




# Mathematical Theory of Nonlinear Single-Phase Poroelasticity

C. J. van Duijn<sup>1,2</sup> · Andro Mikelić<sup>3</sup> 

Received: 28 October 2020 / Accepted: 11 February 2023 / Published online: 10 March 2023  
© The Author(s) 2023

## Abstract

In this paper, we study the equations of nonlinear poroelasticity derived from mixture theory. They describe the quasi-static mechanical behavior of a fluid saturated porous medium. The nonlinearity arises from the compressibility of the fluid and from the dependence of porosity and permeability on the divergence of the displacement. We point some limitations of the model. In our approach, we discretize the quasi-static formulation in time and first consider the corresponding incremental problem. For this, we prove existence of a solution using Brézis' theory of pseudo-monotone operators. Generalizing Biot's free energy to the nonlinear setting, we construct a Lyapunov functional, yielding global stability. This allows us to construct bounds that are uniform with respect to the time step. In the case when dissipative interface effects between the fluid and the solid are taken into account, we consider the continuous time case in the limit when the time step tends to zero. This yields existence of a weak free energy solution.

---

Communicated by Alain Goriely.

---

Andro Mikelić passed away on 28 November 2020.

---

A.M. was partially supported by Darcy Center Eindhoven-Utrecht, the Netherlands, by the project UPGEO ( ANR-19-CU05-032 ) of the French National Research Agency (ANR) and by the LABEX MILYON (ANR-10-LABX-0070) of Université de Lyon, within the program "Investissements d'Avenir" (ANR-11-IDEX-0007) operated by the French National Research Agency (ANR).

---

✉ C. J. van Duijn  
C.J.v.Duijn@TUE.nl

Andro Mikelić  
mikelic@univ-lyon1.fr

- <sup>1</sup> Department of Mechanical Engineering, Darcy Center, Eindhoven University of Technology, Eindhoven, The Netherlands
- <sup>2</sup> Department of Earth Sciences, Utrecht University, Utrecht, The Netherlands
- <sup>3</sup> Univ Lyon, Université Claude Bernard Lyon 1, CNRS UMR 5208, Institut Camille Jordan, 43 blvd. du 11 novembre 1918, F-69622 Villeurbanne cedex, France

**Keywords** Quasi-static nonlinear poroelasticity · Free energy · Incremental problem · Pseudo-monotonicity · Continuous time limit

## 1 Introduction

The elastic quasi-static deformation of a fluid saturated porous medium received much attention in the civil engineering literature because of its relevance to many problems of practical interest. In the framework of consolidation in soil mechanics, these problems relate to the physical loading of soil layers or the effect of soil subsidence due to groundwater withdrawal for drinking water supply or industrial and agricultural purposes. Examples and underlying theories are given in the well-known works of Coussy (2004), Lewis and Schrefler (1998) and Verruijt (2015). They build on the classical theory of Terzaghi (1951) and the pioneering approach of Biot (1962) and Tolstoy (1992).

Recently, other examples of elastic deformation of porous media arise in the context of industrial and biomedical applications, such as paper printing (Bosco et al. 2015), bone regeneration (Cowin 1999; Cardoso et al. 2013), blood flow (Prosi et al. 2005; Čanić et al. 2006) and car filters (Marciniak-Czochra and Mikelić 2015; Mikelić and Tambača 2016).

In its simplest form, assuming both the fluid and the porous material (grains) to be incompressible and assuming the porous medium to be homogeneous and linearly elastic with small strains, the mathematical formulation reads (see Bear and Bachmat 1990; Verruijt 2015 or van Duijn et al. 2019):

$$\operatorname{div} \partial_t \mathbf{u} + \operatorname{div} \left( \frac{\mathbb{K}}{\eta_f} (\rho_f \mathbf{g} - \nabla p) \right) = q \quad (1)$$

and

$$-\operatorname{div} \sigma = \mathbf{F}, \quad (2)$$

where

$$\sigma = \mathcal{G}e(\mathbf{u}) - \alpha p \mathbb{I}, \quad (3)$$

with

$$\mathcal{G}E = 2\mu E + \lambda \operatorname{Tr}(E) \mathbb{I}, \quad \text{for symmetric matrices } E. \quad (4)$$

In these equations,  $\mathbf{u}$  [m] denotes skeleton displacement,  $\mathbb{K}$  [m<sup>2</sup>] intrinsic permeability (a symmetric positive-definite rank-2 tensor),  $\eta_f$  [Pa s] fluid viscosity,  $p$  [Pa] fluid pressure and  $q$  [1/s] sources/sinks. Further,  $\sigma$  [Pa] is the total stress,  $\mathbf{F}$  a given body force (generally linked to gravitational effects),  $\mathcal{G}$  the symmetric, positive-definite, rank-4 Gassmann tensor,  $e(\mathbf{u})$  the linearized strain tensor and  $\alpha \in (0, 1]$  Biot's effective

stress parameter. Finally,  $\mu$  [Pa] and  $\lambda$  [Pa] are Lamé’s parameters. Using for  $\mathcal{G}$  the specific form (4), i.e., Hooke’s law, assumes that the skeleton is mechanically isotropic.

The linear quasi-static Biot system, as well as its dynamical analog, was also derived by means of a multiscale approach, where the starting point is the linear fluid–structure interaction at the pore level. We refer to the monographs Sanchez-Palencia (1980) and Mei and Vernescu (2010) for derivations using two-scale expansions and to Mikelić and Wheeler (2012) for a rigorous mathematical derivation by means of homogenization. The derivations using multiscale analysis confirm Biot’s models in the linear setting. Hence, from different points of view system (1)–(4) is well accepted.

In the engineering literature, one writes  $\alpha = 1 - K/K_g$ , where  $K$  is the drained bulk modulus of the porous skeleton and  $K_g$  the bulk modulus of the grains. Since it is assumed that  $K_g = +\infty$ , we will set  $\alpha = 1$  in (3).

From a mathematical perspective, Eqs. (1)–(4) received much attention. Here, we mention the pioneering paper by Auriault and Sanchez-Palencia (1977) and the work of Ženíšek (1984), who were the first to demonstrate existence and uniqueness. More recent studies include Showalter (2000), Owczarek (2010) and Marciniak-Czochra and Mikelić (2015). Later, Cao et al. (2013) considered a nonlinear extension of (1), by replacing the permeability tensor  $\mathbb{K}$  by the product  $\mathbb{K}k(\text{div } \mathbf{u})$ . The function  $k(\cdot)$  is a relative permeability depending on the volumetric strain  $\text{div } \mathbf{u}$ . From (1), we notice that the overall mixture of two incompressible phases is not incompressible itself.

Though system (1)–(4) is linear, its mathematical complexity lies in the fact that it is of quasi-static nature. In particular (2)–(4) allow to control the size of the volumetric strain only through the size of the data. Some authors circumvent this by introducing a time dependence in (2)–(4) as well. For instance, Bociu et al. (2016) replace  $\mathbf{u}$  in (3) by  $\mathbf{u} + \delta \partial_t \mathbf{u}$ , where  $\delta \geq 0$  is a visco-elastic parameter. Their study allows  $\delta = 0$ , hence it includes the true quasi-static case as well. A different regularization was proposed by Murad and Cushman (1996) who replaced (3) and (4) by

$$\sigma = 2\mu e(\mathbf{u}) + (\lambda \text{div } \mathbf{u} + \lambda^* \text{div } \partial_t \mathbf{u} - \alpha p)\mathbb{I}, \tag{5}$$

with  $\lambda^* > 0$ . This form arises in the non-equilibrium theory, where the fluid pressure and the solid pressure differ by  $\lambda^* \text{div } \partial_t \mathbf{u}$ .

In this paper, we propose to study the quasi-static formulation in which we replace Eq. (1) by the nonlinear fluid phase mass balance based on the mixture theory of Bedford and Drumheller (1978) and Bedford and Drumheller (1983), see, e.g., Rutqvist et al. (2001) and Lewis and Schrefler (1998):

$$n \partial_t \rho + \rho \text{div } \partial_t \mathbf{u} + \text{div } \mathbf{j} = Q, \tag{6}$$

where  $\mathbf{j}$  denotes the Darcy mass flux

$$\mathbf{j} = \frac{\mathbb{K}k(n)\rho}{\eta_f} (\rho \mathbf{g} - \nabla p). \tag{7}$$

Here,  $n$  denotes porosity,  $\rho = \rho_f$  [kg/m<sup>3</sup>] fluid density,  $k$  relative permeability and  $Q$  [kg/m<sup>3</sup> s] sources/sinks.

In Eqs. (6)–(7), the porosity  $n$  is a given function of the volumetric strain: i.e.,

$$n = n(\text{div } \mathbf{u}). \tag{8}$$

An explicit expression for (8) is derived from the Lagrangian solid mass balance equation. This is shown in Sect. 2. Through (8), the relative permeability depends on  $\text{div } \mathbf{u}$ .

Since  $n$  is the volume fraction of voids in the porous medium, it should satisfy the natural bounds

$$0 < n < 1. \tag{9}$$

However, in Sect. 2 we show by means of a counter example that the porosity can attain negative—and thus physically unrealistic—values. Therefore, the bounds in (9) are a major concern in the mathematical model.

To close system (2)–(4), (6)–(7), we introduce a constitutive relation for the fluid density in terms of the pressure. Assuming weak compressibility, we write

$$\rho = \rho(p) = \rho_0(1 + \beta(p - p_0)). \tag{10}$$

Further, we propose an explicit expression for the relative permeability in terms of the porosity

$$k = k(n). \tag{11}$$

In (10),  $\rho_0$  and  $p_0$  are reference values for, respectively, density and pressure and  $\beta$  [ $\text{Pa}^{-1}$ ] is the fluid compressibility coefficient. The relative permeability in (11) satisfies

$$\begin{cases} k \in C^1[0, 1], \\ k(0) > 0 \text{ and } k' > 0 \text{ in } ([0, 1). \end{cases} \tag{12}$$

A well-known example is the Kozeny–Carman formula, see for instance Bear and Bachmat (1990),

$$k(n) = k_0 \frac{n^3}{(1 - n)^2} \quad (k_0 > 0), \tag{13}$$

in a realistic porosity interval, bounded away from  $n = 0$  and  $n = 1$ . Thus, taking  $k$  such that (12) holds and

$$k(n) = \begin{cases} \geq \frac{k_0}{2} \frac{n_*^3}{(1 - n_*)^2}, & \text{for } n \leq n_*, \\ k_0 \frac{n^3}{(1 - n)^2}, & \text{for } 0 < n_* < n < n^* < 1, \\ \leq 2k_0 \frac{(n^*)^3}{(1 - n^*)^2}, & \text{for } n \geq n^*, \end{cases} \tag{14}$$

for appropriately chosen  $0 < n_* < n^* < 1$ , gives a relative permeability satisfying (13) in the interval  $[n_*, n^*]$ .

We notice that Eq. (6), coupled with (2) and (4), is nonlinear due to the relation  $k = k(n)$  and the products involving time derivatives. Assuming constant fluid phase density in the poroelastic mixture is therefore an important simplification. This is studied in Cao et al. (2013) and Bociu et al. (2016).

In studying system (2)–(4), (6)–(11), a crucial role is played by its free energy. The idea is to generalize Biot's free energy (Biot 1962), which is quadratic in strain and fluid density, to the nonlinear poroelastic setting. This free energy serves as a Lyapunov functional. This approach is linked to general entropy methods for PDEs. For a detailed survey covering various fields of applications, we refer to Evans (2004) and to the recent book by Jüngel (2016). An interesting application of the entropy method is discussed in Mikelić (2010), Cao and Pop (2016) and Milišić (2018), where the authors consider dynamic capillary pressure effects in two-phase porous media flow.

This paper is organized as follows. In Sect. 2, we present details of the model formulation. The starting point is the mass balance for the fluid and the solid phase. The latter implies an explicit expression for (8). Introducing a lower bound for the porosity, we modify the fluid mass balance so that a Lyapunov functional can be constructed for the modified system. This modification is such that the fluid equation reduces to its original form in the physical range of the fluid density  $\rho$  and solid volumetric strain  $\mathcal{E}$ . Section 2 is concluded by a weak formulation of the modified system.

In Sect. 3, we consider, for the relaxation parameter  $\lambda^* \geq 0$ , the incremental version of the modified system. Using Brézis' theory of pseudo-monotone operators, existence is demonstrated. Applying the Lyapunov functional yields global (in time) estimates. Next, in Sect. 4, we use these estimates to solve the time-continuous problem when  $\lambda^* > 0$ . In both Sects. 3 and 4, we borrow ideas from Roubiček (2005). Finally, in Sect. 5 we present a discussion and conclusions.

## 2 Problem Formulation

In a number of steps, we construct in this section the equations that serve as starting point for the analysis. The general setting of the problem is as follows:

Let  $\Omega \subset \mathbb{R}^m$  ( $m=2,3$ ) denote a smooth bounded domain, occupied by a linear elastic skeleton. The skeleton material (grains) is assumed incompressible: i.e., the bulk modulus of the grains is infinitely large. The voids in the porous structure are completely filled with a slightly compressible fluid, in the sense that the fluid pressure  $p$  and density  $\rho$  are related by (10).

### 2.1 Balance Equations

For given  $\xi \in \Omega$ , let  $\mathbf{x}(\xi, t)$  denote the location of a solid particle at time  $t > 0$ , that started at  $\mathbf{x}(\xi, 0) = \xi$ . Then the skeleton velocity  $\mathbf{v}_s$  is given by  $\mathbf{v}_s = \partial_t \mathbf{x}|_{\xi}$ .

Restricting themselves to small displacements  $\mathbf{u}$  (within the elastic regime), Rutqvist et al. (2001) and Lewis and Schrefler (1998) argue that in the mass balance equation for the fluid and solid, the material derivative  $\frac{D}{Dt} = \partial_t + \mathbf{v}_s \cdot \nabla$  can be replaced by the partial derivative  $\partial_t$ . This is made explicit by a scaling argument in van Duijn et al. (2019). The resulting Lagrangian form of the mass balances reads:

$$n\partial_t\rho + \rho\operatorname{div}\mathbf{v}_s + \operatorname{div}\mathbf{j} = Q \quad (\text{fluid phase}) \quad (15)$$

and

$$\partial_t(1-n) + (1-n)\operatorname{div}\mathbf{v}_s = 0 \quad (\text{solid phase}), \quad (16)$$

where  $\mathbf{j}$  is mass flux (7).

Within the same approximation, one may write

$$\operatorname{div}\mathbf{v}_s = \partial_t\operatorname{div}\mathbf{u}.$$

Using this in (15) and (16) gives

$$n\partial_t\rho + \rho\partial_t\operatorname{div}\mathbf{u} + \operatorname{div}\mathbf{j} = Q \quad (17)$$

and

$$\partial_t(1-n) + (1-n)\operatorname{div}\partial_t\mathbf{u} = 0. \quad (18)$$

Integrating (18) in time from  $t = 0$ , say, to  $t > 0$ , we have

$$1-n = (1-n_0)e^{-\operatorname{div}(\mathbf{u}-\mathbf{U}_0)} \quad \text{for } t > 0.. \quad (19)$$

Here,  $\mathbf{U}_0$  is the initial displacement and  $n_0$  the initial porosity. With  $n_0 \in (0, 1)$  in  $\Omega$ , expression (19) ensures

$$n < 1 \quad \text{in } \Omega \quad \text{for all } t > 0. \quad (20)$$

To avoid technical complications, we restrict ourselves to  $n_0 = \text{constant}$  in  $\Omega$ .

For small displacements  $\mathbf{u} - \mathbf{U}_0$ , expression (19) is approximated by

$$n = n_0 + (1-n_0)\operatorname{div}(\mathbf{u}-\mathbf{U}_0). \quad (21)$$

**Remark 1** Frequently, the linear form (21) is used for values of  $\operatorname{div}\mathbf{u}$  in a neighborhood of  $\operatorname{div}\mathbf{U}_0$ : i.e., in practical circumstances (21) is applied when  $\mathcal{E}_* < \operatorname{div}(\mathbf{u}-\mathbf{U}_0) < \mathcal{E}^*$ , where  $\mathcal{E}_* < 0 < \mathcal{E}^*$  are appropriately chosen.

Throughout the paper, we redefine

$$\mathbf{u} := \mathbf{u} - \mathbf{U}_0, \tag{22}$$

where  $\mathbf{U}_0 \in H_0^1(\Omega)^m \cap H^2(\Omega)^m$  is the initial displacement. Redefining accordingly

$$\mathbf{F} := \mathbf{F} + \operatorname{div} (\mathcal{G}e(\mathbf{U}_0)), \tag{23}$$

we obtain for the fluid pressure  $p$  and the skeleton displacement  $\mathbf{u}$  the system:

$$n \partial_t \rho + \rho \operatorname{div} \partial_t \mathbf{u} + \operatorname{div} \left( \frac{\mathbb{K}k(n)\rho}{\eta_f} (\rho \mathbf{g} - \nabla p) \right) = \mathcal{Q}, \tag{24}$$

$$- \operatorname{div} (\mathcal{G}e(\mathbf{u}) - p\mathbb{I}) = \mathbf{F}, \tag{25}$$

where

$$\rho = \rho(p) = \rho_0(1 + \beta(p - p_0)), \tag{26}$$

$$n = n(\operatorname{div} \mathbf{u}) = 1 - (1 - n_0)e^{-\operatorname{div} \mathbf{u}} \tag{27}$$

$$\approx n_0 + (1 - n_0)\operatorname{div} \mathbf{u} \quad (\text{small strains}). \tag{28}$$

**Remark 2** Concerning the initial displacement  $\mathbf{U}_0$ , we note that only  $\operatorname{div} \mathbf{U}_0$ , the initial volumetric strain, is used. However, when discussing the free energy, one needs in addition that  $\mathbf{U}_0$  is such that the corresponding elastic energy is finite. For simplicity, we suppose  $\mathbf{U}_0 \in H^2(\Omega)^m$ .

