(\boldsymbol{\psi})$, $P(\boldsymbol{\psi})$ and $Q(\boldsymbol{\psi})$ denote diagonal matrices whose diagonal entries are $c_i(\psi)$, $p_i(\psi)$ and $q_i(\psi)$, respectively. Thus, $C(\boldsymbol{\psi}) = P(\boldsymbol{\psi}) - Q(\boldsymbol{\psi})$ represents the Jacobian of $\boldsymbol{\theta}(\boldsymbol{\psi})$ almost everywhere; $P(\boldsymbol{\psi})$ is almost everywhere the Jacobian of $\boldsymbol{\theta}_1(\boldsymbol{\psi})$; and $Q(\boldsymbol{\psi})$ is almost everywhere the Jacobian of $\boldsymbol{\theta}_2(\boldsymbol{\psi})$.

Hereafter, let $\mathbf{0}$ and O denote the zero vector and the zero matrix of appropriate size, respectively. The following easy property is stated here for later reference.

LEMMA 1. *Let $c_i(\psi)$ satisfy the Assumptions* C1 *and* C2, *and let $p_i(\psi)$ and $q_i(\psi)$ be the Jordan decomposition of $c_i(\psi)$, $i = 1, 2, \ldots, N$. For all $\boldsymbol{\varphi}, \boldsymbol{\psi} \in \mathbb{R}^N$ one has*

$$P(\boldsymbol{\psi})(\boldsymbol{\psi} - \boldsymbol{\varphi}) - [\boldsymbol{\theta}_1(\boldsymbol{\psi}) - \boldsymbol{\theta}_1(\boldsymbol{\varphi})] \geq \mathbf{0},$$
$$Q(\boldsymbol{\psi})(\boldsymbol{\psi} - \boldsymbol{\varphi}) - [\boldsymbol{\theta}_2(\boldsymbol{\psi}) - \boldsymbol{\theta}_2(\boldsymbol{\varphi})] \geq \mathbf{0}.$$

**2.3. Diffusive flux matrix.** Without loss of generality, it will be assumed that matrix $T$ is irreducible. This may not be the case when, at any time, two or more subdomains are not connected by strictly positive diffusive flux coefficients. In such a circumstance the considerations that follow apply separately to each such subdomain where the corresponding matrix $T$ is irreducible.

From (4) one sees that the main diagonal elements of matrix $T$ are $t_{i,i} = \Delta t \sum_{j \in \mathcal{F}_i} \mathcal{K}_j^{\mathrm{n,m-1}}/\delta_j > 0$; whereas the possibly nonzero off-diagonal elements of $T$ are the coefficients of $\psi_{i,\wp(i,j)}$, corresponding to $\wp(i,j) \leq N$, and are given by $t_{i,\wp(i,j)} = -\Delta t \, \mathcal{K}_j^{\mathrm{n,m-1}}/\delta_j \leq 0$, $j \in \mathcal{F}_i$, so that matrix $T$ is symmetric.

In case that the Neumann boundary conditions are specified everywhere along each boundary face, then $N = N_{\mathrm{V}}$, and therefore, $\wp(i,j) \leq N$ for all $i$ and for all $j \in \mathcal{F}_i$. Consequently, $\sum_{\nu=1}^{N} t_{i,\nu} = 0$ for all $i = 1, 2, \ldots, N$, implying that matrix $T$ is singular and positive semidefinite. Moreover,

$$\sum_{i=1}^{N} (T\boldsymbol{\psi})_i = 0 \text{ for all } \boldsymbol{\psi} \in \mathbb{R}^N. \tag{10}$$

When the Diriclet boundary condition is specified in at least one control volume, then $N < N_{\mathrm{V}}$, and the resulting matrix $T$ is a nonsingular, positive definite $M$-matrix.

In general, in order to account for Dirichlet, Neumann, and mixed Neumann–Dirichlet boundary conditions, matrix $T$ is assumed to be symmetric and (at least) positive semidefinite, satisfying either one of the following properties:

T1      $T$ is a symmetric $M$-matrix (i.e., a Stieltjes matrix), or

T2      $T$ is singular, and $T+D$ is a Stieltjes matrix for all diagonal matrices $D \gneqq \mathrm{O}$ (i.e., $D \geq \mathrm{O}$ and $D \neq \mathrm{O}$).

When $T$ is T2, the following compatibility assumption is required on $\boldsymbol{b}$:

$$\sum_{i=1}^{N} \theta_{r,i} < \sum_{i=1}^{N} b_i < \sum_{i=1}^{N} \theta_{s,i}. \tag{11}$$

Inequalities (11) are requirements for system (5) to be physically and mathematically compatible. This assumption, in fact, states that when the flow boundary conditions are specified everywhere along the boundary faces, the resulting total fluid volume within the system must be larger than the total residual volume and smaller than the maximum water volume allowed by soil porosity [4].

**3. Nested iterations.** Assume that matrix $T$ satisfies either T1 or T2. If $T$ satisfies T2, assume also that inequalities (11) hold true. If the moisture capacities satisfy both Assumptions C1 and C2, then system (5) can be written as

$$(12) \qquad \boldsymbol{\theta}_1(\boldsymbol{\psi}) - \boldsymbol{\theta}_2(\boldsymbol{\psi}) + T\boldsymbol{\psi} = \boldsymbol{b}.$$

By choosing $\boldsymbol{\psi}^0 \leq \boldsymbol{\psi}^*$, a sequence of *outer* iterates $\{\boldsymbol{\psi}^k\}$ is derived from (12) by linearizing $\boldsymbol{\theta}_2(\boldsymbol{\psi})$ as follows,

$$\boldsymbol{\theta}_1(\boldsymbol{\psi}^k) - [\boldsymbol{\theta}_2(\boldsymbol{\psi}^{k-1}) + Q(\boldsymbol{\psi}^{k-1})(\boldsymbol{\psi}^k - \boldsymbol{\psi}^{k-1})] + T\boldsymbol{\psi}^k = \boldsymbol{b},$$

so that the outer iterates are solutions of the following mildly nonlinear systems:

$$(13) \qquad \boldsymbol{\theta}_1(\boldsymbol{\psi}^k) + \left(T - Q^{k-1}\right)\boldsymbol{\psi}^k = \boldsymbol{d}^{k-1}, \qquad k = 1, 2, \ldots,$$

where $Q^{k-1} = Q(\boldsymbol{\psi}^{k-1})$ and $\boldsymbol{d}^{k-1} = \boldsymbol{b} + \boldsymbol{\theta}_2(\boldsymbol{\psi}^{k-1}) - Q^{k-1}\boldsymbol{\psi}^{k-1}$.

Next, for all $k = 1, 2, \ldots$, by setting $\boldsymbol{\psi}^{k,0} = \boldsymbol{\psi}^{k-1}$, a sequence of *inner* iterates $\{\boldsymbol{\psi}^{k,\ell}\}$ is derived from (13) by linearizing $\boldsymbol{\theta}_1(\boldsymbol{\psi})$ as follows,

$$[\boldsymbol{\theta}_1(\boldsymbol{\psi}^{k,\ell-1}) + P(\boldsymbol{\psi}^{k,\ell-1})(\boldsymbol{\psi}^{k,\ell} - \boldsymbol{\psi}^{k,\ell-1})] + \left(T - Q^{k-1}\right)\boldsymbol{\psi}^{k,\ell} = \boldsymbol{d}^{k-1},$$

so that the inner iterates are determined from the following linear systems:

$$(14) \qquad \left(P^{k,\ell-1} + T - Q^{k-1}\right)\boldsymbol{\psi}^{k,\ell} = \boldsymbol{f}^{k,\ell-1}, \quad \ell = 1, 2, \ldots,$$

where $P^{k,\ell-1} = P(\boldsymbol{\psi}^{k,\ell-1})$ and $\boldsymbol{f}^{k,\ell-1} = \boldsymbol{d}^{k-1} - \boldsymbol{\theta}_1(\boldsymbol{\psi}^{k,\ell-1}) + P^{k,\ell-1}\boldsymbol{\psi}^{k,\ell-1}$.