In the next sections, we will develop the mathematical theory for system (24)–(28).

The issue of negative porosity in (27) (or, for that matter, a porosity exceeding one in approximation (28)), is discussed next.

### 2.2 Negative Porosity

We consider a simplified version of the linear problem (1)–(4) and show that  $\operatorname{div} \mathbf{u}$  can attain values for which the porosity from (27)–(28) becomes negative.

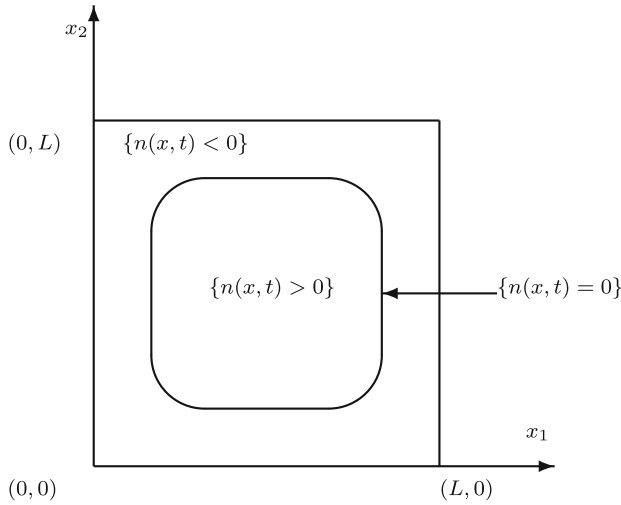
For simplicity, we give the construction in  $\mathbb{R}^2$ .

Let  $\Omega = (0, L)^2$  for some  $L > 0$ . We suppose, as in the rest of this paper, that  $\operatorname{div} \mathbf{u}|_{t=0} = 0$ . Further we set  $\mathbf{F} = 0$  in (25). Using (4) in (25) gives

$$- \operatorname{div} (2\mu e(\mathbf{u}) + (\lambda \operatorname{div} \mathbf{u} - p)\mathbb{I}) = 0 \quad \text{in } \Omega, \tag{29}$$

Proceeding as in Verruijt (2015), when he discusses the Mandel problem, we take the divergence of (29) to obtain

$$\Delta((2\mu + \lambda)\operatorname{div} \mathbf{u} - p) = 0 \quad \text{in } \Omega. \tag{30}$$



**Fig. 1** Sketch of level set of porosity  $n$  at some  $t > 0$ . The region where  $n > 0$  shrinks with increasing time and disappears after a finite time  $T_p > 0$

Hence the function

$$H = (2\mu + \lambda)\text{div } \mathbf{u} - p$$

is harmonic in  $\Omega$ .

The idea is to prescribe boundary conditions for Eqs. (24) and (25) so that  $H|_{\partial\Omega}$  is given. For instance, if we set along the four edges, see Fig. 1,

$$\begin{cases} \{x_1 = 0\} : u_2 = 0, \sigma_{11} = \Sigma^{1,0} \text{ and } p = 0; \\ \{x_1 = L\} : u_2 = 0, \sigma_{11} = \Sigma^{1,L} \text{ and } p = 0; \\ \{x_2 = 0\} : u_1 = 0, \sigma_{22} = \Sigma^{2,0} \text{ and } p = 0; \\ \{x_2 = L\} : u_1 = 0, \sigma_{22} = \Sigma^{2,L} \text{ and } p = 0, \end{cases} \tag{31}$$

and use

$$\sigma_{11} = 2\mu \frac{\partial u_1}{\partial x_1} + \lambda \text{div } \mathbf{u} - p,$$

we have

$$\Sigma^{1,0} = (2\mu + \lambda) \frac{\partial u_1}{\partial x_1} \text{ at } \{x_1 = 0\},$$

implying

$$H = \Sigma^{1,0} \text{ at } \{x_1 = 0\}.$$



Repeating this along the other edges gives

$$H|_{\partial\Omega} = \Sigma^b,$$

where  $\Sigma^b$  denotes the given value of  $\sigma$  along the edges.

Then, we have

**Proposition 1** *Let  $\mathcal{E} = \text{div} \mathbf{u}$  denote the volumetric stress and let  $n(\mathcal{E})$  be given by (27). Suppose there exists a constant  $\Sigma > 0$  such that  $\Sigma^b \leq -\Sigma$ . Then, for  $\Sigma$  sufficiently large, there exists a  $T_p = T_p(\Sigma) > 0$  such that*

$$n(\mathcal{E}(x, t)) < 0 \text{ for } t > T_p \text{ and } x \in \overline{\Omega}. \tag{32}$$

**Proof** Note that the sign of  $\Sigma^b$  implies compression of the medium. Restricting ourselves to the linear case (1) in a homogeneous and isotropic porous medium in which sources/sinks and gravity are absent, we have

$$\partial_t \text{div} \mathbf{u} - \frac{K}{\eta_f} \Delta p = 0 \text{ in } \Omega, t > 0. \tag{33}$$

Since

$$\Delta p = (2\mu + \lambda) \Delta(\text{div} \mathbf{u})$$

and

$$\text{div} \mathbf{u}|_{\partial\Omega} = \frac{\Sigma_b}{2\mu + \lambda} \leq -\frac{\Sigma}{2\mu + \lambda},$$

we have for  $\mathcal{E} = \text{div} \mathbf{u}$  the problem

$$\begin{cases} \partial_t \mathcal{E} = \frac{(2\mu + \lambda)K}{\eta_f} \Delta \mathcal{E} & \text{in } \Omega, t > 0; \\ \mathcal{E}|_{\partial\Omega} \leq -\frac{\Sigma}{2\mu + \lambda} & \text{for } t > 0; \\ \mathcal{E}|_{\{t=0\}} = 0 & \text{in } \Omega. \end{cases} \tag{34}$$

By the strong maximum principle,  $\mathcal{E} < \bar{\mathcal{E}}$  in  $\Omega$  and for  $t > 0$ , where  $\bar{\mathcal{E}}$  is the solution of problem (34) with  $\bar{\mathcal{E}} = -\Sigma/(2\mu + \lambda)$  on  $\partial\Omega$ . Writing  $\bar{\mathcal{E}}$  as a Fourier series, one observes that

$$\bar{\mathcal{E}}(x, t) \rightarrow -\frac{\Sigma}{2\mu + \lambda} \text{ as } t \rightarrow +\infty,$$

uniformly in  $x \in \overline{\Omega}$ .

Thus, if

$$(1 - n_0)e^{\Sigma/(2\mu + \lambda)} > 1,$$

or

$$\Sigma > (2\mu + \lambda) \ln \frac{1}{1 - n_0},$$

the result is immediate. □

This example shows that there is a problem with the model. A modification is needed to prevent the porosity (27), or (28), to become negative. Of course, one could argue that this is outside the scope of the model or outside the range of practical applications, since linear elasticity and small strains are supposed. However, since it is not clear how to ensure that indeed small displacements/strains are guaranteed, one needs to impose a porosity modification to prevent negative values.

### 2.3 Modification of Balance Equations

In a number of steps, we modify Eq. (24) so that it becomes well-posed in a mathematical sense and reduces to its original form in the physical range of the unknowns.

First, to satisfy the natural bounds (9), we replace the porosity approximation (28) by a smooth increasing function  $\bar{n} : \mathbb{R} \rightarrow \mathbb{R}$  such that

$$\bar{n}(\mathcal{E}) = \begin{cases} \lim_{\mathcal{E} \rightarrow -\infty} \bar{n}(\mathcal{E}) = \delta_0 > 0, \\ n_0 + (1 - n_0)\mathcal{E}, \text{ for } \mathcal{E}_* \leq \mathcal{E} \leq \mathcal{E}^*; \\ \lim_{\mathcal{E} \rightarrow +\infty} \bar{n}(\mathcal{E}) = 1. \end{cases} \tag{35}$$

Here,  $\mathcal{E}_*$  and  $\mathcal{E}^*$  are practical values chosen such that  $-n_0/(1 - n_0) < \mathcal{E}_* < 0 < \mathcal{E}^* < 1$  and  $\delta_0 = n(\mathcal{E}_*)/2$ , see Fig. 2 for a sketch.

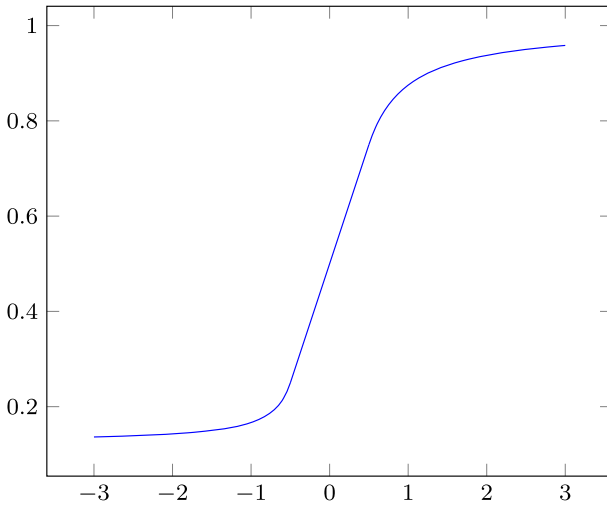
This construction ensures that the modified porosity  $\bar{n}(\mathcal{E})$  remains in the physical range  $(0, 1)$  and coincides with the linear approximation in the interval  $(\mathcal{E}_*, \mathcal{E}^*)$ . Realistic porosity measurements are always done away from the bounds  $n = 0$  and  $n = 1$ , see, e.g., Bear and Bachmat (1990).

We choose to study Eq. (24) with the fluid density as primary unknown. Hence, we need to express the pressure  $p$  in terms of  $\rho$ . Using (26), we have explicitly

$$p = p(\rho) := p_0 + \frac{\rho - \rho_0}{\beta\rho_0}. \tag{36}$$

When considering (24), one clearly has in mind that  $\rho$  takes values near the reference  $\rho_0$ . However, the mathematical nature of the equations does not guarantee this behavior. Hence, a second modification is needed, now for  $\rho$  in the second and third term of the left-hand side of (24). Disregarding gravity, we replace (24) by the modified fluid mass balance equation

$$\bar{n}(\mathcal{E})\partial_t \rho + d(\rho)\partial_t \mathcal{E} - \operatorname{div} \left( k(\mathcal{E})\mathcal{D}(\rho)\mathbb{K}\nabla \rho \right) = Q, \tag{37}$$



**Fig. 2** Sketch of porosity cutoff  $\bar{n}(\mathcal{E})$ , with  $n_0 = 0.5$ ,  $\mathcal{E}_* = -0.5$ ,  $\mathcal{E}^* = 0.5$  and  $\delta_0 = 0.125$

where  $\bar{n}(\mathcal{E})$  is given by (35) and  $k(\mathcal{E}) = k(\bar{n}(\mathcal{E}))$ . Further,  $d, \mathcal{D} : \mathbb{R} \rightarrow \mathbb{R}$  are chosen such that

$$\left. \begin{aligned} d(\rho) &= \rho, \\ \mathcal{D}(\rho) &= \frac{\rho}{\eta_f \beta \rho_0}, \end{aligned} \right\} \text{ for } |\rho - \rho_0| \leq \rho_0 - \rho_*, \tag{38}$$

where  $\rho_* \in (0, \rho_0)$  is a small constant. Outside this range we take for  $d$  and  $\mathcal{D}$  extensions that suit the mathematical analysis. We clarify this at a later point in this section.

**Remark 3** The composite function  $k(\mathcal{E}) = k(\bar{n}(\mathcal{E}))$  satisfies:  $k \in C^1(\mathbb{R}) \cap L^\infty(\mathbb{R})$  and  $k > k(\delta_0) > 0, k' > 0$  in  $\mathbb{R}$ .

The balance of forces (25) is modified by adding the regularizing term  $\lambda^* \partial_t \mathcal{E}$ , as in expression (5). This gives

$$-\operatorname{div} \left( \mathcal{G}e(\mathbf{u}) + (\lambda^* \partial_t \mathcal{E} - p)\mathbb{I} \right) = \mathbf{F}, \tag{39}$$

where  $\mathcal{E} = \operatorname{div} \mathbf{u}$  and where  $\lambda^* \geq 0$ .

We consider system (37), (39) in the set

$$Q_T = \{(x, t) : x \in \Omega, 0 < t < T\},$$

where  $T > 0$  is arbitrarily chosen. To avoid technical complications, we take  $\partial\Omega \in C^1$  throughout the rest of this paper.

As initial conditions, we have

$$\mathcal{E}|_{t=0} = 0 \quad \text{and} \quad \rho|_{t=0} = \rho^0 \quad \text{in } \Omega, \tag{40}$$

where  $\rho^0 : \Omega \rightarrow (0, +\infty)$  is taken near the reference value  $\rho_0$ . Along the boundary, we prescribe

$$\mathbf{u}|_{\partial\Omega} = 0, \quad \nabla \rho \cdot \boldsymbol{\nu}|_{\partial\Omega} = 0, \quad \text{for } 0 < t \leq T. \tag{41}$$

where  $\boldsymbol{\nu}$  is the outward unit normal at  $\partial\Omega$ .

### 2.4 Lyapunov Functional

In this section, we derive an expression for the free energy which acts as a Lyapunov functional for system (37), (39). This generalizes the free energy introduced originally by Biot (1962).

Let  $\{\mathbf{u}, \rho\}$  be a smooth solution of Eqs. (37), (39) that satisfies conditions (40) and (41). Further, let  $g : \mathbb{R} \rightarrow \mathbb{R}$  be a smooth, strictly increasing and globally Lipschitz function satisfying  $g(\rho_0) = 0$ .

We first multiply equation (39) by  $\partial_t \mathbf{u}$  and integrate the result in  $\Omega$ . This gives

$$\begin{aligned} & \frac{1}{2} \frac{d}{dt} \int_{\Omega} \mathcal{G}e(\mathbf{u}) : e(\mathbf{u}) \, dx + \lambda^* \int_{\Omega} (\partial_t \mathcal{E})^2 \, dx - \frac{d}{dt} \int_{\Omega} \mathbf{F} \cdot \mathbf{u} \, dx \\ & - \int_{\Omega} p(\rho) \partial_t \mathcal{E} \, dx = - \int_{\Omega} \partial_t \mathbf{F} \cdot \mathbf{u} \, dx. \end{aligned} \tag{42}$$

Next, we multiply (37) by  $g(\rho)$  and integrate the result in  $\Omega$ . This results in

$$\begin{aligned} & \int_{\Omega} \bar{n}(\mathcal{E}) g(\rho) \partial_t \rho \, dx + \int_{\Omega} d(\rho) g(\rho) \partial_t \mathcal{E} \, dx - \int_{\Omega} k(\mathcal{E}) \mathcal{D}(\rho) g'(\rho) \mathbb{K} \nabla \rho \cdot \nabla \rho \, dx \\ & = \int_{\Omega} Qg(\rho) \, dx. \end{aligned} \tag{43}$$

With

$$G(\rho) = \int_{\rho_0}^{\rho} g(z) \, dz, \tag{44}$$

the first term in (43) can be written as

$$\int_{\Omega} \bar{n}(\mathcal{E}) \partial_t G(\rho) \, dx = \partial_t \int_{\Omega} \bar{n}(\mathcal{E}) G(\rho) \, dx - \int_{\Omega} \bar{n}'(\mathcal{E}) G(\rho) \partial_t \mathcal{E} \, dx. \tag{45}$$

Note that  $G$  is a nonnegative, convex function with  $G(\rho_0) = 0$ .

We substitute (45) back into (43). Adding the resulting expression and (42) yields

$$\begin{aligned} & \frac{d}{dt} \int_{\Omega} \left( \frac{1}{2} \mathcal{G}e(\mathbf{u}) : e(\mathbf{u}) + \bar{n}(\mathcal{E}) G(\rho) - \mathbf{F} \cdot \mathbf{u} \right) \, dx + \lambda^* \int_{\Omega} (\partial_t \mathcal{E})^2 \, dx \\ & + \int_{\Omega} k(\mathcal{E}) \mathcal{D}(\rho) g'(\rho) \mathbb{K} \nabla \rho \cdot \nabla \rho \, dx \end{aligned}$$

$$+ \int_{\Omega} \left\{ d(\rho)g(\rho) - \bar{n}'(\mathcal{E})G(\rho) - p(\rho) \right\} \partial_t \mathcal{E} \, dx = \int_{\Omega} Qg(\rho) \, dx - \int_{\Omega} \partial_t \mathbf{F} \cdot \mathbf{u} \, dx. \tag{46}$$

Before considering the general nonlinear case described by this expression, we first show its implication for the simplified linear setting where we have

$$\bar{n}(\mathcal{E}) = n_0, \quad d(\rho) = \rho_0, \quad k(\mathcal{E}) = 1 \quad \text{and} \quad \mathcal{D} = \frac{1}{\eta_f \beta}.$$

Then,

$$\int_{\Omega} \left\{ d(\rho)g(\rho) - \bar{n}'(\mathcal{E})G(\rho) - p(\rho) \right\} \partial_t \mathcal{E} \, dx \tag{47}$$

in expression (46) simplifies to

$$\int_{\Omega} \left\{ \rho_0 g(\rho) - p(\rho) \right\} \partial_t \mathcal{E} \, dx. \tag{48}$$

Since

$$\int_{\Omega} \partial_t \mathcal{E} \, dx = 0,$$

expression (48) vanishes if  $g(\rho)$  is chosen such that

$$\rho_0 g(\rho) - p(\rho) = \text{constant} = -p_0.$$

This gives

$$g(\rho) = \frac{\rho - \rho_0}{\beta \rho_0^2}$$

and

$$G(\rho) = \frac{(\rho - \rho_0)^2}{2\beta \rho_0^2}.$$

Using these expressions in (46) yields

$$\begin{aligned} & \frac{d}{dt} \int_{\Omega} \left\{ \frac{1}{2} \mathcal{G}e(\mathbf{u}) : e(\mathbf{u}) + \frac{n_0}{\beta \rho_0^2} (\rho - \rho_0)^2 - \mathbf{F} \cdot \mathbf{u} \right\} \, dx + \lambda^* \int_{\Omega} (\partial_t \mathcal{E})^2 \, dx \\ & + \int_{\Omega} \frac{1}{\eta_f \beta^2 \rho_0^2} \mathbb{K} \nabla \rho \cdot \nabla \rho \, dx = \int_{\Omega} Qg(\rho) \, dx - \int_{\Omega} \partial_t \mathbf{F} \cdot \mathbf{u} \, dx. \end{aligned} \tag{49}$$

Hence,

$$\mathcal{L}(\mathbf{u}, \rho) = \int_{\Omega} \left( \frac{1}{2} \mathcal{G}e(\mathbf{u}) : e(\mathbf{u}) + \frac{n_0}{2\beta\rho_0^2} (\rho - \rho_0)^2 - \mathbf{F} \cdot \mathbf{u} \right) dx \tag{50}$$

acts as a Lyapunov functional for the linear form of system (37), (39). The first term denotes the elastic energy of the skeleton, the second term the compression energy of the fluid, and the third term the work done by the force  $\mathbf{F}$ .

Expression (50) coincides with Biot’s original free energy expression from Biot (1962).

Next, we return to the nonlinear case (46). As a first step, we restrict ourselves to the physical range of the porosity. Then, integral (47) becomes

$$\int_{\Omega} \left\{ d(\rho)g(\rho) - (1 - n_0)G(\rho) - p(\rho) \right\} \partial_t \mathcal{E} \, dx. \tag{51}$$

This integral vanishes if  $g(\rho)$  is chosen such that

$$d(\rho)g(\rho) - (1 - n_0)G(\rho) - p(\rho) = -p_0 \tag{52}$$

Differentiating the expression yields a first-order equation for  $g$ . Thus for (51) to vanish,  $g$  should satisfy the initial value problem

$$\begin{cases} d(\rho)g'(\rho) + (d'(\rho) - (1 - n_0))g = \frac{1}{\rho_0\beta}, \text{ for } \rho \in \mathbb{R}; \\ g(\rho_0) = 0. \end{cases} \tag{53}$$

We first consider this problem in the interval  $|\rho - \rho_0| < \bar{\rho} := \rho_0 - \rho_*$  where  $d(\rho) = \rho$ . Then, (53) reduces to

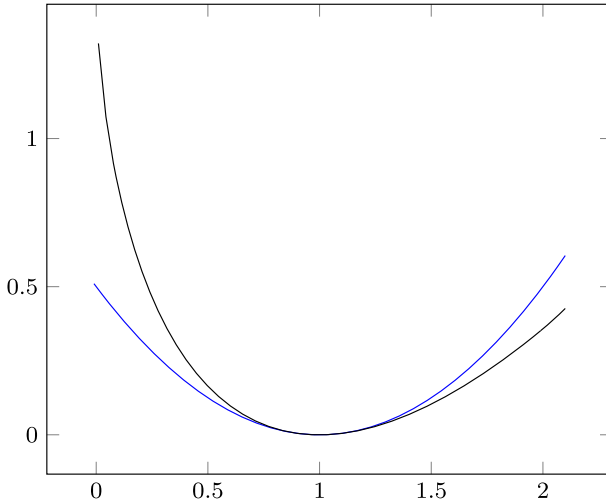
$$\begin{cases} \rho g' + n_0 g = \frac{1}{\beta\rho_0}, \\ g(\rho_0) = 0. \end{cases} \tag{54}$$

Direct integration results in

$$g(\rho) = \frac{1}{\beta n_0 \rho_0} \left( 1 - \left( \frac{\rho_0}{\rho} \right)^{n_0} \right). \tag{55}$$

A second integration yields (Fig. 3)

$$G(\rho) = \int_{\rho_0}^{\rho} g(\xi) \, d\xi = \frac{1}{\beta n_0 (1 - n_0) \rho_0} \left( (1 - n_0)\rho - \rho_0^{n_0} \rho^{1-n_0} + n_0 \rho_0 \right). \tag{56}$$



**Fig. 3** Sketch of the free energy  $\beta G(\rho/\rho_0)$ . The linear case is in blue. The nonlinear case, see (56) and (59) with  $n_0 = 1/3$  and  $\rho_*/\rho_0 = 0.01$ , is in black

When  $|\rho - \rho_0| > \bar{\rho}$ , the function  $d(\rho)$  has not yet been defined. We do this by first extending  $g(\rho)$  for  $|\rho - \rho_0| > \bar{\rho}$  and then by solving  $d(\rho)$  from (52): i.e.,

$$d(\rho) = \frac{(1 - n_0)G(\rho) + p(\rho) - p_0}{g(\rho)}. \tag{57}$$

Clearly, (55) cannot be used for  $\rho \leq 0$ . Instead, we extend (55) in a linear  $C^1$ -manner for  $|\rho - \rho_0| > \bar{\rho}$ . With  $\tilde{\rho} = \rho_0 + \bar{\rho} = 2\rho_0 - \rho_*$ , we set

$$g(\rho) = \begin{cases} \frac{1}{\beta n_0 \rho_0} \left\{ 1 - \left(\frac{\rho_0}{\rho_*}\right)^{n_0} + \frac{\rho - \rho_*}{\rho_*} \left(\frac{\rho_0}{\rho_*}\right)^{n_0} \right\} & \text{for } \rho < \rho_*, \\ \frac{1}{\beta n_0 \rho_0} \left\{ 1 - \left(\frac{\rho_0}{\tilde{\rho}}\right)^{n_0} + \frac{\rho - \tilde{\rho}}{\tilde{\rho}} \left(\frac{\rho_0}{\tilde{\rho}}\right)^{n_0} \right\} & \text{for } \rho > \tilde{\rho}, \end{cases} \tag{58}$$

yielding

$$G(\rho) = \begin{cases} G(\rho_*) + \frac{\rho - \rho_*}{\beta n_0 \rho_0} \left( 1 - \left(\frac{\rho_0}{\rho_*}\right)^{n_0} + \frac{\rho - \rho_*}{2\rho_*} \left(\frac{\rho_0}{\rho_*}\right)^{n_0} \right) & \text{for } \rho < \rho_*; \\ G(\tilde{\rho}) + \frac{\rho - \tilde{\rho}}{\beta n_0 \rho_0} \left( 1 - \left(\frac{\rho_0}{\tilde{\rho}}\right)^{n_0} + \frac{\rho - \tilde{\rho}}{2\tilde{\rho}} \left(\frac{\rho_0}{\tilde{\rho}}\right)^{n_0} \right) & \text{for } \rho > \tilde{\rho}. \end{cases} \tag{59}$$

Substituting expressions (58) and (59) in (57) yields the desired extension for  $d(\rho)$  when  $|\rho - \rho_0| > \bar{\rho}$ . Thus,

$$d(\rho) = \begin{cases} \rho & \text{for } |\rho - \rho_0| \leq \bar{\rho}, \\ \text{(57) with } g \text{ and } G \text{ given by (58) and (59)} & \text{for } |\rho - \rho_0| > \bar{\rho}. \end{cases} \tag{60}$$

Hence, the triple  $\{g(\rho), G(\rho), d(\rho)\}$  constructed above satisfies (52). For this choice, the integral (51) drops from expression (46). So far, we considered for the porosity the linear approximation  $n(\mathcal{E}) = n_0 + (1 - n_0)\mathcal{E}$ . To deal with the full cutoff (35), we introduce a second modification. The starting point is (47). This integral vanishes if

$$d(\rho)g(\rho) - \bar{n}'(\mathcal{E})G(\rho) = p(\rho) - p_0. \tag{61}$$

Keeping  $g$  as in (55), (58) and  $G$  as in (56), (59), we now modify  $d(\rho)$ , calling it  $D(\rho, \mathcal{E})$ , such that

$$D(\rho, \mathcal{E}) = \frac{\bar{n}'(\mathcal{E})}{g(\rho)}G(\rho) + \frac{p(\rho) - p_0}{g(\rho)}. \tag{62}$$

Using (57) in this expression gives

$$D(\rho, \mathcal{E}) = d(\rho) + (\bar{n}'(\mathcal{E}) - (1 - n_0))\frac{G(\rho)}{g(\rho)}. \tag{63}$$

Clearly, for  $|\rho - \rho_0| < \bar{\rho}$  and  $\mathcal{E}_* < \mathcal{E} < \mathcal{E}^*$ , this expression reduces to

$$D(\rho, \mathcal{E}) = \rho.$$

Finally, we use in the Darcy mass flux term  $\mathbf{j}$  from Eq. (37)

$$\mathcal{D}(\rho) = \frac{1}{\eta_f \rho_0 \beta} \begin{cases} \tilde{\rho}, & \text{for } \rho \geq \tilde{\rho}; \\ \rho, & \text{for } \rho_* < \rho < \tilde{\rho}; \\ \rho_*, & \text{for } \rho \leq \rho_*. \end{cases} \tag{64}$$

Thus, in the end we consider the “second” modified fluid mass balance equation

$$\bar{n}(\text{div } \mathbf{u})\partial_t \rho + D(\rho, \text{div } \mathbf{u}) \text{div } \partial_t \mathbf{u} = \text{div} \left( k(\bar{n}(\text{div } \mathbf{u}))\mathcal{D}(\rho)\mathbb{K}\nabla \rho \right) + \mathcal{Q}. \tag{65}$$

System (39), (65) serves as starting point of the analysis. The function  $D(\rho, \mathcal{E})$  in (65) generalizes the fluid density. It is chosen so that

$$J(\mathbf{u}, \rho) = \frac{1}{2} \int_{\Omega} \mathcal{G}e(\mathbf{u}) : e(\mathbf{u}) \, dx + \int_{\Omega} \bar{n}(\text{div } \mathbf{u})G(\rho) \, dx - \int_{\Omega} \mathbf{F} \cdot \mathbf{u} \, dx \tag{66}$$

acts as a Lyapunov functional for the system. The function  $G : \mathbb{R} \rightarrow \mathbb{R}$  satisfies  $G(\rho_0) = 0$ ,  $G(\rho) > 0$  if  $\rho \neq \rho_0$  and  $G$  is strictly convex, with quadratic behavior for large values of  $|\rho|$ . It is explicitly given by (56) and (59).



### 2.5 Summary of Equations and Weak Formulation

The problem describing the nonlinear poroelastic behavior of a fluid saturated porous medium is to find the displacement  $\mathbf{u} : \overline{Q}_T \rightarrow \mathbb{R}^m$  and the fluid density  $\rho : \overline{Q}_T \rightarrow \mathbb{R}$  satisfying

(i) the balance equations

$$\bar{n}(\mathcal{E})\partial_t \rho + D(\rho, \mathcal{E})\partial_t \mathcal{E} = \operatorname{div} \left( k(\mathcal{E})\mathcal{D}(\rho)\mathbb{K}\nabla \rho \right) + Q, \tag{67}$$

$$- \operatorname{div} \left( \mathcal{G}e(\mathbf{u}) + \lambda^* \partial_t \mathcal{E} \mathbb{I} - p(\rho)\mathbb{I} \right) = \mathbf{F}, \tag{68}$$

in  $Q_T = (0, T) \times \Omega$  and

(ii) the initial-boundary conditions (40)–(41).

The coefficients in Eqs. (67)–(68) were introduced in this section. Specifically,

$\bar{n}(\mathcal{E})$  and  $k(\mathcal{E})$  satisfy (35) and Remark 3,

$D(\rho, \mathcal{E})$ ,  $\mathcal{D}(\rho)$  and  $p(\rho)$  are given by (62), (64) and (36),

and  $\lambda^* \geq 0$ .

We recast this classical formulation in the following weak form.

**Definition 1** We call a triple  $(\mathbf{u}, \mathcal{E}, \rho) \in L^\infty(0, T; H^1(\Omega)^m) \times L^\infty(0, T; H^1_{loc}(\Omega)) \times (L^2(0, T; H^1(\Omega)) \cap L^\infty(0, T; L^2(\Omega)))$ ,  $\partial_t \mathcal{E} \in L^2(Q_T) \cap L^\infty(0, T; H^1_{loc}(\Omega))$  a **weak free energy** solution if (i)

$$\begin{aligned} & - \int_0^T \int_\Omega \rho \bar{n}(\mathcal{E}) \partial_t \Phi \, dx dt - \int_\Omega n_0 \rho^0(x) \Phi(x, 0) \, dx \\ & + \int_0^T \int_\Omega \partial_t \mathcal{E} \left( D(\rho, \mathcal{E}) - \rho \bar{n}'(\mathcal{E}) \right) \Phi \, dx dt \\ & + \int_0^T \int_\Omega k(\mathcal{E}) \mathcal{D}(\rho) \mathbb{K} \nabla \rho \cdot \nabla \Phi \, dx dt \\ & = \int_0^T \int_\Omega Q \Phi \, dx dt, \quad \forall \Phi \in H^1(Q_T), \quad \Phi|_{t=T} = 0; \end{aligned} \tag{69}$$

(ii)

$$\mathcal{E} = \operatorname{div} \mathbf{u};$$

(iii)

$$\begin{aligned} & \int_\Omega \mathcal{G}e(\mathbf{u}) : e(\xi) \, dx + \lambda^* \partial_t \int_\Omega \mathcal{E} \operatorname{div} \xi \, dx - \int_\Omega p(\rho) \operatorname{div} \xi \, dx \\ & = \int_\Omega \mathbf{F} \cdot \xi \, dx, \quad \forall \xi \in H_0^1(\Omega)^3 \text{ and for almost all } t \in (0, T]; \end{aligned} \tag{70}$$

(iv)

$$\mathcal{E}|_{t=0} = 0 \quad \text{in } \Omega. \tag{71}$$

(v) For every  $t_1, t_2 \in [0, T]$ ,  $t_1 < t_2$ ,

$$\begin{aligned} & \int_{\Omega} \left( \frac{1}{2} \mathcal{G}e(\mathbf{u}(t_2)) : e(\mathbf{u}(t_2)) + \bar{n}(\mathcal{E}(t_2))G(\rho(t_2)) - \mathbf{F}(t_2) \cdot \mathbf{u}(t_2) \right) dx \\ & + \int_{t_1}^{t_2} \int_{\Omega} \left( \lambda^*(\partial_t \mathcal{E})^2 + k(\mathcal{E})\mathcal{D}(\rho)g'(\rho)\mathbb{K}\nabla\rho \cdot \nabla\rho - Qg(\rho) + \partial_t \mathbf{F} \cdot \mathbf{u} \right) dx dt \\ & \leq \int_{\Omega} \left( \frac{1}{2} \mathcal{G}e(\mathbf{u}(t_1)) : e(\mathbf{u}(t_1)) + \bar{n}(\mathcal{E}(t_1))G(\rho(t_1)) - \mathbf{F}(t_1) \cdot \mathbf{u}(t_1) \right) dx, \end{aligned} \tag{72}$$

where  $g(\rho)$  and  $G(\rho)$  are given, respectively, by (55), (58) and (56), (59).

Here,  $\rho^0 \in L^2(\Omega)$ ,  $Q \in C([0, T]; L^2(\Omega))$  and  $\mathbf{F} \in H^1(0, T; L^2(\Omega)^m)$ .

In Definition 1, we explicitly incorporate energy inequality (72). When dealing with classical solutions, Eqs. (67)–(68) imply the energy balance (see (46), (47) and (66))

$$\begin{aligned} & \partial_t J(\mathbf{u}, \rho) + \int_{\Omega} \lambda^*(\partial_t \mathcal{E})^2 dx + \int_{\Omega} k(\mathcal{E})\mathcal{D}(\rho)g'(\rho)\mathbb{K}\nabla\rho \cdot \nabla\rho dx \\ & = \int_{\Omega} Qg(\rho) dx - \int_{\Omega} \partial_t \mathbf{F} \cdot \mathbf{u} dx. \end{aligned} \tag{73}$$

However, in the weak formulation (69)–(70) we cannot use  $\Phi = g(\rho)$  and  $\xi = \partial_t \mathbf{u}$ , due to lack of smoothness. Therefore, (v) has to be added explicitly. Hence, we consider only those weak solutions satisfying additionally (72). Therefore, they are called **weak free energy** solutions.

In a number of steps, we prove existence of weak solutions when  $\lambda^* > 0$ . We achieve this by first considering the incremental formulation. In this approximation, which is clearly relevant when treating the problem numerically, we obtain existence results which hold for all  $\lambda^* \geq 0$ .