From (13) one can derive the $k$th outer residual, namely $\boldsymbol{r}^k = \boldsymbol{\theta}(\boldsymbol{\psi}^k) + T\boldsymbol{\psi}^k - \boldsymbol{b}$, which, by Lemma 1, satisfies

$$(15) \qquad \boldsymbol{r}^k = -\{Q^{k-1}(\boldsymbol{\psi}^{k-1} - \boldsymbol{\psi}^k) - [\boldsymbol{\theta}_2(\boldsymbol{\psi}^{k-1}) - \boldsymbol{\theta}_2(\boldsymbol{\psi}^k)]\} \leq \boldsymbol{0}.$$

Thus, a straightforward stopping criterion for the outer iterations is $\|\boldsymbol{r}^k\| < \epsilon$, where $\epsilon$ is a prefixed tolerance representing the maximum mass balance error allowed.

Similarly, the $(k,\ell)$th inner residual, $\boldsymbol{r}^{k,\ell} = \boldsymbol{\theta}_1(\boldsymbol{\psi}^{k,\ell}) + \left(T - Q^{k-1}\right)\boldsymbol{\psi}^{k,\ell} - \boldsymbol{d}^{k-1}$, can be derived from (14), and, by Lemma 1, satisfies

$$(16) \qquad \boldsymbol{r}^{k,\ell} = P^{k,\ell-1}(\boldsymbol{\psi}^{k,\ell-1} - \boldsymbol{\psi}^{k,\ell}) - [\boldsymbol{\theta}_1(\boldsymbol{\psi}^{k,\ell-1}) - \boldsymbol{\theta}_1(\boldsymbol{\psi}^{k,\ell})] \geq \boldsymbol{0}.$$

Thus, the stopping criterion for the inner iterations is $\|\boldsymbol{r}^{k,\ell-1}\| < \epsilon$.

The above method is a nonlinear extension of the Newton-type algorithm presented in [4] for piecewise linear systems. This method is summarized into Algorithm 1.

ALGORITHM 1.

---

Choose $\boldsymbol{\psi}^0 \leq \boldsymbol{\psi}^*$
Do $k = 1, 2, \ldots$
   Set $\boldsymbol{\psi}^{k,0} = \boldsymbol{\psi}^{k-1}$
   Do $\ell = 1, 2, \ldots$
      Solve $\left(P^{k,\ell-1} + T - Q^{k-1}\right)\boldsymbol{\psi}^{k,\ell} = \boldsymbol{f}^{k,\ell-1}$
      If $\|\boldsymbol{r}^{k,\ell}\| < \epsilon$, then set $\boldsymbol{\psi}^k = \boldsymbol{\psi}^{k,\ell}$ and exit
   End do
   If $\|\boldsymbol{r}^k\| < \epsilon$, then set $\boldsymbol{\psi} = \boldsymbol{\psi}^k$ and exit
End do

---

**4. Convergence analysis.** Algorithm 1 is well defined and converges to the exact solution of system (5). To prove this statement, some auxiliary results are outlined first.

LEMMA 2. *For any $k \geq 1$, let $\boldsymbol{\psi}^{k,\ell}$ and $\boldsymbol{\psi}^{k,\ell+1}$ be two subsequent inner iterates obtained from* (14). *If $P^{k,\ell} + T - Q^{k-1}$ is an $M$-matrix, then*

$$(17) \qquad \boldsymbol{\psi}^{k,\ell+1} \leq \boldsymbol{\psi}^{k,\ell}.$$

*Proof.* Consider two subsequent inner iterates (14), which can be written as

$$\left(P^{k,\ell} + T - Q^{k-1}\right) \boldsymbol{\psi}^{k,\ell+1} - P^{k,\ell}\boldsymbol{\psi}^{k,\ell} + \boldsymbol{\theta}_1(\boldsymbol{\psi}^{k,\ell}) = \boldsymbol{d}^{k-1},$$
$$\left(P^{k,\ell-1} + T - Q^{k-1}\right) \boldsymbol{\psi}^{k,\ell} - P^{k,\ell-1}\boldsymbol{\psi}^{k,\ell-1} + \boldsymbol{\theta}_1(\boldsymbol{\psi}^{k,\ell-1}) = \boldsymbol{d}^{k-1}.$$

By equating the left-hand sides, one obtains

$$\left(P^{k,\ell} + T - Q^{k-1}\right) \left(\boldsymbol{\psi}^{k,\ell+1} - \boldsymbol{\psi}^{k,\ell}\right) = -\boldsymbol{r}^{k,\ell} \leq \boldsymbol{0};$$

hence, because $P^{k,\ell} + T - Q^{k-1}$ is an $M$-matrix, one has $\left(P^{k,\ell} + T - Q^{k-1}\right)^{-1} > \mathrm{O}$, and consequently, inequality (17) is implied.  □

LEMMA 3. *For $k \geq 1$, let $\boldsymbol{\psi}^k$ be the $k$th outer iterate obtained from* (13), *and let $\boldsymbol{\psi}^{k+1,\ell}$ be the subsequent $\ell$th inner iterate satisfying* (14). *If $P^{k+1,\ell-1} + T - Q^k$ is an $M$-matrix, then*

$$(18) \qquad \boldsymbol{\psi}^k \leq \boldsymbol{\psi}^{k+1,\ell}.$$

*Proof.* From (14) one has that $\boldsymbol{\psi}^{k+1,\ell}$ is solution of

$$(19) \qquad \left(P^{k+1,\ell-1} + T - Q^k\right) \boldsymbol{\psi}^{k+1,\ell} - P^{k+1,\ell-1}\boldsymbol{\psi}^{k+1,\ell-1} + \boldsymbol{\theta}_1(\boldsymbol{\psi}^{k+1,\ell-1}) = \boldsymbol{d}^k;$$

whereas $\boldsymbol{\psi}^k$ is solution of (13). By subtracting (13) from (19), one has

$$\left(P^{k+1,\ell-1} + T - Q^k\right) \left(\boldsymbol{\psi}^{k+1,\ell} - \boldsymbol{\psi}^k\right) = \boldsymbol{r}^{k+1,\ell} - \boldsymbol{r}^k \geq \boldsymbol{0}.$$

Thus, since $P^{k+1,\ell-1} + T - Q^k$ is an $M$-matrix, one has $\left(P^{k+1,\ell-1} + T - Q^k\right)^{-1} > \mathrm{O}$, and consequently, inequality (18) is implied.  □

The following results show that Algorithm 1 is well defined, monotone, and converging in both cases that $T$ is either T1 or T2. The two cases are analyzed separately next.

THEOREM 1. *Let $\boldsymbol{\theta}(\boldsymbol{\psi})$ be given by* (6), *and let $T$ be an irreducible diffusive flux matrix satisfying* T1. *If the moisture capacity satisfies both Assumptions* C1 *and* C2, *then Algorithm* 1 *is well defined and the inner and the outer iterates are monotone. Specifically, for all $k = 1, 2, \ldots$, one has*

$$(20) \qquad \boldsymbol{\psi}^{k,\ell+1} \leq \boldsymbol{\psi}^{k,\ell}, \quad \ell = 1, 2, \ldots,$$
$$(21) \qquad \boldsymbol{\psi}^k \leq \boldsymbol{\psi}^{k+1,\ell}, \quad \ell = 1, 2, \ldots,$$
$$(22) \qquad \boldsymbol{\psi}^k \leq \boldsymbol{\psi}^{k+1}.$$

*Moreover, the outer iterates $\boldsymbol{\psi}^k$ converge to the exact solution of system* (5).