### 3 Existence of a Solution to the Incremental Problem

In this section, we study the time discretized form of (67), (68).

In doing so we use the function  $g = g(\rho)$ , defined by (55) and (58) as the primary unknown. This is allowed since  $g : \mathbb{R} \rightarrow \mathbb{R}$  is smooth and strictly increasing. The switch to  $g$  is done for mathematical convenience, because it allows us to obtain Lyapunov functional estimates in a straightforward way. We start with some definitions. Let

$$p(g) := p(\rho(g)) \quad \text{and} \quad \mathcal{D}(g) := \mathcal{D}(\rho(g))\rho'(g). \tag{74}$$

Further, since

$$G(\rho(z)) = \int_{\rho_0}^{\rho(z)} g(\xi) \, d\xi = \int_0^z \zeta \rho'(\zeta) \, d\zeta, \quad z \in \mathbb{R}, \tag{75}$$

let

$$\left. \begin{aligned} G(g) &:= \int_0^g \zeta \rho'(\zeta) \, d\zeta \\ &\text{and, from (62),} \\ D(g, \mathcal{E}) &= \frac{\bar{n}'(\mathcal{E})}{g} G(g) + \frac{p(g) - p_0}{g}. \end{aligned} \right\} \tag{76}$$

Note that the first term in  $D(g, \mathcal{E})$  is bounded with respect to  $\mathcal{E}$  and grows linearly in  $g$  for large  $|g|$ . The second (pressure) term is bounded with respect to  $g$  since

$$p(g) - p_0 = p(\rho(g)) - p_0 = \frac{\rho(g) - \rho_0}{\beta \rho_0}.$$

Using these definitions in (67) and (68), we find in terms of  $g$

$$\bar{n}(\mathcal{E}) \partial_t \rho(g) + D(g, \mathcal{E}) \partial_t \mathcal{E} = \operatorname{div} \left( k(\mathcal{E}) \mathcal{D}(g) \mathbb{K} \nabla g \right) + Q, \tag{77}$$

$$- \operatorname{div} \left( \mathcal{G} e(\mathbf{u}) + \lambda^* \partial_t \mathcal{E} \mathbb{I} - p(g) \mathbb{I} \right) = \mathbf{F}, \tag{78}$$

in  $Q_T$ .

Next, we turn to the time discretized form of Eqs. (77) and (78).

Let  $\tau \in (0, 1)$  denote the time discretization step and  $N \in \mathbb{N}$  a large integer such that  $N\tau = T$ . At each discrete time  $t_j = j\tau$ , with  $j = 0, 1, \dots, N$ , we set

$$\mathbf{F}^j(x) = \mathbf{F}(x, j\tau), \quad Q^j(x) = Q(x, j\tau), \quad x \in \Omega.$$

Let  $\mathbf{u}^{j-1}$  and  $g^{j-1}$  denote, respectively, the displacement and transformed density at  $t_{j-1}$  for some  $j \in \{1, 2, \dots, N\}$ : i.e.,

$$\mathbf{u}^{j-1}(x) = \mathbf{u}(x, t_{j-1}), \quad g^{j-1}(x) = g(x, t_{j-1}), \quad x \in \Omega.$$

Then,  $\mathbf{u}$  and  $g$  at time  $t_j$  are obtained as solutions of the incremental problem (writing  $\mathbf{U} = \mathbf{u}^{j-1}$ ,  $\mathcal{E} = g^{j-1}$  and  $V = H_0^1(\Omega)^m \times H^1(\Omega)$ ):

**Problem (PD):** Given  $(\mathbf{U}, \mathcal{E}) \in V$ , find  $(\mathbf{u}, g) \in V$  such that

$$\begin{aligned} \int_{\Omega} \frac{\bar{n}(\operatorname{div} \mathbf{U})}{\tau} (\rho(g) - \rho(\mathcal{E})) \psi \, dx + \int_{\Omega} D_{\tau}(g, \operatorname{div} \mathbf{u}, \operatorname{div} \mathbf{U}) \operatorname{div} \frac{\mathbf{u} - \mathbf{U}}{\tau} \psi \, dx \\ + \int_{\Omega} k(\operatorname{div} \mathbf{u}) \mathcal{D}(g) \mathbb{K} \nabla g \cdot \nabla \psi \, dx = \int_{\Omega} Q^j \psi \, dx, \quad \forall \psi \in H^1(\Omega); \end{aligned} \tag{79}$$

$$\int_{\Omega} \mathcal{G}e(\mathbf{u}) : e(\xi) \, dx + \frac{\lambda^*}{\tau} \int_{\Omega} \operatorname{div}(\mathbf{u} - \mathbf{U}) \operatorname{div} \xi \, dx - \int_{\Omega} p(g) \operatorname{div} \xi \, dx = \int_{\Omega} \mathbf{F}^j \cdot \xi \, dx, \quad \forall \xi \in H_0^1(\Omega)^m. \tag{80}$$

The coefficient  $D_{\tau}$  in Eq. (79) is given by

$$D_{\tau}(g, \operatorname{div} \mathbf{u}, \operatorname{div} \mathbf{U}) = \frac{\bar{n}(\operatorname{div} \mathbf{u}) - \bar{n}(\operatorname{div} \mathbf{U})}{\operatorname{div} \mathbf{u} - \operatorname{div} \mathbf{U}} \frac{G(g)}{g} + \frac{p(g) - p_0}{g}. \tag{81}$$

This expression results from  $D(g, \mathcal{E})$  in (76), when the derivative  $\bar{n}'(\mathcal{E})$  is replaced by the finite difference  $\frac{\bar{n}(\operatorname{div} \mathbf{u}) - \bar{n}(\operatorname{div} \mathbf{U})}{\operatorname{div} \mathbf{u} - \operatorname{div} \mathbf{U}}$ . The specific choice of (81) appears convenient in the estimates concerning the time-discrete Lyapunov functional.

Using the weak topology of the space  $H_0^1(\Omega)^m \times H^1(\Omega)$ , serious difficulties arise with the coefficients  $\bar{n}$ ,  $D_{\tau}$  and  $k$  depending on  $\operatorname{div} \mathbf{u}$ . To remedy this, we introduce a mollifier  $\mathcal{Y}_{\varepsilon}$ , where  $\varepsilon$  is a small positive parameter (see, e.g., Roubiřek 2005, page 203), and replace  $\operatorname{div} \mathbf{u}$  in the nonlinearities by the convolution  $\operatorname{div} \mathbf{u} \star \mathcal{Y}_{\varepsilon} = \mathbf{u} \star \nabla \mathcal{Y}_{\varepsilon}$ . Using this substitution one can treat nonlinear coefficients containing  $\operatorname{div} \mathbf{u}$  as lower order terms in the equations. This allows us to use the theory of pseudo-monotone operators.

Applying this convolution, the regularized form of Problem (PD) reads:

Problem (PD) $_{\varepsilon}$ : Given  $(\mathbf{U}, \mathcal{E}) \in V$ , find  $(\mathbf{u}_{\varepsilon}, g_{\varepsilon}) \in V$  such that, with  $\mathcal{E}_{\varepsilon} = \mathbf{u}_{\varepsilon} \star \nabla \mathcal{Y}_{\varepsilon}$ ,

$$\int_{\Omega} \frac{\bar{n}(\operatorname{div} \mathbf{U})}{\tau} (\rho(g_{\varepsilon}) - \rho(\mathcal{E})) \psi \, dx + \int_{\Omega} \left( \frac{\bar{n}(\mathcal{E}_{\varepsilon}) - \bar{n}(\operatorname{div} \mathbf{U})}{\tau g_{\varepsilon}} G(g_{\varepsilon}) + \frac{p(g_{\varepsilon}) - p_0}{\tau g_{\varepsilon}} \right) \operatorname{div}(\mathbf{u}_{\varepsilon} - \mathbf{U}) \psi \, dx + \int_{\Omega} k(\mathcal{E}_{\varepsilon}) \mathcal{D}(g_{\varepsilon}) \mathbb{K} \nabla g_{\varepsilon} \cdot \nabla \psi \, dx = \int_{\Omega} Q^j \psi \, dx, \quad \forall \psi \in H^1(\Omega), \tag{82}$$

$$\int_{\Omega} \mathcal{G}e(\mathbf{u}_{\varepsilon}) : e(\xi) \, dx + \frac{\lambda^*}{\tau} \int_{\Omega} \operatorname{div}(\mathbf{u}_{\varepsilon} - \mathbf{U}) \operatorname{div} \xi \, dx - \int_{\Omega} p(g_{\varepsilon}) \operatorname{div} \xi \, dx = \int_{\Omega} \mathbf{F}^j \cdot \xi \, dx, \quad \forall \xi \in H_0^1(\Omega)^m. \tag{83}$$

Note that the denominator  $\tau g_{\varepsilon}$  in (82) originates from the time step  $\tau$  in the discretization and the term  $g$  in the denominator of (76). We have the following existence result

**Proposition 2** *Let  $\varepsilon > 0$  be a small positive constant. Under the assumptions of Definition 1, Problem (PD) $_{\varepsilon}$  admits at least one solution  $(\mathbf{u}_{\varepsilon}, g_{\varepsilon}) \in V$ .*

**Proof** We start by introducing a nonlinear operator  $\mathcal{A}$ , defined on  $V$  and with values in its dual  $V'$ . It results from adding (82) and (83). We write the resulting relation, with  $(\mathbf{u}, g) \in V$ , as

$$\mathcal{A}(\mathbf{u}, g) = \mathbf{b}, \tag{84}$$

where

$$\begin{aligned}
 \langle \mathcal{A}(\mathbf{u}, g), (\xi, \psi) \rangle &:= \frac{1}{\tau} \int_{\Omega} \mathcal{G}e(\mathbf{u}) : e(\xi) \, dx + \frac{\lambda^*}{\tau^2} \int_{\Omega} \operatorname{div}(\mathbf{u} - \mathbf{U}) \operatorname{div} \xi \, dx \\
 &- \int_{\Omega} \frac{p(g)}{\tau} \operatorname{div} \xi \, dx + \int_{\Omega} k(\mathbf{u} \star \nabla \gamma_{\varepsilon}) \mathcal{D}(g) \mathbb{K} \nabla g \cdot \nabla \psi \, dx \\
 &+ \int_{\Omega} \frac{\bar{n}(\operatorname{div} \mathbf{U})}{\tau} (\rho(g) - \rho(\mathcal{E})) \psi \, dx \\
 &+ \int_{\Omega} \left( \frac{\bar{n}(\mathbf{u} \star \nabla \gamma_{\varepsilon}) - \bar{n}(\operatorname{div} \mathbf{U})}{\tau g} G(g) + \frac{p(g) - p_0}{\tau g} \right) \operatorname{div}(\mathbf{u} - \mathbf{U}) \psi \, dx, \quad \forall (\xi, \psi) \in V.
 \end{aligned} \tag{85}$$

and

$$\langle b, (\xi, \psi) \rangle := \int_{\Omega} \mathbf{F}^j \cdot \xi \, dx + \int_{\Omega} Q^j \psi \, dx, \quad \forall (\xi, \psi) \in V. \tag{86}$$

The idea is to show that  $\mathcal{A}$  is a perturbed monotone operator: i.e.,  $\mathcal{A}$  is monotone in its principal part containing derivatives of  $\mathbf{u}$  and  $g$ . To be precise, we show that  $\mathcal{A}$  is pseudo-monotone and coercive. This allows to apply Brézis’ theorem to (84) (see Chapter 2 in monographs Lions 1969; Roubiřek 2005 or Chapters 26 and 27 in Zeidler (1990)) to conclude existence for Problem  $(\mathbf{PD})_{\varepsilon}$ .

For the comfort of the reader, we recall that an operator  $\mathcal{A} : V \rightarrow V'$  is pseudo-monotone if and only if  $\mathcal{A}$  is bounded and

$$\left. \begin{aligned}
 \{\mathbf{u}_r, g_r\} \rightharpoonup \{\mathbf{u}, g\} \text{ weakly in } V, \\
 \limsup_{r \rightarrow +\infty} \langle \mathcal{A}(\mathbf{u}_r, g_r), (\mathbf{u}_r, g_r) - (\mathbf{u}, g) \rangle \leq 0,
 \end{aligned} \right\} \Rightarrow \begin{cases} \forall (\mathbf{v}, h) \in V, \\ \langle \mathcal{A}(\mathbf{u}, g), (\mathbf{u}, g) - (\mathbf{v}, h) \rangle \leq \\ \liminf_{r \rightarrow +\infty} \langle \mathcal{A}(\mathbf{u}_r, g_r), (\mathbf{u}_r, g_r) - (\mathbf{v}, h) \rangle, \end{cases} \tag{87}$$

The boundedness of  $\mathcal{A}$  is immediate. To show (87), we follow Chapter 2 from Roubiřek (2005) or Chapter 17 from Schweizer (2018) and rewrite  $\mathcal{A}$  in a form having a principal part containing partial derivatives of  $\mathbf{u}$  (in  $e(\mathbf{u})$  and  $\operatorname{div} \mathbf{u}$ ) and  $\nabla g$ , and a lower order part containing  $\mathbf{u}$  and  $g$ . Specifically, we introduce the operator  $\mathcal{B} : V \times V \rightarrow V'$  by

$$\begin{aligned}
 \langle \mathcal{B}((\mathbf{w}, \ell), (\mathbf{u}, g)), (\xi, \psi) \rangle &= \frac{1}{\tau} \int_{\Omega} \mathcal{G}e(\mathbf{u}) : e(\xi) \, dx + \frac{\lambda^*}{\tau^2} \int_{\Omega} \operatorname{div}(\mathbf{u} - \mathbf{U}) \operatorname{div} \xi \, dx \\
 &- \int_{\Omega} \frac{p(\ell)}{\tau} \operatorname{div} \xi \, dx + \int_{\Omega} \left( \frac{\bar{n}(\mathbf{w} \star \nabla \gamma_{\varepsilon}) - \bar{n}(\operatorname{div} \mathbf{U})}{\tau \ell} G(\ell) + \frac{p(\ell) - p_0}{\tau \ell} \right) \operatorname{div}(\mathbf{u} - \mathbf{U}) \psi \, dx \\
 &+ \int_{\Omega} \frac{\bar{n}(\operatorname{div} \mathbf{U})}{\tau} (\rho(\ell) - \rho(\mathcal{E})) \psi \, dx + \int_{\Omega} k(\mathbf{w} \star \nabla \gamma_{\varepsilon}) \mathcal{D}(\ell) \mathbb{K} \nabla g \cdot \nabla \psi \, dx, \quad \forall (\xi, \psi) \in V.
 \end{aligned} \tag{88}$$

We observe that  $\mathcal{B}((\mathbf{u}, g), (\mathbf{u}, g)) = \mathcal{A}(\mathbf{u}, g)$ . The introduction of  $\mathcal{B}$  is useful because it reflects the monotonicity of the principal part of  $\mathcal{A}(\mathbf{u}, g)$ . This is a direct consequence

of

$$\langle \mathcal{B}((\mathbf{w}, \ell), (\mathbf{u}_1, g_1)) - \mathcal{B}((\mathbf{w}, \ell), (\mathbf{u}_2, g_2)), (\mathbf{u}_1, g_1) - (\mathbf{u}_2, g_2) \rangle \geq 0, \tag{89}$$

with equality if and only if  $\mathbf{u}_1 = \mathbf{u}_2$  and  $g_1 = g_2$ . Inequality (89) is checked by a short computation in (88).

To show (87) we consider a sequence  $\{\mathbf{u}_r, g_r\} \subset V$  such that

$$(\mathbf{u}_r, g_r) \rightharpoonup (\mathbf{u}, g) \text{ weakly in } V \text{ and } \limsup_{r \rightarrow +\infty} \langle \mathcal{A}(\mathbf{u}_r, g_r), (\mathbf{u}_r, g_r) - (\mathbf{u}, g) \rangle \leq 0. \tag{90}$$

As in Roubiřek (2005) we set  $(\mathbf{u}_\delta, g_\delta) = (1 - \delta)(\mathbf{u}, g) + \delta(\mathbf{v}, h)$ , where  $\delta \in [0, 1]$  and  $(\mathbf{v}, h) \in V$ . Using the monotonicity from (89), we obtain

$$\begin{aligned} \delta \langle \mathcal{A}(\mathbf{u}_r, g_r), (\mathbf{u}, g) - (\mathbf{v}, h) \rangle &\geq -\langle \mathcal{A}(\mathbf{u}_r, g_r), (\mathbf{u}_r, g_r) - (\mathbf{u}, g) \rangle \\ &+ \langle \mathcal{B}((\mathbf{u}_r, g_r), (\mathbf{u}_\delta, g_\delta)), (\mathbf{u}_r, g_r) - (\mathbf{u}, g) \rangle + \delta \langle \mathcal{B}((\mathbf{u}_r, g_r), (\mathbf{u}_\delta, g_\delta)), \\ &(\mathbf{u}, g) - (\mathbf{v}, h) \rangle. \end{aligned} \tag{91}$$

The sequence  $(\mathbf{u}_r, g_r)$  is bounded in  $V$  and there exists a subsequence which strongly converges in  $L^5(\Omega)^m$  and (a.e.) in  $\Omega$ , to  $(\mathbf{u}, g)$ . Hence, it suffices to pass to the limit along this subsequence. In (91), the terms containing the operator  $\mathcal{B}$  are fixed with respect to the gradients. Hence,

$$\lim_{r \rightarrow +\infty} \langle \mathcal{B}((\mathbf{u}_r, g_r), (\mathbf{u}_\delta, g_\delta)), (\mathbf{u}_r, g_r) - (\mathbf{u}, g) \rangle = 0, \tag{92}$$

and

$$\lim_{r \rightarrow +\infty} \langle \mathcal{B}((\mathbf{u}_r, g_r), (\mathbf{v}, h)), (\xi, \psi) \rangle = \langle \mathcal{B}((\mathbf{u}, g), (\mathbf{v}, h)), (\xi, \psi) \rangle, \tag{93}$$

for any  $(\xi, \psi) \in V$ . With these results, we are in a position to pass to the limit  $r \rightarrow +\infty$  in inequality (91). It yields

$$\begin{aligned} \delta \liminf_{r \rightarrow +\infty} \langle \mathcal{A}(\mathbf{u}_r, g_r), (\mathbf{u}, g) - (\mathbf{v}, h) \rangle &\geq -\limsup_{r \rightarrow +\infty} \langle \mathcal{A}(\mathbf{u}_r, g_r), (\mathbf{u}_r, g_r) - (\mathbf{u}, g) \rangle \\ &+ \delta \langle \mathcal{B}((\mathbf{u}, g), (\mathbf{u}_\delta, g_\delta)), (\mathbf{u}, g) - (\mathbf{v}, h) \rangle. \end{aligned} \tag{94}$$

By the pseudo-monotonicity hypothesis (90), inequality (94) implies

$$\begin{aligned} \liminf_{r \rightarrow +\infty} \langle \mathcal{A}(\mathbf{u}_r, g_r), (\mathbf{u}, g) - (\mathbf{v}, h) \rangle &\geq \langle \mathcal{B}((\mathbf{u}, g), (\mathbf{u}, g)), (\mathbf{u}, g) - (\mathbf{v}, h) \rangle \\ &= \langle \mathcal{A}(\mathbf{u}, g), (\mathbf{u}, g) - (\mathbf{v}, h) \rangle, \forall (\mathbf{v}, h) \in V. \end{aligned} \tag{95}$$

We use this inequality to conclude

$$\liminf_{r \rightarrow +\infty} \langle \mathcal{A}(\mathbf{u}_r, g_r), (\mathbf{u}_r, g_r) - (\mathbf{v}, h) \rangle$$

$$\begin{aligned}
 &\geq \liminf_{r \rightarrow +\infty} \langle \mathcal{A}(\mathbf{u}_r, g_r), (\mathbf{u}_r, g_r) - (\mathbf{u}, g) \rangle \\
 &\quad + \liminf_{r \rightarrow +\infty} \langle \mathcal{A}(\mathbf{u}_r, g_r), (\mathbf{u}, g) - (\mathbf{v}, h) \rangle \\
 &= \liminf_{r \rightarrow +\infty} \underbrace{\langle \mathcal{B}((\mathbf{u}_r, g_r), (\mathbf{u}, g)), (\mathbf{u}_r, g_r) - (\mathbf{u}, g) \rangle}_{=0 \text{ by (92)}} \\
 &\quad + \liminf_{r \rightarrow +\infty} \underbrace{\langle \mathcal{B}((\mathbf{u}_r, g_r), (\mathbf{u}_r, g_r)) - \mathcal{B}((\mathbf{u}_r, g_r), (\mathbf{u}, g)), (\mathbf{u}_r, g_r) - (\mathbf{u}, g) \rangle}_{\geq 0 \text{ by (89)}} \\
 &\quad + \liminf_{r \rightarrow +\infty} \langle \mathcal{A}(\mathbf{u}_r, g_r), (\mathbf{u}, g) - (\mathbf{v}, h) \rangle \geq \liminf_{r \rightarrow +\infty} \langle \mathcal{A}(\mathbf{u}_r, g_r), (\mathbf{u}, g) - (\mathbf{v}, h) \rangle \\
 &\geq \langle \mathcal{A}(\mathbf{u}, g), (\mathbf{u}, g) - (\mathbf{v}, h) \rangle, \quad \forall \{\mathbf{v}, h\} \in V.
 \end{aligned}$$

This completes the proof of the *pseudo-monotonicity*.

It remains to prove *coercivity*. We evaluate directly the term

$$\langle \mathcal{A}(\mathbf{u}, g), (\mathbf{u} - \mathbf{U}, g) \rangle.$$

Taking  $\xi = \mathbf{u} - \mathbf{U}$  and  $\psi = g$  in (85), the cross terms involving the product  $p(g) \operatorname{div}(\mathbf{u} - \mathbf{U})$  cancel and the term  $p_0 \operatorname{div}(\mathbf{u} - \mathbf{U})/\tau$  drops out after integration. What remains is

$$\begin{aligned}
 &\langle \mathcal{A}(\mathbf{u}, g), (\mathbf{u} - \mathbf{U}, g) \rangle \\
 &= \frac{1}{\tau} \int_{\Omega} \mathcal{G}e(\mathbf{u}) : e(\mathbf{u} - \mathbf{U}) \, dx + \frac{\lambda^*}{\tau^2} \int_{\Omega} \operatorname{div}(\mathbf{u} - \mathbf{U})^2 \, dx \\
 &\quad + \int_{\Omega} \frac{\bar{n}(\operatorname{div} \mathbf{U})}{\tau} (\rho(g) - \rho(\mathcal{E}))g \, dx \\
 &\quad + \int_{\Omega} \frac{\bar{n}(\mathbf{u} \star \nabla \mathcal{Y}_\varepsilon) - \bar{n}(\operatorname{div} \mathbf{U})}{\tau} G(g) \, dx + \int_{\Omega} \mathbb{K}k(\mathbf{u} \star \nabla \mathcal{Y}_\varepsilon) \mathcal{D}(g) |\nabla g|^2 \, dx.
 \end{aligned} \tag{96}$$

The third and fourth terms in the right-hand side need special attention.

Since  $\rho = \rho(g)$  is a  $C^1$  monotonically increasing function, we have the elementary inequality

$$x(\rho(x) - \rho(y)) \geq \int_y^x \zeta \rho'(\zeta) \, d\zeta, \quad \forall x, y \in \mathbb{R}. \tag{97}$$

Using this inequality and the expression for  $G$  (see (76)) in these terms gives

$$\begin{aligned}
 &\bar{n}(\operatorname{div} \mathbf{U})(\rho(g) - \rho(\mathcal{E}))g + (\bar{n}(\mathbf{u} \star \nabla \mathcal{Y}_\varepsilon) - \bar{n}(\operatorname{div} \mathbf{U})) \int_0^g \zeta \rho'(\zeta) \, d\zeta \\
 &\geq \bar{n}(\mathbf{u} \star \nabla \mathcal{Y}_\varepsilon) \int_0^g \zeta \rho'(\zeta) \, d\zeta - \bar{n}(\operatorname{div} \mathbf{U}) \int_0^{\mathcal{E}} \zeta \rho'(\zeta) \, d\zeta.
 \end{aligned} \tag{98}$$

Applying Korn’s inequality, see Theorem 1.33 from Roubiřek (2005), and inserting inequality (98) into equality (96) yields

$$\langle \mathcal{A}(\mathbf{u}, g), (\mathbf{u} - \mathbf{U}, g) \rangle \geq \frac{C_1}{\tau} \|\mathbf{u}\|_{H_0^1(\Omega)^m}^2 + C_2 \|\nabla g\|_{L^2(\Omega)^m}^2 - \frac{C_3}{\tau} + \frac{C_4}{\tau} \int_{\Omega} \underbrace{\left( \int_0^g \zeta \rho'(\zeta) \, d\zeta \right) dx}_{\approx Cg^2 \text{ for large } |g|}, \tag{99}$$

where  $C_i, i = 1, \dots, 4$  are positive constants. This proves the coercivity.

Having established pseudo-monotonicity and coercivity of the operator  $\mathcal{A}$ , we are in position to apply Br ezis’ theorem. This concludes the assertion of the proposition.  $\square$

**Theorem 1** *Problem (PD) admits at least one solution  $(\mathbf{u}, g) \in V$ .*

**Proof** For each  $\varepsilon > 0$ , let  $(\mathbf{u}_\varepsilon, g_\varepsilon)$  be a solution of Problem  $(\mathbf{PD}_\varepsilon)$  as obtained in Proposition 2. From the coercivity part of the proof of Proposition 2 and Eq. (84), it follows that

$$\|\mathbf{u}_\varepsilon\|_{H_0^1(\Omega)^m} + \|g_\varepsilon\|_{H^1(\Omega)} \leq C, \tag{100}$$

where  $C$  is independent of  $\varepsilon$ . Estimate (100) yields weak compactness in  $H^1$ . However, this is not enough to prove that  $\mathbf{u}_\varepsilon \star \nabla \gamma_\varepsilon$  converges strongly in  $L^2$  and (a.e.) on  $\Omega$  as  $\varepsilon \rightarrow 0$ . The remedy is to consider the momentum Eq. (83), which gives us improved regularity through the elasticity term. Since  $p(g_\varepsilon)$  is bounded in  $H^1(\Omega)$ , uniformly with respect to  $\varepsilon$ , we conclude that

$$\|\mathbf{u}_\varepsilon\|_{H^2(\Omega)^m} \leq C, \tag{101}$$

where  $C$  does not depend on  $\varepsilon$ . Using estimates (100)–(101), there is a subsequence  $(\mathbf{u}_\varepsilon, g_\varepsilon)$ , denoted by the same subscript, and a pair  $(\mathbf{u}, g) \in (H_0^1(\Omega)^m \cap H^2(\Omega)^m) \times H^1(\Omega)$  such that

$$\mathbf{u}_\varepsilon \rightarrow \mathbf{u} \text{ strongly in } H_0^1(\Omega)^m, \tag{102}$$

$$\operatorname{div} \mathbf{u}_\varepsilon \rightarrow \operatorname{div} \mathbf{u} \text{ strongly in } L^2(\Omega) \text{ and (a.e.) on } \Omega, \tag{103}$$

$$g_\varepsilon \rightharpoonup g \text{ weakly in } H^1(\Omega), \tag{104}$$

$$g_\varepsilon \rightarrow g \text{ strongly in } L^2(\Omega) \text{ and (a.e.) on } \Omega, \tag{105}$$

as  $\varepsilon \rightarrow 0$ . The convergence properties allow to pass to the limit in system (82)–(83). Hence, the pair  $(\mathbf{u}, g)$  satisfies the equations of Problem (PD), which proves the theorem.  $\square$

To complete the study of the incremental problem, we need to estimate the behavior of solutions after at least  $O(1/\tau)$  times steps. Here, we use the discrete version of Lyapunov functional (66).



In Problem **(PD)**, where the discrete time step  $\tau$  enters as parameter, one find after one step  $(\mathbf{u}^1, g^1)$  from the initial values  $(\operatorname{div} \mathbf{u}, \rho)|_{t=0} = (0, \rho^0)$ . The idea is to repeat this procedure for an arbitrary number of steps. If  $M \in \mathbb{N}$ ,  $M \leq N = T/\tau$ , then  $(\mathbf{u}^M, g^M)$  denotes the time discretized approximation of the original quasi-static equation, at  $t = t_M = M\tau$ .

The corresponding Lyapunov functional at  $t = t_M$  reads

$$J^M = \int_{\Omega} \left( \frac{1}{2} \mathcal{G}e(\mathbf{u}^M) : e(\mathbf{u}^M) - \mathbf{F}^M \cdot \mathbf{u}^M + \bar{n}(\operatorname{div} \mathbf{u}^M)G(g^M) \right) dx. \tag{106}$$

It satisfies

**Theorem 2** For each  $M \in \mathbb{N}$ ,  $M \leq N = T/\tau$ , we have

$$J^M + \tau \sum_{j=1}^M \int_{\Omega} \left( \lambda^* \left( \frac{\operatorname{div}(\mathbf{u}^j - \mathbf{u}^{j-1})}{\tau} \right)^2 + \frac{\mathbf{F}^j - \mathbf{F}^{j-1}}{\tau} \cdot \mathbf{u}^{j-1} + k(\operatorname{div} \mathbf{u}^j)\mathcal{D}(g^j)\mathbb{K}\nabla g^j \cdot \nabla g^j - Q^j g^j \right) dx \leq J^0. \tag{107}$$

Here,

$$J^0 = n_0 \int_{\Omega} G(g^0) dx, \quad g^0 = g(\rho^0).$$

**Proof** At time  $t = t_j$ , with  $j = 1, \dots, N$ , the equations in Problem **(PD)** read

$$\begin{aligned} & \int_{\Omega} \mathcal{G}e(\mathbf{u}^j) : e(\xi) dx + \frac{\lambda^*}{\tau} \int_{\Omega} \operatorname{div}(\mathbf{u}^j - \mathbf{u}^{j-1}) \operatorname{div} \xi dx - \int_{\Omega} p(g^j) \operatorname{div} \xi dx \\ & = \int_{\Omega} \mathbf{F}^j \cdot \xi dx, \quad \forall \xi \in H_0^1(\Omega)^m, \tag{108} \\ & \int_{\Omega} \left( \frac{\bar{n}(\operatorname{div} \mathbf{u}^{j-1})}{\tau} (\rho(g^j) - \rho(g^{j-1})) + \frac{\bar{n}(\operatorname{div} \mathbf{u}^j) - \bar{n}(\operatorname{div} \mathbf{u}^{j-1})}{\tau g^j} G(g^j) \right) \psi dx \\ & + \int_{\Omega} \frac{p(g^j) - p_0}{\tau g^j} \operatorname{div}(\mathbf{u}^j - \mathbf{u}^{j-1}) \psi dx + \int_{\Omega} k(\operatorname{div} \mathbf{u}^j)\mathcal{D}(g^j)\mathbb{K}\nabla g^j \cdot \nabla \psi dx \\ & = \int_{\Omega} Q^j \psi dx, \quad \forall \psi \in H^1(\Omega). \end{aligned} \tag{109}$$

Note that in Eq. (109) we have used explicitly the form of  $D_{\tau}$  from (81). Next, we take  $\xi = (\mathbf{u}^j - \mathbf{u}^{j-1})/\tau$  in (108) and  $\psi = g^j$  in (109). The resulting two equalities are added and summed up with respect to  $j$  up from  $j = 1$  to  $j = M$ . Using the observations

- (i) cross terms containing pressure cancel;

(ii)

$$\sum_{j=1}^M \mathcal{G}e(\mathbf{u}^j) : e(\mathbf{u}^j - \mathbf{u}^{j-1}) \geq \frac{1}{2} (\mathcal{G}e(\mathbf{u}^M) : e(\mathbf{u}^M) - \mathcal{G}e(\mathbf{u}^0) : e(\mathbf{u}^0));$$

(iii)

$$\begin{aligned} & \sum_{j=1}^M \left( \bar{n}(\operatorname{div} \mathbf{u}^{j-1}) g^j (\rho(g^j) - \rho(g^{j-1})) + (\bar{n}(\operatorname{div} \mathbf{u}^j) - \bar{n}(\operatorname{div} \mathbf{u}^{j-1})) G(g^j) \right) \\ & \geq \bar{n}(\operatorname{div} \mathbf{u}^M) G(g^M) - \bar{n}(\operatorname{div} \mathbf{u}^0) G(g^0), \end{aligned}$$

where (97) is used;

(iv)

$$\sum_{j=1}^M \mathbf{F}^j \cdot (\mathbf{u}^j - \mathbf{u}^{j-1}) = \mathbf{F}^M \cdot \mathbf{u}^M - \mathbf{F}^0 \cdot \mathbf{u}^0 - \sum_{j=0}^{M-1} (\mathbf{F}^{j+1} - \mathbf{F}^j) \cdot \mathbf{u}^j,$$

one finds inequality (107). The reduced expression for  $J^0$  results from  $\mathbf{u}|_{t=0} = 0$ .

□

Having established existence for the discrete Problem (PD) in Theorem 1 and a Lyapunov estimate in Theorem 2, we are now in a position to obtain estimates that are uniform in the time step  $\tau$ .

**Proposition 3** *There exists a constant  $C > 0$  such that*

$$\|\mathbf{u}^M\|_{H^1(\Omega)^m}^2 + \|g^M\|_{L^2(\Omega)}^2 \leq C, \tag{110}$$

and

$$\tau \sum_{j=1}^M \int_{\Omega} \left( \lambda^* \left( \frac{\operatorname{div}(\mathbf{u}^j - \mathbf{u}^{j-1})}{\tau} \right)^2 + |\nabla g^j|^2 \right) dx \leq C, \tag{111}$$

for all  $M$  and  $\tau$  such that  $1 \leq M \leq N = T/\tau$ , with  $\tau$  sufficiently small.