*Proof.* Assumptions C1 and C2 are required to formulate Algorithm 1. Moreover, since $c_i(\psi) > 0$ for all $\psi \leq \psi_i^*$, one has $p_i(\psi) > 0$ for all $i$ and for all $\psi \in \mathbb{R}$. Then one proceeds by double induction.

$k = 1$: In this case one has $Q^0 = Q(\boldsymbol{\psi}^0 \leq \boldsymbol{\psi}^*) = \mathrm{O}$ and $P^{1,\ell-1} \gneqq \mathrm{O}$. Thus, $P^{1,\ell-1} + T - Q^0$ are $M$-matrices for all $\ell = 1, 2, \ldots$. Consequently, Lemma 2 implies inequality (20). Hence, the sequence $\{\boldsymbol{\psi}^{1,\ell}\}$ is monotonically decreasing, and accordingly, the sequence $\{\boldsymbol{\theta}_1(\boldsymbol{\psi}^{1,\ell})\}$ is monotonically decreasing as well. Moreover, $T^{-1} > \mathrm{O}$ componentwise. To prove boundedness of the first inner iterates, let $\boldsymbol{\psi}^b = T^{-1}[\boldsymbol{b} - \boldsymbol{\theta}_1(\boldsymbol{\psi}^{1,1})]$ so that from (14) one has

$$\boldsymbol{\psi}^{1,\ell} - \boldsymbol{\psi}^b = T^{-1}[\boldsymbol{\theta}_1(\boldsymbol{\psi}^{1,1}) - \boldsymbol{\theta}_1(\boldsymbol{\psi}^{1,\ell}) + \boldsymbol{r}^{1,\ell}] \geq \boldsymbol{0}.$$

Hence, the monotonically decreasing sequence $\{\boldsymbol{\psi}^{1,\ell}\}$ is bounded below by $\boldsymbol{\psi}^b$ and is therefore converging to, say, $\bar{\boldsymbol{\psi}}^1$. Consequently, $P^{1,\ell} \to P(\bar{\boldsymbol{\psi}}^1)$ and $\boldsymbol{\theta}_1(\boldsymbol{\psi}^{1,\ell}) \to \boldsymbol{\theta}_1(\bar{\boldsymbol{\psi}}^1)$. Thus, (16) yields $\boldsymbol{r}^{1,\ell} \to \boldsymbol{0}$, implying that $\bar{\boldsymbol{\psi}}^1 = \boldsymbol{\psi}^1$ is the exact solution of (13).

$k > 1$: One assumes that the $(k-1)$th outer cycle has been successfully completed so that $\boldsymbol{\psi}^{k-1}$ is a solution of (13), and $Q^{k-1} \leq P^{k-1}$.

　$\ell = 1$: Since $T$ is T1, one has that $P^{k,0} + T - Q^{k-1} = P^{k-1} + T - Q^{k-1}$ is an M-matrix. Hence, $\boldsymbol{\psi}^{k,1}$ can be uniquely determined from (14).

　$\ell > 1$: One assumes that $P^{k,\ell-2} + T - Q^{k-1}$ is an $M$-matrix so that $\boldsymbol{\psi}^{k,\ell-1}$ can be uniquely determined from (14), and because of Lemma 3, it satisfies $\boldsymbol{\psi}^{k,\ell-1} \geq \boldsymbol{\psi}^{k-1}$. Consequently, $P^{k,\ell-1} \geq P^{k-1} \geq Q^{k-1}$, and hence, the iteration matrix $P^{k,\ell-1} + T - Q^{k-1}$ is an M-matrix. Thus, $\boldsymbol{\psi}^{k,\ell}$ can be uniquely determined from (14).

This shows that $P^{k,\ell-1} + T - Q^{k-1}$ are M-matrices for all $\ell = 1, 2, \ldots$. Hence, Lemma 2 implies inequality (20) and Lemma 3 yields inequality (21). This latter inequality shows that the monotonically decreasing sequence $\{\boldsymbol{\psi}^{k,\ell}\}$ is bounded below by $\boldsymbol{\psi}^{k-1}$ and is therefore converging to, say, $\bar{\boldsymbol{\psi}}^k$. Consequently, $P^{k,\ell} \to P(\bar{\boldsymbol{\psi}}^k)$ and $\boldsymbol{\theta}_1(\boldsymbol{\psi}^{k,\ell}) \to \boldsymbol{\theta}_1(\bar{\boldsymbol{\psi}}^k)$. Thus, (16) yields $\boldsymbol{r}^{k,\ell} \to \boldsymbol{0}$, implying that $\bar{\boldsymbol{\psi}}^k = \boldsymbol{\psi}^k$ is the exact solution of (13). Finally, inequality (22) is an immediate consequence of (21).

To prove convergence, one needs to show that the monotonically increasing sequence $\{\boldsymbol{\psi}^k\}$ is bounded above. To this purpose, let $\boldsymbol{\psi}^t = T^{-1}[\boldsymbol{b} - \boldsymbol{\theta}_r]$, and recall that $\boldsymbol{\theta}(\boldsymbol{\psi}^k) > \boldsymbol{\theta}_r$ and $\boldsymbol{r}^k \leq \boldsymbol{0}$. Hence, from (13) one has

$$\boldsymbol{\psi}^t - \boldsymbol{\psi}^k = T^{-1}[\boldsymbol{\theta}(\boldsymbol{\psi}^k) - \boldsymbol{\theta}_r - \boldsymbol{r}^k] > \boldsymbol{0}.$$

Thus, the monotonically increasing sequence $\{\boldsymbol{\psi}^k\}$ is bounded above by $\boldsymbol{\psi}^t$ and is therefore converging to, say, $\bar{\boldsymbol{\psi}}$. Consequently, $Q^k \to Q(\bar{\boldsymbol{\psi}})$ and $\boldsymbol{\theta}_2(\boldsymbol{\psi}^k) \to \boldsymbol{\theta}_2(\bar{\boldsymbol{\psi}})$. Thus, (15) yields $\boldsymbol{r}^k \to \boldsymbol{0}$, implying that $\bar{\boldsymbol{\psi}}$ is the exact solution of (5). $\quad\square$

THEOREM 2. *Let $\boldsymbol{\theta}(\boldsymbol{\psi})$ be given by (6), and let $T$ be an irreducible diffusive flux matrix satisfying* T2. *Assume also that $\boldsymbol{b}$ satisfies inequalities* (11). *If the moisture capacity satisfies both Assumptions* C1 *and* C2, *then Algorithm 1 is well defined, monotone, and convergent. Specifically, for all $k = 1, 2, \ldots$, inequalities* (20)–(22) *hold true.*

*Proof.* Assumptions C1 and C2 are required to formulate Algorithm 1. Moreover, since $c_i(\psi) > 0$ for all $\psi \leq \psi_i^*$, one has $p_i(\psi) > 0$ for all $i$ and for all $\psi \in \mathbb{R}$. Furthermore, if $Q(\boldsymbol{\psi}) = P(\boldsymbol{\psi})$, then $\boldsymbol{\theta}(\boldsymbol{\psi}) = \boldsymbol{\theta}_s$. The proof proceeds by double induction.