**Proof** Combining expression (106) for  $J^M$  and inequality (107) yields for any  $1 \leq M \leq N$

$$\begin{aligned} & \frac{1}{2} \int_{\Omega} \mathcal{G}e(\mathbf{u}^M) : e(\mathbf{u}^M) dx + \int_{\Omega} \bar{n}(\operatorname{div} \mathbf{u}^M) G(g^M) dx \\ & \leq \int_{\Omega} \mathbf{F}^M \cdot \mathbf{u}^M dx + J^0 + \tau \sum_{j=1}^M \int_{\Omega} Q^j g^j dx \end{aligned}$$

$$\begin{aligned}
 & +\tau \sum_{j=1}^M \int_{\Omega} \frac{\mathbf{F}^j - \mathbf{F}^{j-1}}{\tau} \cdot \mathbf{u}^{j-1} \, dx \leq \frac{\delta}{2} \|\mathbf{u}^M\|_{L^2(\Omega)^m}^2 + \frac{1}{2\delta} \|\mathbf{F}^M\|_{L^2(\Omega)^m}^2 + J^0 \\
 & + \frac{\tau}{2} \sum_{j=1}^M \|g^j\|_{L^2(\Omega)}^2 + \frac{\tau}{2} \sum_{j=1}^M \|\mathbf{u}^{j-1}\|_{L^2(\Omega)^m}^2 \\
 & + \frac{\tau}{2} \sum_{j=1}^M \|Q^j\|_{L^2(\Omega)}^2 + \frac{\tau}{2} \sum_{j=1}^M \left\| \frac{\mathbf{F}^j - \mathbf{F}^{j-1}}{\tau} \right\|_{L^2(\Omega)^m}^2.
 \end{aligned}$$

By the assumptions on  $Q$  and  $\mathbf{F}$ , the last two terms are uniformly bounded with respect to  $\tau$  and  $M$ . We estimate the left-hand side from below by applying Korn’s inequality to the first term and the quadratic growth of  $G$  to the second term. Then for  $\delta$  and  $\tau$  sufficiently small, we obtain for the combination

$$\mathcal{U}_j = \|\mathbf{u}^j\|_{H^1(\Omega)^m}^2 + \|g^j\|_{L^2(\Omega)}^2, \quad j = 0, \dots, M,$$

the inequality

$$\mathcal{U}_M \leq C_1 + C_2\tau \sum_{j=0}^{M-1} \mathcal{U}_j,$$

where  $C_1$  and  $C_2$  do not dependent on  $\tau$  and  $M$ . Next, we apply the *discrete Gronwall inequality*<sup>1</sup>, see footnote, to find

$$\mathcal{U}_M \leq C_1 e^{C_2(M-1)\tau} < C_1 e^{C_2T} \quad \text{for all } 1 \leq M \leq N.$$

The second estimate follows directly from Theorem 2. □

However, to pass to the limit  $\tau \rightarrow 0$  in the nonlinearities, one needs more information on the behavior of the ratios  $\{\operatorname{div} (u^j - u^{j-1})/\tau\}$  and  $\{(g^j - g^{j-1})/\tau\}$ . In fact, we must establish relative compactness of the sequences  $\{\operatorname{div} \mathbf{u}^j\}$  and  $\{g^j\}$ .

We start with a local  $H^1$ -estimate for  $\mathcal{E}^j = \operatorname{div} \mathbf{u}^j$ .

**Lemma 1** *Let  $\varphi \in C_0^\infty(\Omega)$  and  $\tau > 0$  sufficiently small. Then, there exists a constant  $C = C(\varphi)$  such that*

$$\tau \sum_{j=1}^N \|\varphi \mathcal{E}^j\|_{H^1(\Omega)}^2 + \frac{\lambda^*}{2\mu + \lambda} \max_{1 \leq M \leq N} \|\varphi \mathcal{E}^M\|_{H^1(\Omega)}^2 \leq C. \tag{112}$$

**Proof** Let

$$L^j = (2\mu + \lambda)\mathcal{E}^j - p(g^j) + \lambda^* \frac{\operatorname{div} (\mathbf{u}^j - \mathbf{u}^{j-1})}{\tau}, \quad j = 1, \dots, M. \tag{113}$$

<sup>1</sup> *Discrete version of Gronwall’s lemma:* Let  $\{\mathcal{U}_n\}$  and  $\{w_n\}$  be nonnegative sequences satisfying  $\mathcal{U}_n \leq A + \sum_{j=0}^{n-1} \mathcal{U}_j w_j$ . Then for all  $n$ ,  $\mathcal{U}_n \leq A \exp\{\sum_{j=0}^{n-1} w_j\}$ .

Inequality (110) implies

$$\tau \sum_{j=1}^M \left( \|\mathbf{u}^j\|_{H^1(\Omega)^m}^2 + \|g^j\|_{L^2(\Omega)}^2 \right) \leq \tau MC \leq TC.$$

Combined with (111) this gives for  $L^j$

$$\tau \sum_{j=1}^M \|L^j\|_{L^2(\Omega)}^2 \leq C. \tag{114}$$

As in the counterexample for negative porosity, we take the divergence of the time-discrete momentum equation. This yields

$$-\Delta L^j = \operatorname{div} \mathbf{F}^j \quad \text{in } \Omega. \tag{115}$$

In general, however, there are no boundary conditions for  $L^j$  available. Here, we must rely on local estimates to obtain (112). Let us first write the equation for  $\varphi L^j \in H^2(\Omega) \cap H_0^1(\Omega)$ :

$$\Delta(\varphi L^j) = -\varphi \operatorname{div} \mathbf{F}^j + 2\nabla\varphi \cdot \nabla L^j + L^j \Delta\varphi.$$

Its weak form reads

$$\int_{\Omega} \nabla(\varphi L^j) \nabla \zeta \, dx = \int_{\Omega} \operatorname{div} \mathbf{F}^j \varphi \zeta \, dx + \int_{\Omega} L^j (2\nabla\varphi \cdot \nabla \zeta + \zeta \Delta\varphi) \, dx, \quad \forall \zeta \in H_0^1(\Omega). \tag{116}$$

Taking  $\zeta = \varphi L^j$  results in

$$\begin{aligned} \int_{\Omega} |\nabla(\varphi L^j)|^2 \, dx &= - \int_{\Omega} \mathbf{F}^j \cdot \nabla(\varphi L^j) \varphi \, dx - \int_{\Omega} \mathbf{F}^j \cdot \nabla\varphi \varphi L^j \, dx \\ &\quad + \int_{\Omega} (L^j)^2 \varphi \Delta\varphi \, dx - \int_{\Omega} 2L^j \nabla\varphi \cdot \nabla(\varphi L^j) \, dx. \end{aligned}$$

With  $C = C(\varphi)$  denoting a generic constant depending on  $\varphi$ , we have

$$\|\varphi L^j\|_{H^1(\Omega)}^2 \leq C(\|\mathbf{F}^j\|_{L^2(\Omega)}^2 + \|L^j\|_{L^2(\Omega)}^2), \tag{117}$$

for  $1 \leq j \leq M \leq N$ . Combining this inequality with (114) gives

$$\tau \sum_{j=1}^M \|\varphi L^j\|_{H^1(\Omega)}^2 \leq C(\tau \sum_{j=1}^M \|\mathbf{F}^j\|_{L^2(\Omega)^m}^2 + 1) \leq C. \tag{118}$$

Next, we multiply expression (113) by  $\tau\varphi$  and write it as

$$\tau(2\mu + \lambda)\varphi\mathcal{E}^j + \lambda^*\varphi(\mathcal{E}^j - \mathcal{E}^{j-1}) = \tau L^j\varphi + \tau p(g^j)\varphi \in H^1(\Omega).$$

Taking the  $H^1$ -inner product of this expression with  $\varphi\mathcal{E}^j$  gives

$$\begin{aligned} &(2\mu + \lambda)\tau\|\varphi\mathcal{E}^j\|_{H^1(\Omega)}^2 + \lambda^*(\varphi(\mathcal{E}^j - \mathcal{E}^{j-1}), \varphi\mathcal{E}^j)_{H^1(\Omega)} \\ &= \tau(\varphi(L^j + p(g^j)), \varphi\mathcal{E}^j)_{H^1(\Omega)} \end{aligned}$$

or

$$\begin{aligned} &(2\mu + \lambda)\frac{\tau}{2}\|\varphi\mathcal{E}^j\|_{H^1(\Omega)}^2 + \lambda^*(\varphi(\mathcal{E}^j - \mathcal{E}^{j-1}), \varphi\mathcal{E}^j)_{H^1(\Omega)} \\ &\leq \frac{\tau}{2(2\mu + \lambda)}\|\varphi(L^j + p(g^j))\|_{H^1(\Omega)}^2. \end{aligned} \tag{119}$$

Using the identity

$$\sum_{j=1}^M a^j(a^j - a^{j-1}) = \frac{(a^M)^2}{2} - \frac{(a^0)^2}{2} + \frac{1}{2}\sum_{j=1}^M (a^j - a^{j-1})^2,$$

when summing up (119) gives

$$\begin{aligned} &\tau\sum_{j=1}^M \|\varphi\mathcal{E}^j\|_{H^1(\Omega)}^2 + \frac{\lambda^*}{2\mu + \lambda}\|\varphi\mathcal{E}^M\|_{H^1(\Omega)}^2 \\ &\leq + \frac{\tau}{(2\mu + \lambda)^2}\sum_{j=1}^M \|\varphi(L^j + p(g^j))\|_{H^1(\Omega)}^2. \end{aligned}$$

Combining this inequality with (111) and (118), results in the estimate of the lemma. □

We conclude this section with an estimate for  $(\rho(g^j) - \rho(g_{j-1}))/\tau$ . However, since in Eqs. (67) or (69) the (discrete) time derivative is multiplied by  $\bar{n}(\mathcal{E})$ , we look for an estimate for

$$\mathcal{N}^j = \bar{n}(\mathcal{E}^j)\rho(g^j). \tag{120}$$

With the results of Proposition 3 and Lemma 1, the space-time compactness of  $\mathcal{N}$  will imply the same property of  $g$ .

We summarize our findings in the next proposition

**Proposition 4** *For given  $\tau > 0$  and  $j = 1, \dots, N$ , let  $(\mathbf{u}_\tau(t_j), g_\tau(t_j)) \in V$  denote a solution of Problem (PD). Then, we have*

$$\max_{1 \leq j \leq N} (\|\mathbf{u}_\tau(t_j)\|_{H^1(\Omega)^m} + \|g_\tau(t_j)\|_{L^2(\Omega)}) \leq C, \tag{121}$$

$$\tau \sum_{j=1}^N \int_{\Omega} \left( \lambda^* \left( \frac{\operatorname{div}(\mathbf{u}_{\tau}(t_j) - \mathbf{u}_{\tau}(t_{j-1}))}{\tau} \right)^2 + |\nabla g_{\tau}(t_j)|^2 \right) dx \leq C, \tag{122}$$

$$\tau \sum_{j=1}^N \|\varphi \operatorname{div} \mathbf{u}_{\tau}(t_j)\|_{H^1(\Omega)}^2 + \lambda^* \max_{1 \leq j \leq N} \|\varphi \operatorname{div} \mathbf{u}_{\tau}(t_j)\|_{H^1(\Omega)}^2 \leq C(\varphi), \tag{123}$$

$$\tau \sum_{j=1}^N \left( \left\| \frac{\mathcal{N}^j - \mathcal{N}^{j-1}}{\tau} \right\|_{H^{-2}(\Omega)}^2 + \|\varphi \mathcal{N}^j\|_{H^1(\Omega)}^2 \right) \leq C(\varphi), \tag{124}$$

where

$$\mathcal{N}^j = \bar{n}(\operatorname{div} \mathbf{u}_{\tau}(t_j))\rho(g_{\tau}(t_j))$$

and where  $\varphi \in C_0^{\infty}(\Omega)$ .

**Proof** We only need to prove estimate (124). Rewriting Eq. (79), we have

$$\begin{aligned} \int_{\Omega} \frac{\mathcal{N}^j - \mathcal{N}^{j-1}}{\tau} \psi \, dx &= \int_{\Omega} \frac{\bar{n}(\mathcal{E}^{j-1})(\rho(g^j) - \rho(g^{j-1}))}{\tau} \psi \, dx \\ &\quad + \int_{\Omega} \frac{(\bar{n}(\mathcal{E}^j) - \bar{n}(\mathcal{E}^{j-1}))\rho(g^j)}{\tau} \psi \, dx \\ &= \int_{\Omega} \frac{\bar{n}(\mathcal{E}^j) - \bar{n}(\mathcal{E}^{j-1})}{\tau} (\rho(g^j) - G(g^j)) \psi \, dx \\ &\quad - \int_{\Omega} \frac{(\mathcal{E}^j - \mathcal{E}^{j-1})(p(g^j) - p_0)}{\tau g^j} \psi \, dx \\ &\quad + \int_{\Omega} Q^j \psi \, dx - \int_{\Omega} k(\mathcal{E}^j) \mathcal{D}(g^j) \mathbb{K} \nabla g^j \nabla \psi \, dx, \quad \text{for } \psi \in H_0^2(\Omega). \end{aligned}$$

Recalling that for  $m \leq 3$ ,  $H^2(\Omega) \subset L^{\infty}(\Omega)$ , we have

$$\begin{aligned} &\left\| \frac{\mathcal{N}^j - \mathcal{N}^{j-1}}{\tau} \right\|_{H^{-2}(\Omega)}^2 \\ &\leq C \left( \left\| \frac{(\mathcal{E}^j - \mathcal{E}^{j-1})}{\tau} \right\|_{L^2(\Omega)}^2 \|g^j\|_{L^2(\Omega)}^2 + \|Q^j\|_{L^2(\Omega)}^2 + \|\nabla g^j\|_{L^2(\Omega)^m}^2 \right) \end{aligned}$$

and the full estimate reads

$$\begin{aligned} &\tau \sum_{j=1}^N \left\| \frac{\mathcal{N}^j - \mathcal{N}^{j-1}}{\tau} \right\|_{H^{-2}(\Omega)}^2 \\ &\leq C \left( \max_{1 \leq j \leq N} \|g^j\|_{L^2(\Omega)}^2 \tau \sum_{j=1}^N \left\| \frac{\mathcal{E}^j - \mathcal{E}^{j-1}}{\tau} \right\|_{L^2(\Omega)}^2 + 1 + \tau \sum_{j=1}^N \|\nabla g^j\|_{L^2(\Omega)^m}^2 \right) \leq C \end{aligned} \tag{125}$$

The local estimate for the space derivatives is given by

$$\begin{aligned} & \tau \sum_{j=1}^N \|\nabla(\varphi \mathcal{N}^j)\|_{L^{3/2}(\Omega)}^2 \\ & \leq C \left( \max_{1 \leq j \leq N} \|\nabla(\varphi \mathcal{E}^j)\|_{L^2(\Omega)}^2 \tau \sum_{j=1}^N \|g^j\|_{L^6(\Omega)}^2 + C + \tau \sum_{j=1}^N \|\nabla g^j\|_{L^2(\Omega)^m}^2 \right) \leq C. \end{aligned} \tag{126}$$

This results in estimate (124). □

### 4 Existence for Continuous Time Problem with $\lambda^* > 0$

In Proposition 4, where the time step  $\tau$  enters as a parameter, one finds  $\{(\mathbf{u}_\tau(t_j), g_\tau(t_j))\}_{j=1, \dots, N}$  from the “initial value”  $\text{div } \mathbf{u}(0) = 0$  and  $g(0) = g^0$ . Here  $N = O(1/\tau)$  and  $g^0 = g(\rho^0)$ . This procedure yields a time discretized approximation of the original quasi-static equations.

In this section, we investigate the limit  $\tau \searrow 0$ . Here a crucial role is played by the parameter  $\lambda^*$ , which is needed to control the behavior in time of  $\mathcal{E} = \text{div } \mathbf{u}$ .

Using the discrete solution  $(\mathbf{u}_\tau(t_j), g_\tau(t_j))$ , we construct two approximations that hold for all  $0 \leq t \leq T$ . The first is the piecewise constant approximation

$$(\bar{\mathbf{u}}_\tau(t), \bar{g}_\tau(t)) = (\mathbf{u}_\tau(t_j), g_\tau(t_j)) \quad \text{for } j\tau \leq t < (j+1)\tau. \tag{127}$$

The second is the Rothe interpolant, which is the piecewise linear time-continuous approximation

$$\begin{aligned} (\tilde{\mathbf{u}}_\tau(t), \tilde{g}_\tau(t)) &= \left(j + 1 - \frac{t}{\tau}\right) (\mathbf{u}_\tau(t_j), g_\tau(t_j)) + \left(\frac{t}{\tau} - j\right) (\mathbf{u}_\tau(t_{j+1}), g_\tau(t_{j+1})), \\ &\text{for } j\tau \leq t \leq (j+1)\tau. \end{aligned} \tag{128}$$

In (127) and (128), the index  $j$  runs from  $j = 0$  to  $j = N - 1$ .

Applying Proposition 4 yields for both approximations, with  $\natural \in \{-, \sim\}$ ,

$$\max_{0 \leq t \leq T} (\|\mathbf{u}_\tau^\natural(t)\|_{H^1(\Omega)^m}^2 + \|g_\tau^\natural(t)\|_{L^2(\Omega)}^2) dt \leq C, \tag{129}$$

$$\int_0^T \int_\Omega |\nabla g_\tau^\natural(t)|^2 dx dt \leq C, \tag{130}$$

$$\int_0^T \|\varphi \mathcal{E}_\tau^\natural(t)\|_{H^1(\Omega)}^2 dt \leq C, \tag{131}$$

$$\lambda^* \max_{0 \leq t \leq T} \|\varphi \mathcal{E}_\tau^\natural(t)\|_{H^1(\Omega)}^2 \leq C, \tag{132}$$

$$\int_0^T \|\varphi \mathcal{N}_\tau^\natural(t)\|_{W^{1,3/2}(\Omega)}^2 \leq C, \tag{133}$$

where  $\mathcal{E}_\tau^{\natural} = \operatorname{div} \mathbf{u}_\tau^{\natural}, \bar{\mathcal{N}}_\tau = \bar{n}(\bar{\mathcal{E}}_\tau)\rho(\bar{g}_\tau)$  and  $\tilde{\mathcal{N}}_\tau(t) = (j+1-t/\tau)\mathcal{N}^j + (t/\tau-j)\mathcal{N}^{j+1}$ . Further, we have

$$\begin{aligned} \partial_t \tilde{\mathcal{N}}_\tau &= \frac{\mathcal{N}^{j+1} - \mathcal{N}^j}{\tau} \quad \text{and} \quad \partial_t \tilde{\mathcal{E}}_\tau = \frac{\mathcal{E}^{j+1} - \mathcal{E}^j}{\tau}, \\ \text{for } t_j \leq t \leq t_{j+1} \quad \text{and} \quad j &= 0, \dots, N-1. \end{aligned}$$

Hence, by (122)

$$\int_0^T \int_\Omega \lambda^* |\partial_t \tilde{\mathcal{E}}_\tau(t)|^2 \, dx dt \leq C \tag{134}$$

and

$$\int_0^T \|\partial_t \tilde{\mathcal{N}}_\tau(t)\|_{H^{-2}(\Omega)}^2 \, dt \leq C. \tag{135}$$

In what follows, we rely heavily on the material and theory collected in Roubiček (2005, Chapters 7 and 8). Since the piecewise constant approximation  $(\bar{\mathbf{u}}_\tau(t), \bar{g}_\tau(t))$  is discontinuous in time, its time derivative is only a measure. To deal with this, we introduce the space  $\mathcal{M}(0, T; L^2(\Omega))$  of regular Borel measures in  $[0, T]$  with values in  $L^2(\Omega)$ , which is the dual space of  $C([0, T]; L^2(\Omega))$ . With  $\delta(t_j)$  denoting the Dirac measure concentrated in  $t_j$ , we have

$$\begin{aligned} \|\partial_t \bar{\mathcal{E}}_\tau\|_{\mathcal{M}(0,T;L^2(\Omega))} &= \left\| \sum_{j=1}^N (\mathcal{E}^j - \mathcal{E}^{j-1})\delta(t_j) \right\|_{\mathcal{M}(0,T;L^2(\Omega))} \\ &= \tau \sum_{j=1}^N \left\| \frac{\mathcal{E}^j - \mathcal{E}^{j-1}}{\tau} \right\|_{L^2(\Omega)} = \|\partial_t \tilde{\mathcal{E}}_\tau\|_{L^1(0,T;L^2(\Omega))} \\ &\leq \sqrt{T} \|\partial_t \tilde{\mathcal{E}}_\tau\|_{L^2(0,T;L^2(\Omega))} \leq C. \end{aligned} \tag{136}$$

Analogously

$$\|\partial_t \bar{\mathcal{N}}_\tau\|_{\mathcal{M}(0,T;H^{-2}(\Omega))} \leq C, \tag{137}$$

where  $\mathcal{M}(0, T; H^{-2}(\Omega))$  is the dual space of  $C([0, T]; H_0^2(\Omega))$ .

For the convergence of the time-continuous approximation (128), we use estimates (129)–(135) and the well-known weak and weak\* compactness theorems. The result is that there exists a quadruple  $\{\tilde{\mathbf{u}}, \tilde{g}, \tilde{\mathcal{E}}, \tilde{\mathcal{N}}\}$  such that along a subsequence  $\tau \searrow 0$  we have

$$\tilde{\mathbf{u}}_\tau \rightharpoonup \tilde{\mathbf{u}} \quad \text{weak* in } L^\infty(0, T; H_0^1(\Omega)^m), \tag{138}$$

$$\tilde{g}_\tau \rightharpoonup \tilde{g} \quad \text{weakly in } L^2(0, T; H^1(\Omega)), \tag{139}$$

$$\tilde{\mathcal{E}}_\tau \rightharpoonup \tilde{\mathcal{E}} \quad \text{weakly in } L^2(0, T; H^1(\omega)), \tag{140}$$



$$\partial_t \tilde{\mathcal{E}}_\tau \rightharpoonup \partial_t \tilde{\mathcal{E}} \text{ weakly in } L^2(0, T; L^2(\Omega)), \tag{141}$$

$$\tilde{\mathcal{N}}_\tau \rightharpoonup \tilde{\mathcal{N}} \text{ weakly in } L^2(0, T; W^{1,3/2}(\omega)), \tag{142}$$

$$\partial_t \tilde{\mathcal{N}}_\tau \rightharpoonup \partial_t \tilde{\mathcal{N}} \text{ weakly in } L^2(0, T; H^{-2}(\Omega)). \tag{143}$$

Concerning the convergence of  $(\bar{\mathbf{u}}_\tau, \bar{g}_\tau)$ , we use estimates (129)–(133), now combined with (136)–(137). Moreover, applying (Roubiřek 2005, Corollary 7.9), we use that the spaces

$$W^{1,2,\mathcal{M}}(0, T; H^1(\omega), L^2(\omega)) = \{z \in L^2(0, T; H^1(\omega)) \mid \frac{dz}{dt} \in \mathcal{M}(0, T; L^2(\omega))\}$$

and  $W^{1,2,\mathcal{M}}(0, T; W^{1,3/2}(\omega), H^{-2}(\omega))$  are compactly embedded in  $L^2(0, T; L^2(\omega))$ , for any smooth bounded subset  $\omega$  of  $\Omega$ . The result is that there exists  $(\bar{\mathbf{u}}, \bar{g}, \bar{\mathcal{E}}, \bar{\mathcal{N}})$  such that along a subsequence  $\tau \searrow 0$  one has the same convergence as in (138)–(140) and (142). The convergence in (141) and (143) is now replaced by weak–\* convergence in  $\mathcal{M}(0, T; L^2(\Omega))$  for  $\partial_t \tilde{\mathcal{E}}_\tau$  and in  $\mathcal{M}(0, T; H^{-2}(\Omega))$  for  $\partial_t \tilde{\mathcal{N}}_\tau$ .

Furthermore, the estimates allow us to conclude

$$\bar{\mathcal{E}}_\tau \rightarrow \bar{\mathcal{E}} \text{ strongly in } L^2((0, T) \times \omega) \text{ and (a.e) on } (0, T) \times \omega, \tag{144}$$

$$\bar{\mathcal{N}}_\tau \rightarrow \bar{\mathcal{N}} \text{ strongly in } L^2((0, T) \times \omega) \text{ and (a.e) on } (0, T) \times \omega. \tag{145}$$

As a consequence,

$$\rho(\bar{g}_\tau) = \frac{\bar{\mathcal{N}}_\tau}{\bar{n}(\bar{\mathcal{E}}_\tau)} \rightarrow \frac{\bar{\mathcal{N}}}{\bar{n}(\bar{\mathcal{E}})} \tag{146}$$

and

$$\bar{g}_\tau = \rho^{-1}\left(\frac{\bar{\mathcal{N}}_\tau}{\bar{n}(\bar{\mathcal{E}}_\tau)}\right) \rightarrow \rho^{-1}\left(\frac{\bar{\mathcal{N}}}{\bar{n}(\bar{\mathcal{E}})}\right) = \bar{g}. \tag{147}$$

strongly in  $L^2((0, T) \times \omega)$  and a.e. on  $(0, T) \times \omega$ . This in turn implies

$$\begin{cases} \rho(\bar{g}_\tau) \rightarrow \rho(\bar{g}) \text{ strongly in } L^2((0, T) \times \omega) \text{ and (a.e) on } (0, T) \times \omega; \\ \mathcal{D}(\bar{g}_\tau) \rightarrow \mathcal{D}(\bar{g}) \text{ strongly in } L^2((0, T) \times \omega) \text{ and (a.e) on } (0, T) \times \omega. \end{cases} \tag{148}$$

Inherited from  $\bar{\mathcal{E}}_\tau = \operatorname{div} \bar{\mathbf{u}}_\tau$ , the convergence properties imply

$$\bar{\mathcal{E}} = \operatorname{div} \bar{\mathbf{u}} \text{ a.e. in } (0, T) \times \Omega. \tag{149}$$

As in (Roubiček 2005, pages 224–226), one shows that  $\tilde{\mathbf{u}} = \bar{\mathbf{u}}$  and  $\tilde{g} = \bar{g}$ . Then, (149) implies that  $\bar{\mathcal{E}} = \tilde{\mathcal{E}}$ . Alternatively, this follows from estimate (134) which gives

$$\begin{aligned} \int_0^T \|\bar{\mathcal{E}}_\tau(t) - \tilde{\mathcal{E}}_\tau(t)\|_{L^2(\Omega)}^2 dt &= \frac{\tau^3}{2} \sum_{j=1}^{N-1} \|\mathcal{E}^j - \mathcal{E}^{j+1}\|_{L^2(\Omega)}^2 \\ &= C\tau^2 \|\partial_t \tilde{\mathcal{E}}_\tau\|_{L^2(0,T;L^2(\Omega))}^2 = C\tau^2. \end{aligned} \tag{150}$$

Similarly,

$$\int_0^T \|\bar{\mathcal{N}}_\tau(t) - \tilde{\mathcal{N}}_\tau(t)\|_{H^{-2}(\Omega)}^2 dt = C\tau^2, \tag{151}$$

which yields  $\bar{\mathcal{N}} = \tilde{\mathcal{N}}$ .

From this point on, we denote the limit, as  $\tau \searrow 0$ , by the quadruple  $(\mathbf{u}, g, \mathcal{E}, \mathcal{N})$ , where

$$\mathcal{E} = \operatorname{div} \mathbf{u} \quad \text{and} \quad \mathcal{N} = \bar{n}(\mathcal{E})\rho(g).$$

We are now in a position to prove the main existence result for a weak solution of the time-continuous case.

**Theorem 3** *Let  $\lambda^* > 0$ . Then, there exists at least one weak free energy solution  $(\mathbf{u}, \mathcal{E}, \rho)$  satisfying Definition 1.*

**Proof** In the proof, we use approximations (127) and (128), and their convergence properties.

Let  $\tau > 0$ , sufficiently small, and let  $t \in (\tau, T)$ . Then,  $t_j \leq t < t_{j+1}$  for some  $j \in \{1, \dots, N - 1\}$  and  $\bar{\mathbf{u}}_\tau(t) = \mathbf{u}^j$  and  $\bar{g}_\tau(t) = g^j$ .

We first consider the momentum balance Eq. (80).

The starting point is Problem (PD). Using Eq. (108), we have for any  $\xi \in H_0^1(\Omega)^m$

$$\begin{aligned} \int_\Omega \mathcal{G}e(\bar{\mathbf{u}}_\tau) : e(\xi) dx &= \int_\Omega \mathcal{G}e(\mathbf{u}^j) : e(\xi) dx \\ &= -\frac{\lambda^*}{\tau} \int_\Omega (\mathcal{E}^j - \mathcal{E}^{j-1}) \operatorname{div} \xi dx + \int_\Omega p(g^j) \operatorname{div} \xi dx + \int_\Omega \mathbf{F}^j \cdot \xi dx \\ &= -\lambda^* \int_\Omega \partial_t \tilde{\mathcal{E}}_\tau(t - \tau) \operatorname{div} \xi dx + \int_\Omega p(\bar{g}_\tau) \operatorname{div} \xi dx + \int_\Omega \mathbf{F}_\tau \cdot \xi dx. \end{aligned} \tag{152}$$

Here, we introduced

$$\mathbf{F}_\tau(t) = \mathbf{F}(t_j) = \mathbf{F}^j \quad \text{for } t_j \leq t < t_{j+1} \quad \text{and } j = 0, \dots, N - 1.$$

Multiplying Eq. (152) by  $\alpha \in C_0^\infty(0, T)$  and integrating the result over  $(\tau, T)$ , yields

$$\int_\tau^T \left\{ \int_\Omega \mathcal{G}e(\bar{\mathbf{u}}_\tau) : e(\xi) dx \right\} \alpha(t) dt + \lambda^* \int_\tau^T \left\{ \int_\Omega \partial_t \tilde{\mathcal{E}}_\tau(t - \tau) \operatorname{div} \xi dx \right\} \alpha(t) dt$$

$$+ \int_{\tau}^T \left\{ \int_{\Omega} p(\bar{g}_{\tau}) \operatorname{div} \xi, dx \right\} \alpha(t) dt = \int_{\tau}^T \left\{ \int_{\Omega} \mathbf{F}_{\tau} \cdot \xi, dx \right\} \alpha(t) dt. \tag{153}$$

Next, we send  $\tau \searrow 0$  along the appropriate subsequence to have convergence of the terms containing  $\mathbf{u}$ ,  $\mathcal{E}$  and  $\mathbf{F}$ . What remains is the pressure term. We recall that  $p(g)$  is the composite function  $(p \circ \rho)(g)$ , where  $p(\rho)$  is given by (36) and  $\rho(g)$  is defined through (53) and (55). Since  $\bar{g}_{\tau} \rightarrow g$  strongly in  $L^2((0, T) \times \omega)$ , see (147), we have similarly

$$p(\bar{g}_{\tau}) = (p \circ \rho)(\bar{g}_{\tau}) \rightarrow (p \circ \rho)(\bar{g}) = (p \circ \rho)(g) = p(g)$$

strongly in  $L^2((0, T) \times \omega)$  and a.e. in  $(0, T) \times \omega$ .

This concludes the first part of the proof.

Next, we tackle the mass balance Eq. (69).

We first put Eq. (77) in the form

$$\partial_t \mathcal{N} - \partial_t \bar{n}(\mathcal{E}) \rho(g) + D(g, \mathcal{E}) \partial_t \mathcal{E} - \operatorname{div} \left( k(\mathcal{E}) \mathcal{D}(g) \mathbb{K} \nabla g \right) = Q$$

and apply the discretization of Problem (PD). Similarly to (153) this gives for any  $\psi \in C_0^{\infty}(\Omega)$  and  $\alpha \in C^{\infty}[0, T]$

$$\begin{aligned} & \int_{\tau}^T \int_{\Omega} \left( \partial_t \tilde{\mathcal{N}}_{\tau}(t - \tau) - \partial_t \tilde{v}_{\tau}(t - \tau) \left( \rho(\bar{g}_{\tau}) - \frac{G(\bar{g}_{\tau})}{\bar{g}_{\tau}} \right) + \partial_t \tilde{\mathcal{E}}_{\tau}(t - \tau) \frac{p(\bar{g}_{\tau}) - p_0}{\bar{g}_{\tau}} \right) \\ & \psi(x) \alpha(t) dx dt + \\ & \int_{\tau}^T \int_{\Omega} k(\bar{\mathcal{E}}_{\tau}) \mathcal{D}(\bar{g}_{\tau}) \mathbb{K} \nabla \bar{g}_{\tau} \cdot \nabla \psi(x) \alpha(t) dx dt = \int_{\tau}^T \int_{\Omega} Q_{\tau} \psi(x) \alpha(t) dx dt, \end{aligned} \tag{154}$$

where

$$\begin{aligned} \tilde{v}_{\tau}(t) &= \left( j + 1 - \frac{t}{\tau} \right) \bar{n}(\mathcal{E}^j) + \left( \frac{t}{\tau} - j \right) \bar{n}(\mathcal{E}_{j+1}), \\ &\text{and} \\ Q_{\tau}(t) &= Q(t_j) = Q^j \\ &\text{for } j\tau \leq t < (j + 1)\tau. \end{aligned} \tag{155}$$

The boundedness of  $\bar{n}'$  implies

$$\| \partial_t \tilde{v}_{\tau} \|_{L^2((0, T) \times \Omega)} \leq C \tag{156}$$

and inherited from (131)

$$\| \tilde{v}_{\tau} \|_{L^2(0, T; H^1(\omega))} \leq C. \tag{157}$$

Hence,

$$\partial_t \tilde{v}_\tau \rightharpoonup \partial_t \bar{n}(\mathcal{E}) \text{ weakly in } L^2((0, T) \times \Omega) \tag{158}$$

and

$$\tilde{v}_\tau \rightarrow \bar{n}(\mathcal{E}) \text{ strongly in } L^2((0, T) \times \omega) \text{ and a.e. in } (0, T) \times \omega. \tag{159}$$

We are now in position to pass to the limit  $\tau \searrow 0$  in (154) and obtain

$$\begin{aligned} & \int_0^T \langle \partial_t (\bar{n}(\mathcal{E})\rho(g)), \psi \rangle_{H^{-2}(\Omega), H_0^2(\Omega)} \alpha(t) \, dt \\ & - \int_0^T \int_\Omega \left( \partial_t \bar{n}(\mathcal{E})(\rho(g) - \frac{G(g)}{g}) + \partial_t \mathcal{E} \frac{p(g) - p_0}{g} \right) \psi(x) \alpha(t) \, dx \, dt \\ & + \int_0^T \int_\Omega k(\mathcal{E}) \mathcal{D}(g) \mathbb{K} \nabla \bar{g} \cdot \nabla \psi(x) \alpha(t) \, dx \, dt = \int_0^T \int_\Omega Q \psi(x) \alpha(t) \, dx \, dt \end{aligned} \tag{160}$$

or

$$\begin{aligned} & \partial_t \left( \bar{n}(\mathcal{E})\rho(g) \right) - \partial_t \bar{n}(\mathcal{E})\rho(g) + D(\rho, \mathcal{E})\partial_t \mathcal{E} - \operatorname{div} \left( k(\mathcal{E})\mathcal{D}(g)\mathbb{K}\nabla g \right) \\ & = Q \text{ in } \mathcal{D}'((0, T) \times \Omega), \end{aligned} \tag{161}$$

It remains to check the initial and boundary conditions and the energy inequality (72). First, we notice that (140)–(141) imply

$$\tilde{\mathcal{E}}_\tau \rightharpoonup \mathcal{E} \text{ weakly in } W^{1,2,2}(0, T; H^1(\omega), L^2(\omega)), \tag{162}$$

where  $W^{1,2,2}(0, T; H^1(\omega), L^2(\omega)) = \{z \in L^2(0, T; H^1(\omega)) \mid \partial_t z \in L^2(0, T; L^2(\omega))\}$ . In this space, the trace in time  $\mathcal{E} \rightarrow \mathcal{E}(0)$  is a weakly continuous map from  $W^{1,2,2}(0, T; H^1(\omega), L^2(\omega))$  to  $L^2(\omega)$ . Hence,

$$\tilde{\mathcal{E}}_\tau(0) \rightharpoonup \mathcal{E}(0) \text{ weakly in } L^2(\omega), \tag{163}$$

where  $\mathcal{E}(0) = \operatorname{div} \mathbf{u}^0 = 0$ .