$k = 1$: In this case one has $Q^0 = Q(\boldsymbol{\psi}^0 \leq \boldsymbol{\psi}^*) = \mathrm{O}$ and $P^{1,\ell-1} \gneqq \mathrm{O}$. Thus, $P^{1,\ell-1} + T - Q^0$ are $M$-matrices for all $\ell = 1, 2, \ldots$. Consequently, Lemma 2 implies inequality (20). Hence, the sequence $\{\boldsymbol{\psi}^{1,\ell}\}$ is monotonically decreasing. To

prove boundness of the first inner iterates, consider the first inner residuals $\boldsymbol{r}^{1,\ell}$ that are given by

$$\boldsymbol{\theta}_1(\boldsymbol{\psi}^{1,\ell}) + T\boldsymbol{\psi}^{1,\ell} - \boldsymbol{b} = \boldsymbol{r}^{1,\ell} \geq \boldsymbol{0}.$$

Moreover, the sequence $\{\boldsymbol{\theta}_1(\boldsymbol{\psi}^{1,\ell})\}$ is monotonically decreasing and bounded below by $\boldsymbol{\theta}_r$, hence, converging to, say, $\boldsymbol{\theta}_1^1$. Furthermore, by (10) one has that the sum of the inner residuals, namely,

$$\sum_{i=1}^{N} [\theta_{1,i}(\psi_i^{1,\ell}) - b_i] \geq 0,$$

is nonnegative, monotonically decreasing, and hence, bounded. Thus, the sequence $\{\boldsymbol{r}^{1,\ell}\}$ converges to, say, $\boldsymbol{r}^1$. Consequently, the sequence $\{T\boldsymbol{\psi}^{1,\ell}\}$ is itself bounded and converges to $\boldsymbol{b} + \boldsymbol{r}^1 - \boldsymbol{\theta}_1^1$. Since $T$ is irreducible, the monotonically decreasing sequence $\{\boldsymbol{\psi}^{1,\ell}\}$ is either bounded or componentwise unbounded. Unbounded growth of $\{\boldsymbol{\psi}^{1,\ell}\}$ to $-\infty$, however, yields $\{\boldsymbol{\theta}_1(\boldsymbol{\psi}^{1,\ell})\} \to \boldsymbol{\theta}_r$, and consequently, the first inequality (11) implies

$$\lim_{\ell \to \infty} \sum_{i=1}^{N} [\theta_{1,i}(\psi_i^{1,\ell}) - b_i] = \sum_{i=1}^{N} [\theta_{r,i} - b_i] < 0,$$

which contradicts the previous inequality. Hence, the monotonically decreasing sequence $\boldsymbol{\psi}^{1,\ell}$ must be bounded and is therefore converging to, say, $\bar{\boldsymbol{\psi}}^1$. Consequently, $P^{1,\ell} \to P(\bar{\boldsymbol{\psi}}^1)$ and $\boldsymbol{\theta}_1(\boldsymbol{\psi}^{1,\ell}) \to \boldsymbol{\theta}_1(\bar{\boldsymbol{\psi}}^1)$. Thus, (16) yields $\boldsymbol{r}^{1,\ell} \to \boldsymbol{0}$, implying that $\bar{\boldsymbol{\psi}}^1 = \boldsymbol{\psi}^1$ is the exact solution of (13).

$k > 1$: One assumes that the $(k-1)$th outer cycle has been successfully completed so that $\boldsymbol{\psi}^{k-1}$ is a solution of (13) and $Q^{k-1} \leq P^{k-1}$. Moreover, $Q^{k-1} \neq P^{k-1}$. In fact, if $Q^{k-1} = P^{k-1}$, one has $\boldsymbol{\theta}(\boldsymbol{\psi}^{k-1}) = \boldsymbol{\theta}_s$, and (13) yields

$$\boldsymbol{\theta}_s + T\boldsymbol{\psi}^{k-1} - \boldsymbol{b} = \boldsymbol{r}^{k-1} \leq \boldsymbol{0},$$

where, by (10) and inequality (11), the sum of the left-hand side is strictly positive. This is a contradiction, implying $Q^{k-1} \lneqq P^{k-1}$.

$\ell = 1$: Since $Q^{k-1} \lneqq P^{k-1}$, one has that $P^{k,0} + T - Q^{k-1} = P^{k-1} + T - Q^{k-1}$ is an M-matrix. Consequently, $\boldsymbol{\psi}^{k,1}$ is uniquely determined from (14).

$\ell > 1$: One assumes that $P^{k,\ell-2} + T - Q^{k-1}$ is an $M$-matrix so that $\boldsymbol{\psi}^{k,\ell-1}$ can be uniquely determined from (14), and because of Lemma 3, it satisfies $\boldsymbol{\psi}^{k,\ell-1} \geq \boldsymbol{\psi}^{k-1}$. Consequently, $P^{k,\ell-1} \geq P^{k-1} \gneqq Q^{k-1}$ implies that the iteration matrix $P^{k,\ell-1} + T - Q^{k-1}$ is an M-matrix. Thus, $\boldsymbol{\psi}^{k,\ell}$ can be uniquely determined from (14).

This shows that $P^{k,\ell-1} + T - Q^{k-1}$ are $M$-matrices for all $\ell = 1, 2, \ldots$. Hence, Lemma 2 implies inequality (20), and Lemma 3 yields inequality (21). This latter inequality shows that the monotonically decreasing sequence $\{\boldsymbol{\psi}^{k,\ell}\}$ is bounded below by $\boldsymbol{\psi}^{k-1}$ and is therefore converging to, say, $\bar{\boldsymbol{\psi}}^k$. Consequently, $P^{k,\ell} \to P(\bar{\boldsymbol{\psi}}^k)$ and $\boldsymbol{\theta}_1(\boldsymbol{\psi}^{k,\ell}) \to \boldsymbol{\theta}_1(\bar{\boldsymbol{\psi}}^k)$. Thus, (16) yields $\boldsymbol{r}^{k,\ell} \to \boldsymbol{0}$, implying that $\bar{\boldsymbol{\psi}}^k = \boldsymbol{\psi}^k$ is the exact solution of (13). Finally, inequality (22) is an immediate consequence of (21).

To prove convergence, one also needs to show the boundness of the monotonically increasing sequence $\{\boldsymbol{\psi}^k\}$. To this purpose, consider that the outer residuals satisfy

$$\boldsymbol{\theta}(\boldsymbol{\psi}^k) + T\boldsymbol{\psi}^k - \boldsymbol{b} = \boldsymbol{r}^k \leq \boldsymbol{0}.$$

Thus, the $k$th outer residual is componentwise nonpositive. Moreover, $\boldsymbol{\theta}(\boldsymbol{\psi}^k)$ is monotonically increasing and bounded above by $\boldsymbol{\theta}_s$ and, hence, is converging to, say, $\bar{\boldsymbol{\theta}}$. Furthermore, by (10) one has that the sum of the residuals, namely,

$$\sum_{i=1}^{N}[\theta_i(\psi_i^k) - b_i] \leq 0,$$

is nonpositive, monotonically increasing, and hence, bounded. This proves that the sequence $\{\boldsymbol{r}^k\}$ converges to, say, $\bar{\boldsymbol{r}}$. Consequently, the sequence $\{T\boldsymbol{\psi}^k\}$ is itself bounded and converges to $\boldsymbol{b} + \bar{\boldsymbol{r}} - \bar{\boldsymbol{\theta}}$. Since $T$ is irreducible, the monotonically increasing sequence $\{\boldsymbol{\psi}^k\}$ is either bounded or componentwise unbounded. Unbounded growth of $\{\boldsymbol{\psi}^k\}$ to $+\infty$ componentwise, however, yields $\{\boldsymbol{\theta}(\boldsymbol{\psi}^k)\} \to \boldsymbol{\theta}_s$, and consequently, the second inequality (11) yields

$$\lim_{k\to\infty} \sum_{i=1}^{N}[\theta_i(\psi_i^k) - b_i] = \sum_{i=1}^{N}[\theta_{s,i} - b_i] > 0,$$

which contradicts the previous inequality. Hence, the monotonically increasing sequence $\{\boldsymbol{\psi}^k\}$ must be bounded and is therefore converging to, say, $\bar{\boldsymbol{\psi}}$. Consequently, $Q^k \to Q(\bar{\boldsymbol{\psi}})$ and $\boldsymbol{\theta}_2(\boldsymbol{\psi}^k) \to \boldsymbol{\theta}_2(\bar{\boldsymbol{\psi}})$. Thus, (15) yields that $\boldsymbol{r}^k$ converges to $\boldsymbol{0}$, implying that $\bar{\boldsymbol{\psi}}$ is the exact solution of system (5). $\square$

The above results, other than ensuring convergence of the proposed method, also indicate the convergence behavior of each inner and outer iterate. In particular, inequalities (20)–(22) suggest some crucial checks to be considered within the implementation of Algorithm 1.

**5. Hints and remarks.** The proposed numerical method has a number of interesting properties which are listed below.

**5.1. Initial guess.** In general, faster convergence is attained by taking as initial guess for an iterative procedure the solution from an outer iteration loop or from the previous time step. Therefore, in order to comply with the requirements $\boldsymbol{\psi}^0 \leq \boldsymbol{\psi}^*$ of Algorithm 1, a suggested choice for the initial guess, that takes advantage of the known solution from the previous Picard iteration, is as follows:

$$(23) \qquad \boldsymbol{\psi}^0 = \min(\boldsymbol{\psi}^*, \boldsymbol{\psi}^{\mathrm{n,m-1}}).$$

Note that because $T$ is symmetric and (at least) positive semidefinite, one has that $P^{k,\ell-1} + T - Q^{k-1}$ are positive definite for all $k$ and $\ell$. Hence, each linear system in the inner iterates (14) can be efficiently solved by a preconditioned conjugate gradient method (see, e.g., [13]).

**5.2. Soil storativity.** In general, Richards' equation (1) does not account for the effects of specific storage. Consequently, it cannot be used to model a variety of variably saturated flow problems, including many transient drainage and seepage phenomena in large domains [1, 9]. When accounting for the effects of specific storage,

the governing differential equation is an extension of the classical Richards' equation and is given by

$$(24) \qquad S_s \frac{\theta}{\theta_s} \frac{\partial \psi}{\partial t} + \frac{\partial \theta(\psi)}{\partial t} = \nabla \cdot [K(\psi)\nabla(\psi + z)] + S,$$

where $S_s$ is the specific storativity. The Picard iterations applied to a finite volume discretization of (24) leads to a mildly nonlinear system which, for each $m = 1, 2, \ldots$, can be written in the form (5) with a different right-hand side and with an additional contribution to the main diagonal of $T$ which now is T1. Hence, Algorithm 1 applies to solve a finite volume discretization of (24) in a straightforward manner.

**5.3. One step convergence.** In unsaturated flows, it often happens that Algorithm 1 converges in only one outer iteration. This is the case when the solution of system (5) satisfies $\bar{\psi} \leq \psi^*$. In fact, since $\psi^0 \leq \psi^*$, and because $\theta_2(\psi) = \mathbf{0}$ and $Q(\psi) = \mathrm{O}$ for all $\psi \leq \psi^*$, the linearization performed on $\theta_2(\psi)$ to derive the outer iterations (13) is exact. Consequently, no further outer iterations behind $k = 1$ are required. In this case system (5) is being solved by only one set of inner iterations that coincide with the classical Newton method.

Similarly, if the solution of system (5) satisfies $\bar{\psi} \geq \psi^*$, one has $\theta_1(\psi) = \theta_1(\psi^*) + P(\psi^*)(\psi - \psi^*)$ and $P(\psi) = P(\psi^*)$ for all $\psi \geq \psi^*$. Hence, the linearization performed on $\theta_1(\psi)$ to derive the inner iterations (14) is exact, and consequently, no further iterations behind $\ell = 1$ are required by the inner loops of Algorithm 1. In this case system (5) is being solved by one set of outer iterations with only one inner iteration each. Also, in this case Algorithm 1 simplifies to the classical Newton method.

**5.4. Special cases.** It is worthwhile noting that, if the inner cycle of Algorithm 1 is deliberately restricted to only one iteration, then this method simplifies to the classical Newton method for system (5):

$$(25) \qquad \theta(\psi^{k-1}) + C(\psi^{k-1})(\psi^k - \psi^{k-1}) + T\psi^k = b.$$

Moreover, if both the inner and the outer cycles are restricted to only one iteration, and if the time and the Picard indices are resumed in (25), one has

$$(26) \qquad \theta(\psi^{n,m-1}) + C(\psi^{n,m-1})(\psi^{n,m} - \psi^{n,m-1}) + T^{n,m-1}\psi^{n,m} = b^{n,m-1},$$

which is the linearization approach proposed by Celia, Bouloutas, and Zarba [7]; finally, if the restriction to only one iteration is further extended to the Picard scheme, then (26) simplifies to

$$(27) \qquad C(\psi^{n-1})(\psi^n - \psi^{n-1}) + T^{n-1}\psi^n = b^{n-1} - \theta(\psi^{n-1}),$$

which is a consistent (nonconservative) discretization of the $\psi$-based form of Richards' equation.

Each iteration in both (25) and (26) requires the solution of a linear system. The linearization approach (26) is apparently simpler because it only performs one Newton step per each Picard iteration. The main problem with the linearized formulation (26) resides in the fact that convergence of the Picard iterations, which is essential to obtain a conservative solution, is much too slow and not always guaranteed.

**5.5. Nonlinearity effect.** The Newton method (25), if successful, yields a perfectly conservative solution at each Picard iterate. Hence, convergence of the (slow) Picard iterations is not essential when the Newton method (25) is used. In specific situations, however, due to the nonlinear character of $\boldsymbol{\theta}(\boldsymbol{\psi})$, the Newton method may not converge unless the initial guess is chosen to be *sufficiently close* to the solution (see, e.g., [2]). This can be explained by the fact that the sign of $dc_i(\psi)/d\psi$ is undefined.

For the proposed method, since $P(\boldsymbol{\psi})$ and $Q(\boldsymbol{\psi})$ are nondecreasing functions of $\boldsymbol{\psi}$, where they are defined, both $dp_i(\psi)/d\psi$ and $dq_i(\psi)/d\psi$ are nonnegative. Intuitively, this explains the monotonicity and the superior robustness of Algorithm 1.

**6. Constitutive relationships.** The nested iterative method described above applies to a large variety of constitutive relationships relating the moisture content and the hydraulic conductivity to the pressure head (see, e.g., [3, 24]). Next, two of the most commonly employed models, namely the Brooks–Corey [3] model and the van Genuchten [24] model, will be illustrated in detail.

**6.1. The Brooks–Corey model.** The constitutive relationships proposed by Brooks and Corey [3] are given by

$$(28) \qquad \theta(\psi) = \begin{cases} \theta_r + (\theta_s - \theta_r)\left(\frac{\psi_d}{\psi}\right)^n & \text{if } \psi \leq \psi_d, \\ \\ \theta_s & \text{if } \psi > \psi_d, \end{cases}$$

$$(29) \qquad K(\psi) = K_s\left[\frac{\theta(\psi) - \theta_r}{\theta_s - \theta_r}\right]^{3+\frac{2}{n}},$$

where $\psi_d = -1/\alpha$, and $\theta_s$, $\theta_r$, $K_s$, $\alpha$, and $n$ are material parameters which affect the shape of the soil hydraulic functions and satisfy $0 \leq \theta_r < \theta_s$ and $K_s$, $\alpha$, $n > 0$.

Assuming that the material parameters $K_s$, $\alpha$, $n$ are constant over each control volume $\Omega_i$, the $i$th water mass can be expressed by (6), where the moisture capacity is derived from (28) which yields

$$(30) \qquad c_i(\psi) = \begin{cases} n_i \frac{\theta_{s,i} - \theta_{r,i}}{|\psi_{d,i}|}\left(\frac{\psi_{d,i}}{\psi}\right)^{n_i+1} & \text{if } \psi \leq \psi_{d,i}, \\ \\ 0 & \text{if } \psi > \psi_{d,i}. \end{cases}$$

Equation (30) indicates that $c_i(\psi)$ assumes its maximum value at $\psi_i^* = \psi_{d,i}$. Moreover, $c_i(\psi)$ is strictly positive and monotonically increasing for all $\psi \leq \psi_i^*$. It has a discontinuity of the first kind at $\psi = \psi_i^*$ and vanishes for all $\psi > \psi_i^*$. Thus, $c_i(\psi)$ is a nonnegative function with bounded variations satisfying both Assumptions C1 and C2. Hence, its Jordan decomposition is given by (7), and the corresponding volumes $\theta_{1,i}(\psi)$ and $\theta_{2,i}(\psi)$ are given by (9).