Next, using (142), (143) and (145), we conclude that

$$\tilde{\mathcal{N}}_\tau(0) \rightharpoonup \mathcal{N}(0) \text{ weakly in } H^{-2}(\omega), \tag{164}$$

which justifies the initial condition for  $\mathcal{N}$ . Since  $\mathcal{N} = \bar{n}(\mathcal{E})\rho(g)$ , we have simultaneously the initial conditions for the density  $\rho$  and for  $g$ .

We still miss the flux boundary condition for the mass balance equation (67). The starting point is again Eq. (154), now with  $\psi \in H^1(\Omega)$  and  $\alpha(T) = 0$ . Since

$$\int_\tau^T \int_\Omega \partial_t \tilde{\mathcal{N}}_\tau(x, t - \tau) \psi(x) \alpha(t) \, dx \, dt = - \int_\tau^T \int_\Omega \tilde{\mathcal{N}}_\tau(x, t - \tau) \psi(x) \frac{d}{dt} \alpha(t) \, dx \, dt$$

$$\begin{aligned}
 & - \int_{\Omega} \tilde{\mathcal{N}}_{\tau}(x, 0) \psi(x) \alpha(\tau) \, dx \rightarrow - \int_0^T \int_{\Omega} \tilde{\mathcal{N}}(x, t) \psi(x) \frac{d}{dt} \alpha(t) \, dx dt \\
 & - \alpha(0) \int_{\Omega} \bar{n}(0) \rho^0(x) \psi(x) \, dx
 \end{aligned}$$

and since the strong convergence (144)–(145), together with the weak convergence (140) and (142), implies the same for  $\Omega$ , we may pass to the limit  $\tau \searrow 0$  and conclude that

$$\begin{aligned}
 & - \int_0^T \int_{\Omega} \rho(g) \bar{n}(\mathcal{E}) \partial_t \Phi(x, t) \, dx dt - \int_{\Omega} \bar{n}(0) \rho^0(x) \Phi(x, 0) \, dx \\
 & + \int_0^T \int_{\Omega} \partial_t \mathcal{E} \left( D(\rho(g), \mathcal{E}) \right. \\
 & \left. - \rho(g) \bar{n}'(\mathcal{E}) \right) \Phi(x, t) \, dx dt + \int_0^T \int_{\Omega} k(\mathcal{E}) \mathcal{D}(g) \mathbb{K} \nabla g \cdot \nabla_x \Phi(x, t) \, dx dt \\
 & = \int_0^T \int_{\Omega} Q \Phi(x, t) \, dx dt, \quad \forall \Phi \in H^1(\Omega \times (0, T)), \quad \text{with } \Phi|_{t=T} = 0.
 \end{aligned}$$

To show inequality (72), we follow again (Roubiřek 2005, pp. 223–226). Starting point is inequality (107), which we write for any  $K \in \mathbb{N}$ ,  $K \leq N - 1$ , and for any time step  $\tau$  as

$$\begin{aligned}
 & J^{K+1} - J^K + \tau \int_{\Omega} \left\{ \lambda^* \left( \frac{\mathcal{E}^{K+1} - \mathcal{E}^K}{\tau} \right)^2 + k(\mathcal{E}^{K+1}) \mathcal{D}(g^{K+1}) \mathbb{K} \nabla g^{K+1} \cdot \nabla g^{K+1} \right. \\
 & \left. + \frac{\mathbf{F}^{K+1} - \mathbf{F}^K}{\tau} \cdot \mathbf{u}^K - Q^{K+1} g^{K+1} \right\} dx \leq 0. \tag{165}
 \end{aligned}$$

With the notation from (127) and (128), we rewrite inequality (165):

$$\begin{aligned}
 & \frac{d}{dt} \tilde{J}_{\tau}(t) + \int_{\Omega} \left\{ \lambda^* (\partial_t \tilde{\mathcal{E}}_{\tau}(t))^2 + k(\bar{\mathcal{E}}_{\tau}(t + \tau)) \mathcal{D}(\bar{g}_{\tau}(t + \tau)) \mathbb{K} \nabla \bar{g}_{\tau}(t + \tau) \cdot \nabla \bar{g}_{\tau}(t + \tau) \right. \\
 & \left. + \partial_t \tilde{\mathbf{F}}_{\tau}(t) \cdot \bar{\mathbf{u}}(t) - \bar{Q}_{\tau}(t + \tau) \bar{g}_{\tau}(t + \tau) \right\} dx \leq 0,
 \end{aligned}$$

for  $K\tau \leq t \leq (K + 1)\tau$  and  $0 \leq K \leq N - 1$ .

Integrating this inequality from  $t = t_1$  to  $t = t_2$ , with  $0 \leq t_1 \leq t_2 \leq T$ , using the convergence results from (138)–(145) and the weak lower semi-continuity of the gradient  $\nabla \bar{g}_{\tau}$  in  $L^2(Q_T)$ , yields the desired result.  $\square$

### 5 Discussion and Conclusion

In this paper, we study a model that describes the quasi-static mechanical behavior of a fluid saturated porous medium. In its simplest (linear) form, it is described by

Eqs. (1)–(4), where (1) results from the fluid phase mass balance in the case that the fluid is incompressible.

We follow Rutqvist et al. (2001) and Lewis and Schrefler (1998) and propose a fluid mass balance that is based on the mixture theory of Bedford and Drumheller (1978) and Bedford and Drumheller (1983). This yields Eq. (6) and the resulting nonlinear system is given by (2)–(4) and (6). Note that the time derivative of the fluid density  $\rho$  appears in (6), since the fluid is assumed weakly compressible, see expression (10). Models where the fluid density is constant (see Bociu et al. 2016; Cao et al. 2013) do not contain this source term. Moreover, the porosity  $n$  and the deformation of the medium are related through (8). An expression for this relation is derived from the solid phase mass balance. It is given by (19) or, when the deformation is small, by approximation (21).

It is shown by means of a counterexample that the porosity may admit non-physical, i.e., negative, values. This is made precise in Proposition 1. To obtain a well-posed mathematical problem, the porosity is modified according to cutoff (35). This cutoff is chosen such that it reduces to the correct expression in the physical range. Outside this range, it remains positive. Likewise, a cutoff for the density is introduced through expressions (60) and (64).

The momentum balance Eq. (2)–(4) is modified as well. Following Murad and Cushman (1996), we add the term

$$\lambda^* \operatorname{div} \partial_t \mathbf{u} \quad (\lambda^* \geq 0) \quad (166)$$

to the expression for the total stress. This results in expression (5). Murad and Cushman give a thermodynamically based derivation of the equation in which (166) appears as the difference between the fluid and solid pressures. Having  $\lambda^* > 0$ , (166) acts as a time regularization of the volumetric stress for our quasi-static problem.

An important role in the analysis of the equations is played by the free energy of the system. This free energy acts as a Lyapunov functional. It is given by (66), which generalizes Biot's original expression developed for the linear case (Biot 1962). In the case that the deformation and fluid density are in the physical range, the free energy simplifies to, see also (56),

$$J(\mathbf{u}, \rho) = \int_{\Omega} \left( \frac{1}{2} \mathcal{G}e(\mathbf{u}) : e(\mathbf{u}) - \mathbf{F} \cdot \mathbf{u} + \frac{n(\operatorname{div} \mathbf{u})}{\beta_0 n_0 (1 - n_0) \rho_0} \left( (1 - n_0) \rho - \rho_0^{n_0} \rho^{1-n_0} + n_0 \rho_0 \right) \right) dx. \quad (167)$$

We introduce a weak formulation and prove existence of a solution in a number of steps. Discretizing in time, we first consider the incremental equations. Using Brézis' fundamental theorem for pseudo-monotone operators, see for instance Lions (1969) and Roubiček (2005), we obtain existence for the corresponding incremental problem. The result holds for any  $\lambda^* \geq 0$ . Moreover, using the free energy, estimates that are global in time are derived. These (stability) estimates are crucial when considering the time-continuous, quasi-static, formulation for which we prove existence at the expense of having  $\lambda^* > 0$ . The free energy implies global stability of the solution.

The quasi-static case with  $\lambda^* = 0$  fails to have the time estimate from inequality (111). Therefore, we are not able to extend the results to the quasi-static case. Furthermore, the convergence results are based on compactness arguments. Hence, also uniqueness remains an open problem. Clearly, we have existence (globally in time) for the time-discrete quasi-static case, which is relevant for numerical purposes.

We note that only in the proof of the local  $H^1(\Omega)$ – estimates for  $\operatorname{div} \mathbf{u}$ , we use the fact that the Gassmann tensor has the specific form of Hooke’s law (4). In the incremental problem we could have replaced  $\mathcal{G}$  by a general rank-4, symmetric, positive-definite Gassmann tensor.

Some particular cases of system (24)–(25) were studied before. An interesting example is the consolidation with an irrotational composite flow rate, when the system reduces to a scalar pseudo-parabolic PDE. For details see Holland and Showalter (2018).

We notice also that the model studied in this paper was extensively used by Schrefler et al., see Schrefler et al. (1990) and Lewis and Schrefler (1998) and references therein. It is broadly accepted in the computational poromechanics community. A review of different numerical methods and software is given in Rutqvist et al. (2001) and Minkoff et al. (2003).

**Acknowledgements** The first author expresses his gratitude to Tea Mikelić for retrieving the missing files from her father’s computer. Further, he wants to thank the colleagues from Hasselt University, Sorin Pop and Koondanibha Mitra, for their assistance with the revision.

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article’s Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article’s Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

## References

- Auriault, J.-L., Sanchez-Palencia, E.: Étude du compartement macroscopique d’un milieu poreux saturé déformable. *J. Mécanique* **16**, 575–603 (1977)
- Bear, J., Bachmat, Y.: *Introduction to Modeling of Transport Phenomena in Porous Media*. Springer Science and Business Media, New York (1990)
- Bedford, A., Drumheller, D.S.: A variational theory of immiscible mixtures. *Arch. Ration. Mech. Anal.* **68**, 37–51 (1978)
- Bedford, A., Drumheller, D.S.: Theories of immiscible and structured mixtures. *Int. J. Eng. Sci.* **21**, 863–960 (1983)
- Biot, M.A.: Mechanics of deformation and acoustic propagation in porous media. *J. Appl. Phys.* **33**, 1482 (1962)
- Bociu, L., Guidoboni, G., Sacco, R., Webster, J.T.: Analysis of nonlinear poro-elastic and poro-visco-elastic models. *Arch. Ration. Mech. Anal.* **222**, 1445–1519 (2016)
- Bosco, E., Peerlings, R.H., Geers, M.G.: Predicting hygro-elastic properties of paper sheets based on an idealized model of the underlying fibrous network. *Int. J. Solids Struct.* **56**, 43–52 (2015)
- Čanić, S., Hartley, C.J., Rosenstrauch, D., Tambača, J., Guidoboni, G., Mikelić, A.: Blood flow in compliant arteries: an effective viscoelastic reduced model, numerics and experimental validation. *Ann. Biomed. Eng.* **34**, 575–592 (2006)

- Cao, X., Pop, I.S.: Degenerate two-phase porous media flow model with dynamic capillarity. *J. Differ. Equ.* **260**, 2418–2456 (2016)
- Cao, Y., Chen, S., Meir, A.J.: Analysis and numerical approximations of equations of nonlinear poroelasticity. *Discrete Contin. Dyn. Syst. Ser. B* **18**, 1253–1273 (2013)
- Cardoso, L., Fritton, S.P., Gailani, G., Benalla, M., Cowin, S.C.: Advances in assessment of bone porosity, permeability, and interstitial fluid flow. *J. Biomech.* **46**, 253–265 (2013)
- Coussy, O.: *Poromechanics*. Wiley (2004)
- Cowin, S.C.: Bone poroelasticity. *J. Biomech.* **32**, 218–238 (1999)
- Evans, L.C.: A survey of entropy methods for partial differential equations. *Bull. Am. Math. Soc.* **41**, 409–438 (2004)
- Holland, E., Showalter, R.E.: Poro-visco-elastic compaction in sedimentary basins. *SIAM J. Math. Anal.* **50**, 2295–2316 (2018)
- Jüngel, A.: *Entropy Methods for Diffusive Partial Differential Equations*. BCAM Springer Briefs. Springer (2016)
- Lewis, R.W., Schrefler, B.A.: *The Finite Element Method in the Static and Dynamic Deformation and Consolidation of Porous Media*. Wiley (1998)
- Lions, J.L.: *Quelques méthodes de résolution des problèmes aux limites non linéaires*. Dunod, Paris (1969)
- Marciniak-Czochra, A., Mikelić, A.: A rigorous derivation of the equations for the clamped Biot–Kirchhoff–Love poroelastic plate. *Arch. Ration. Mech. Anal.* **215**, 1035–1062 (2015)
- Mei, C.C., Vernescu, B.: *Homogenization Methods for Multiscale Mechanics*. World scientific (2010)
- Mikelić, A.: A global existence result for the equations describing unsaturated flow in porous media with dynamic capillary pressure. *J. Differ. Equ.* **248**, 1561–1577 (2010)
- Mikelić, A., Tambača, J.: Derivation of a poroelastic flexural shell model. *Multiscale Model. Simul.* **14**, 364–397 (2016)
- Mikelić, A., Wheeler, M.F.: On the interface law between a deformable porous medium containing a viscous fluid and an elastic body. *Math. Models Methods Appl. Sci.* **22**, 1250031 (2012)
- Milišić, J.P.: The unsaturated flow in porous media with dynamic capillary pressure. *J. Differ. Equ.* **264**, 5629–5658 (2018)
- Minkoff, S.E., Stone, C.M., Bryant, S., Peszynska, M., Wheeler, M.F.: Coupled fluid flow and geomechanical deformation modeling. *J. Pet. Sci. Eng.* **38**, 37–56 (2003)
- Murad, M.A., Cushman, J.H.: Multiscale flow and deformation in hydrophilic swelling porous media. *Int. J. Eng. Sci.* **34**, 313–338 (1996)
- Owczarek, S.: A Galerkin method for Biot consolidation model. *Math. Mech. Solids* **15**, 42–56 (2010)
- Prosi, M., Zunino, P., Perktold, K., Quarteroni, A.: Mathematical and numerical models for transfer of low-density lipoproteins through the arterial walls: a new methodology for the model set up with applications to the study of disturbed luminal flow. *J. Biomech.* **38**, 903–917 (2005)
- Roubíček, T.: *Nonlinear Partial Differential Equations with Applications*. Springer (2005)
- Rutqvist, J., Börgesson, L., Chijimatsu, M., Kobayashi, A., Jing, L., Nguyen, T.S., Noorishad, J., Tsang, C.F.: Thermohydromechanics of partially saturated geological media: governing equations and formulation of four finite element models. *Int. J. Rock Mech. Min. Sci.* **38**, 105–127 (2001)
- Sanchez-Palencia, E.: *Non-homogeneous Media and Vibration Theory*, Springer Lecture Notes in Physics, vol. 129. Springer (1980)
- Schrefler, B.A., Simoni, L., Xikui, L., Zienkiewicz, O.C.: Mechanics of partially saturated porous media. In: Desai, D.S., Gioda, G. (eds.) *Numerical Methods and Constitutive Modelling in Geomechanics, Courses and Lectures CISM no. 311*, pp. 169–211. Springer, New York (1990)
- Schweizer, B.: *Partielle Differentialgleichungen*. Springer, Berlin (2018)
- Showalter, R.E.: Diffusion in poro-elastic media. *J. Math. Anal. Appl.* **251**, 310–340 (2000)
- Terzaghi, K.: *Theoretical Soil Mechanics*. Chapman and Hall Limited, London (1951)
- Tolstoy, I. (ed.): *Acoustics, Elasticity, and Thermodynamics of Porous Media*. Twenty-One Papers by M.A. Biot. Acoustical Society of America, New York (1992)
- van Duijn, C.J., Mikelić, A., Wick, T.: A monolithic phase-field model of a fluid-driven fracture in a nonlinear poroelastic medium. *Math. Mech. Solids* **24**, 1530–1555 (2019)
- Verruijt, A.: Theory and problems of poroelasticity (2015). <https://geotechpedia.com>
- Zeidler, E.: *Nonlinear Functional Analysis and Its Applications II/B Nonlinear Monotone Operators*. Springer, New York (1990)
- Ženšek, A.: The existence and uniqueness theorem in Biot’s consolidation theory. *Appl. Math.* **29**, 194–211 (1984)



**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.