From (29), the discrete hydraulic conductivity associated to each control volume is taken to be

$$(31) \qquad K_i(\psi) = K_{s,i}\left[\frac{\theta_i(\psi) - \theta_{r,i}}{\theta_{s,i} - \theta_{r,i}}\right]^{3+\frac{2}{n_i}}$$

so that $\mathcal{K}_j > 0$ for all $j = 1, 2, \ldots, M$, and the resulting matrix $T$ is either T1 or T2.

**6.2. The van Genuchten model.** The most commonly employed constitutive relationships were derived by van Genuchten [24] and are given by

$$
(32) \qquad \theta(\psi) = \begin{cases} \theta_r + \dfrac{\theta_s - \theta_r}{[1+|\alpha\psi|^n]^m} & \text{if } \psi \le 0, \\[2ex] \theta_s & \text{if } \psi > 0, \end{cases}
$$

$$
(33) \qquad K(\psi) = K_s \left[ \frac{\theta - \theta_r}{\theta_s - \theta_r} \right]^{\frac{1}{2}} \left\{ 1 - \left[ 1 - \left( \frac{\theta - \theta_r}{\theta_s - \theta_r} \right)^{\frac{1}{m}} \right]^m \right\}^2,
$$

where $m = 1 - 1/n$, and $\theta_s$, $\theta_r$, $K_s$, $\alpha$, and $n \ge 1$ are material parameters which affect the shape of the soil hydraulic functions and satisfy $0 \le \theta_r < \theta_s$ and $K_s$, $\alpha > 0$.

Assuming that the material parameters are constant over each control volume $\Omega_i$, the water mass within each control volume can be expressed by (6), where the $i$th moisture capacity can be derived from (32) which yields

$$
(34) \qquad c_i(\psi) = \begin{cases} \alpha_i n_i m_i \dfrac{\theta_{s,i} - \theta_{r,i}}{[1+|\alpha_i\psi|^{n_i}]^{m_i+1}} |\alpha_i\psi_i|^{n_i-1} & \text{if } \psi \le 0, \\[2ex] 0 & \text{if } \psi > 0. \end{cases}
$$

Further analysis of (34) indicates that $c_i(\psi)$ is nonnegative and assumes its maximum value where $\frac{dc_i(\psi)}{d\psi} = 0$, that is, at

$$
\psi_i^* = -\frac{1}{\alpha_i} \left( \frac{n_i - 1}{n_i} \right)^{\frac{1}{n_i}}.
$$

Moreover, $c_i(\psi)$ is strictly positive and monotonically increasing for all $\psi < \psi_i^*$, monotonically decreasing for all $\psi \in (\psi_i^*, 0)$, and vanishes for all $\psi \ge 0$. Thus, $c_i(\psi)$ is a nonnegative function with bounded variations satisfying both assumptions C1 and C2. Hence, its Jordan decomposition is given by (7), and the corresponding volumes $\theta_{1,i}(\psi)$ and $\theta_{2,i}(\psi)$ are given by (9).

From (33), the discrete hydraulic conductivity associated to each control volume is taken to be

$$
(35) \qquad K_i(\psi) = K_{s,i} \left[ \frac{\theta_i(\psi) - \theta_{r,i}}{\theta_{s,i} - \theta_{r,i}} \right]^{\frac{1}{2}} \left\{ 1 - \left[ 1 - \left( \frac{\theta_i(\psi) - \theta_{r,i}}{\theta_{s,i} - \theta_{r,i}} \right)^{\frac{1}{m}} \right]^m \right\}^2
$$

so that $\mathcal{K}_j > 0$ for all $j = 1, 2, \ldots, M$, and the resulting matrix $T$ is either T1 or T2.

Algorithm 1 then applies to solve system (5) when either the Brooks–Corey or the van Genuchten model is used to specify the constitutive relationships.

**7. Numerical tests.** In order to evaluate the proposed method, Algorithm 1 is applied to nontrivial, one- and two-dimensional test cases available from the literature. The first one-dimensional test case deals with a sharp moisture front that infiltrates into the soil column [10, 16]. The second one-dimensional test case involves flow into a layered soil with variable initial conditions [20, 21]. The two-dimensional test case concerns flow into very dry heterogeneous soil where a perched water table develops

surrounded by unsaturated soil [12, 17, 21]. These cases represent good challenges for a numerical algorithm due to their nonlinear nature.

In each test case, system (5) is solved by Algorithm 1 by using the $L_2$ norm to evaluate the residual and with an accuracy set to $\epsilon = 10^{-3}$, $\epsilon = 10^{-6}$, and $\epsilon = 10^{-12}$, although this latter value may be unnecessarily stringent. Consequently, the resulting local and global mass balance is exact within the specified accuracy.

For each test, the number of time steps, the total number of outer iterations, and the total number of inner iterations required to cover the entire simulation are reported. Also, the average outer iteration per time step and the average inner iteration per outer iteration are indicated.

In particular, the total number of inner iterations corresponds to the number of linear systems being solved within each run. Thus, this number is directly related to the overall performance.

**7.1. Test problem 1.** This problem considers a soil column of $2.0\,m$ deep discretized with a vertical resolution $\Delta z = 0.00625\,m$. The initial pressure head distribution is $\psi(z, 0) = z - 2$. At the bottom of the column, a water table boundary condition (i.e., $\psi(0, t) = 0$) is imposed, while a time-dependent Dirichlet condition is imposed at the top boundary as follows:

$$\psi(2, t) = \begin{cases} -0.05 + 0.03 \sin(2\pi t/100\,000) & \text{if } 0 < t \le 100\,000, \\ 0.1 & \text{if } 100\,000 < t \le 180\,000, \\ -0.05 + 2\,952.45 \exp(-t/18\,204.8) & \text{if } 180\,000 < t \le 300\,000. \end{cases}$$

The soil hydraulic properties are described by the van Genuchten model. The soil parameters are given in Table 1 (see [10, 16]). Thus, having specified Dirichlet boundary conditions, the resulting matrix $T$ is T1 at every time step.

TABLE 1
*Soil hydraulic properties used in test problem 1.*

| | |
|---|---|
| $\theta_s$ | 0.410 |
| $\theta_r$ | 0.095 |
| $\alpha\ (m^{-1})$ | 1.9 |
| $n$ | 1.31 |
| $K_s\ (m/\text{day})$ | 0.062 |

The present simulation is performed using a large time step size $\Delta t = 1\,000\,s$, and only one Picard iteration per time step is allowed. The second period of the simulation ($100\,000 < t \le 180\,000$) is considered to be very challenging for numerical solvers. The sudden increase of the upper Dirichlet boundary condition to a positive value of $0.1\,m$ (ponding) generates a sharp moisture front that infiltrates into the soil column.

At the beginning of the third time period ($t > 180\,000\,s$) ponding decreases exponentially, reaching asymptotically a final value of $-0.05\,m$, and by the end of the simulation the entire column is close to full saturation.

With the specified tolerance $\epsilon = 10^{-3}$, the resulting mass balance error at every time step is always below $10^{-5}$ soon after the very first inner iteration. Consequently, the exit conditions to both the inner and the outer cycle are immediately satisfied and the corresponding numerical results are equivalent to those that can be obtained from the nonconservative scheme (27). In the second run, with a tolerance $\epsilon = 10^{-6}$, convergence is achieved most of the times with just one inner and one outer iteration, except in the second period ($100\,000 < t \le 180\,000$) when two outer iterations are required. Table 2 summarizes the computational statistics for this problem.

TABLE 2
*Results for test problem* 1.

|  | $\epsilon = 10^{-3}$ | $\epsilon = 10^{-6}$ | $\epsilon = 10^{-12}$ |
|---|---|---|---|
| Number of time steps | 300 | 300 | 300 |
| Total outer iterations | 300 | 385 | 1 335 |
| Total inner iterations | 300 | 388 | 2 148 |
| Average outer iterations | 1 | 1.28 | 4.45 |
| Average inner iterations | 1 | 1.01 | 1.61 |

The computed pressure head obtained with a tolerance $\epsilon = 10^{-12}$ is displayed in Figure 1. This solution is very similar to the one reported in the literature [10, 16].



FIG. 1. *Computed pressure head in test problem* 1.

**7.2. Test problem 2.** This case involves vertical drainage through a layered soil from initially saturated conditions. At time $t = 0$, the pressure head at the base of the column is reduced from 2 to $0\,m$. During the subsequent drainage, a no-flow boundary condition is applied to the top of the column. Although a one-dimensional problem, this is considered to be a challenging test for numerical methods because a sharp discontinuity in the moisture content occurs at the interface between two material layers [20, 21].

During downward draining the middle coarse soil tends to restrict drainage from the upper fine soil, and high saturation levels are maintained in the upper fine soil for a considerable period of time. The Brooks–Corey model is used to prescribe the pressure-moisture relationship. The hydraulic properties of the soils are given in Table 3. The soil profile is soil 1 for $0 < z < 0.6\,m$ and $1.2\,m < z < 2\,m$, and soil 2 for

TABLE 3
*Soil hydraulic properties used in test problem* 2.

|  | Soil 1 | Soil 2 |
|---|---|---|
| $\theta_s$ | 0.35 | 0.35 |
| $\theta_r$ | 0.07 | 0.035 |
| $\alpha\ (cm^{-1})$ | 0.0286 | 0.0667 |
| $n$ | 1.5 | 3.0 |
| $K_s\ (cm/s)$ | $9.81 \times 10^{-5}$ | $9.81 \times 10^{-3}$ |

$0.6\,m < z < 1.2\,m$. Thus, having specified a Dirichlet boundary condition at the base of the bottom boundary, the resulting matrix $T$ is T1 at every time step.

Simulations are performed on a mesh of 150 cells and a time step size $\Delta t = 3\,500\,s$. Only one Picard iteration per time step is allowed.

In the first run, with the specified tolerance $\epsilon = 10^{-3}$, the resulting mass balance error at every time step is always below $10^{-3}$ soon after the very first inner iteration. Consequently, the entire simulation is completed with just one inner and one outer iteration per time step. In the second and in the third run, with a tolerance of $10^{-6}$ and $10^{-12}$, respectively, a few inner and outer iterations become necessary to achieve convergence within the required accuracy. The computational statistics for this problem are summarized in Table 4.

Figure 2 shows the saturation computed at a time of $1\,050\,000\,s$ (approximately 12 days) with a tolerance $\epsilon = 10^{-12}$. The discontinuities at the interfaces between the two material layers are resolved within one grid cell.

**7.3. Test problem 3.** The present test problem involves a two-dimensional flow into initially very dry layered soil of sand and clay with a developing water table [12, 17, 21]. The van Genuchten model is used to prescribe the pressure-moisture relationship. The hydraulic properties of the sand and clay are given in Table 5, and the initial pressure head is set to $-480\,m$.

To achieve a perched water table, a $3\,m \times 1\,m$ region of sand is bounded by clay, as shown in Figure 3. A no-flow boundary condition is applied everywhere except for a

TABLE 4
*Results for test problem* 2.

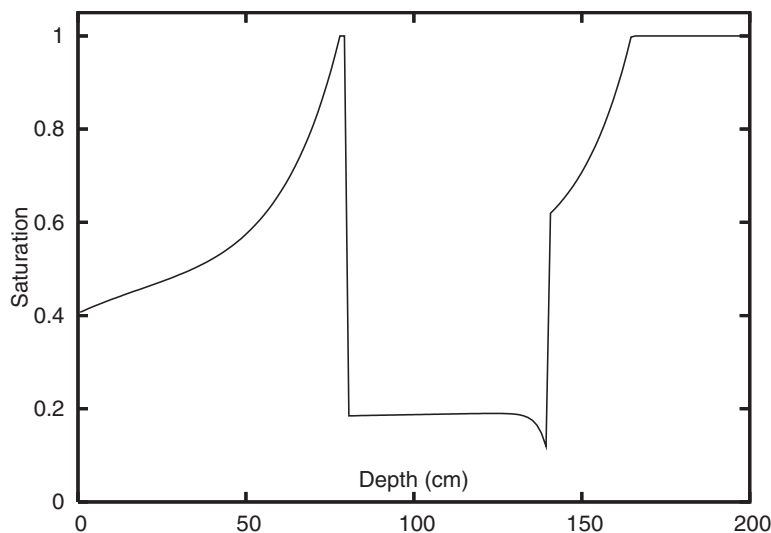|  | $\epsilon = 10^{-3}$ | $\epsilon = 10^{-6}$ | $\epsilon = 10^{-12}$ |
|---|---|---|---|
| Number of time steps | 300 | 300 | 300 |
| Total outer iterations | 300 | 1 260 | 1 443 |
| Total inner iterations | 300 | 1 702 | 4 469 |
| Average outer iterations | 1 | 4.20 | 4.81 |
| Average inner iterations | 1 | 1.35 | 3.10 |



FIG. 2. *Saturation predictions after approximately* 12 *days.*

TABLE 5
*Soil hydraulic properties used in test problem* 3.

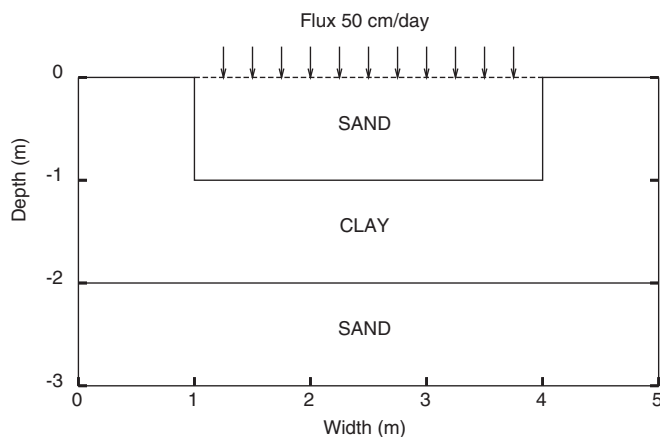|  | Soil 1 | Soil 2 |
|---|---|---|
| $\theta_s$ | 0.3658 | 0.4686 |
| $\theta_r$ | 0.0286 | 0.1060 |
| $\alpha \ (cm^{-1})$ | 0.0280 | 0.0104 |
| $n$ | 2.2390 | 1.3954 |
| $K_s \ (m/s)$ | $6.26 \times 10^{-3}$ | $1.5167 \times 10^{-4}$ |



FIG. 3. *Geometry of perched water table problem.*

water flux rate of $0.5 \, m/\text{day}$ that is applied to the top sand surface. Thus, having spec-ified the Neumann boundary conditions everywhere along the boundary, the resulting matrix $T$ is T2 at every time step. This problem was specifically devised for testing numerical algorithms' ability to survive both very dry conditions and transitions to a saturated state [12].

The flow region is covered by a uniform mesh consisting of a $100 \times 60$ finite volume grid with grid size $\Delta x = \Delta z = 0.05 \, m$. The time step size is set to $\Delta t = 3\,600 \, s$. With such a large $\Delta t$, the model accuracy is enhanced by allowing two Picard iterations per time step through which the diffusive flux matrix is updated. Table 6 summarizes the performance of the proposed algorithm for a simulation period of 24 hours.

After approximately half a day, a water table begins to develop at the interface of the sand and clay layers. Figure 4 shows the computed saturation at $t = 24$ hours that agrees well with published results [12, 17, 21].

TABLE 6
*Results for test problem* 3.

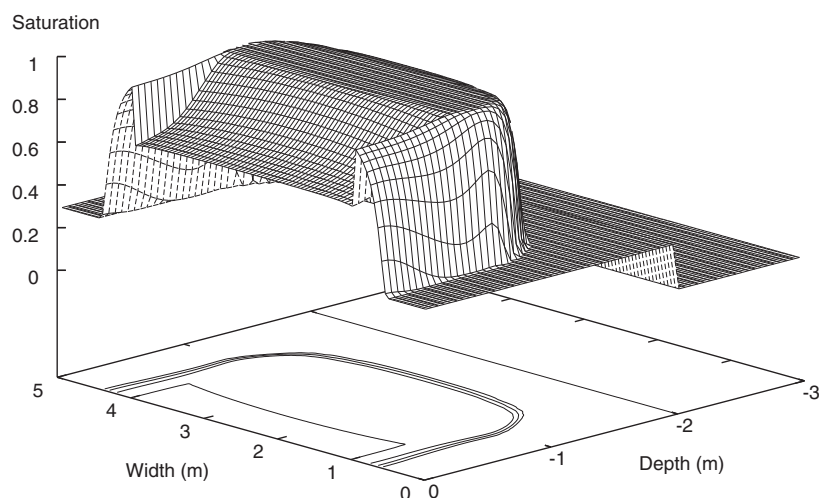|  | $\epsilon = 10^{-3}$ | $\epsilon = 10^{-6}$ | $\epsilon = 10^{-12}$ |
|---|---|---|---|
| Number of time steps | 24 | 24 | 24 |
| Total outer iterations | 58 | 91 | 133 |
| Total inner iterations | 85 | 308 | 482 |
| Average outer iterations | 2.42 | 3.79 | 5.54 |
| Average inner iterations | 1.47 | 3.38 | 3.62 |

Fig. 4. *Computed saturation after* 1 *day.*

**8. Conclusions.** A simple algorithm for solving the nonlinear finite volume discretization of Richards' equation has been derived and analyzed. This method is based on a nested Newton-type linearization and shows similar (quadratic) convergence rate.

It is proved that, under rather general assumptions on the soil properties, the iterates are well defined and monotonically converge to the exact solution. Local and global mass conservation is always assured at each Picard iteration and within a remarkably small number of inner and outer iterations.

Simple, and yet nontrivial, numerical tests have confirmed the efficiency, the robustness, and the usefulness of the proposed algorithm for solving the mixed form of Richards' equation under different flow conditions and for any time step size.

REFERENCES

[1] J. BEAR AND A. VERRUIJT, *Modeling Groundwater Flow and Pollution*, D. Reidel, Dordrecht, Holland, 1987.
[2] L. BERGAMASCHI AND M. PUTTI, *Mixed finite elements and Newton-type linearizations for the solution of Richards' equation*, Int. J. Numer. Methods Engrg., 45 (1999), pp. 1025–1046.
[3] R. H. BROOKS AND A. T. COREY, *Hydraulic properties of porous media*, Hydrology Paper No. 3, Civil Engineering, Colorado State University, Fort Collins, CO, 1964.
[4] L. BRUGNANO AND V. CASULLI, *Iterative solution of piecewise linear systems and applications to flows in porous media*, SIAM J. Sci. Comput., 31 (2009), pp. 1858–1873.
[5] V. CASULLI AND R. A. WALTERS, *An unstructured grid, three-dimensional model based on the shallow water equations*, Internat. J. Numer. Methods Fluids, 32 (2000), pp. 331–348.
[6] V. CASULLI AND P. ZANOLLI, *Semi-implicit numerical modeling of non-hydrostatic free-surface flows for environmental problems*, Math. Comput. Modelling, 36 (2002), pp. 1131–1149.
[7] M. A. CELIA, E. T. BOULOUTAS, AND R. L. ZARBA, *A general mass-conservative numerical solution of the unsaturated flow equation*, Water Resour. Res., 26 (1990), pp. 1483–1496.
[8] V. V. CHISTYAKOV, *On mapping of bounded variation*, J. Dyn. Control Sys., 3 (1997), pp. 261–289.
[9] T. P. CLEMENT, W. R. WISE, AND F. J. MOLZ, *A physically based, two-dimensional, finite-difference algorithm for modeling variably saturated flow*, J. Hydrology, 161 (1994), pp. 71–90.
[10] C. M. F. D'HAESE, M. PUTTI, C. PANICONI, AND N. E. C. VERHOEST, *Assessment of adaptive and heuristic time stepping for variably saturated flow*, Internat. J. Numer. Methods Fluids, 53 (2007), pp. 1173–1193.

[11] C. Fassino and G. Manzini, *Fast-secant algorithms for the non-linear Richards' equation*, Commun. Numer. Methods Engrg., 14 (1998), pp. 921–930.

[12] P. A. Forsyth, Y. S. Wu, and K. Pruess, *Robust numerical methods for saturated-unsaturated flow with dry initial conditions in heterogeneous media*, Adv. Water Resour., 18 (1995), pp. 25-38.

[13] G. H. Golub and C. F. van Loan, *Matrix Computations*, 3rd ed., The Johns Hopkins University Press, Baltimore, MD, 1996.

[14] D. Greenspan and V. Casulli, *Numerical Analysis for Applied Mathematics, Science and Engineering*, Addison Wesley, Redwood City, CA, 1988.

[15] D. Kavetski, P. Binning, and S. W. Sloan, *Adaptive time stepping and error control in a mass conservative numerical solution of the mixed form of Richards equation*, Adv. Water Resour., 24 (2001), pp. 595–605.

[16] D. Kavetski, P. Binning, and S. W. Sloan, *Noniterative time stepping schemes with adaptive truncation error control for the solution of Richards equation*, Water Resour. Res., 38 (2002), pp. 29.1–29.10.

[17] M. R. Kirkland, R. G. Hills, and P. J. Wierenga, *Algorithms for solving Richards equation for variably saturated soils*, Water Resour. Res., 28 (1992), pp. 2049–2058.

[18] K. Krabbenhøft, *An alternative to primary variable switching in saturated–unsaturated flow computations*, Adv. Water Resour., 30 (2007), 483–492.

[19] G. Manzini and S. Ferraris, *Mass-conservative finite volume methods on 2-D unstructured grids for the Richards' equation*, Adv. Water Resour., 27 (2004), pp. 1199–1215.

[20] F. Marinelli and D. S. Durnford, *Semi analytical solution to Richards equation for layered porous media*, J. Irrigation Drainage Engrg., 124 (1998), pp. 290–299.

[21] D. McBride, M. Cross, N. Croft, C. Bennett, and J. Gebhardt, *Computational modelling of variably saturated flow in porous media with complex three-dimensional geometries*, Internat J. Numer. Methods Fluids, 50 (2006), pp. 1085–1117.

[22] C. T. Miller, G. A. Williams, C. T. Kelley, and M. D. Tocci, *Robust solution of Richards' equation for nonuniform porous media*, Water Resour. Res., 34 (1998), pp. 2599–2610.

[23] C. Paniconi and M. Putti, *A comparison of Picard and Newton iteration in the numerical solution of multidimensional variably saturated flow problems*, Water Resour. Res., 30 (1994), pp. 3357–3374.

[24] M. T. van Genuchten, *A Closed-form Equation for Predicting the Hydraulic Conductivity of Unsaturated Soils*, Soil Sci. Soc. Am. J., 44 (1980), pp. 892–898.