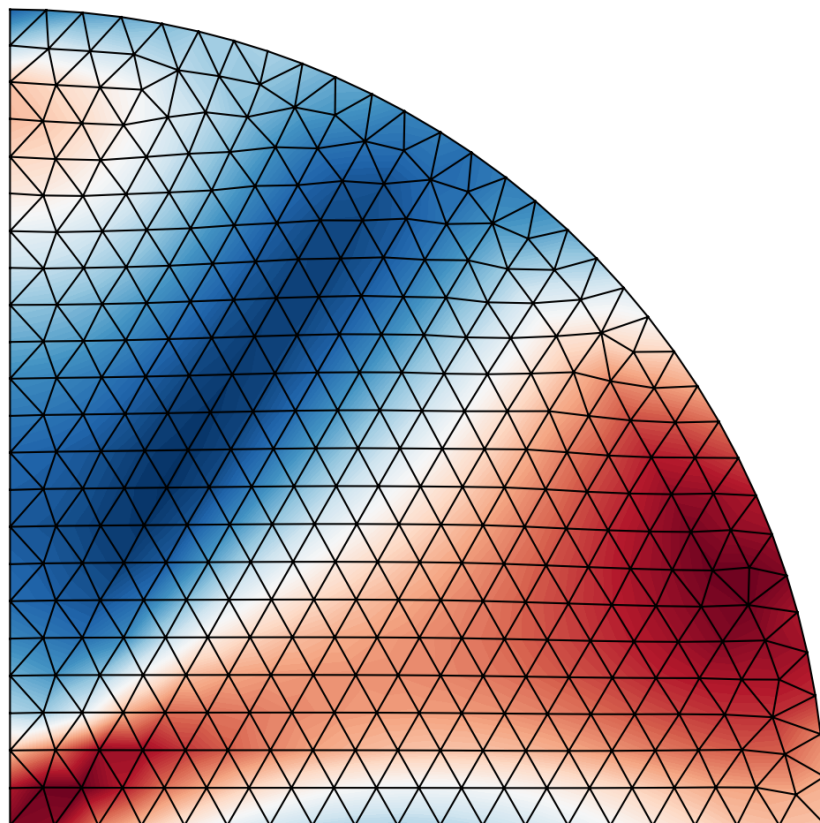


# Analysis of Stabilized Finite Element Methods for a Morpho-Poroelastic Model Applied to Tumor Growth

Duncan den Bakker



# Analysis of Stabilized Finite Element Methods for a Morpho-Poroelastic Model Applied to Tumor Growth

by

Duncan den Bakker

to obtain the degree of Master of Science  
at the Delft University of Technology,  
to be defended publicly on Monday July 27, 2020 at 14:00.

Student number: 4476166  
Project duration: December 1, 2020 – July 27, 2020  
Thesis committee: Prof. dr. ir. F.J. Vermolen, TU Delft, supervisor  
Dr. N.V. Budko, TU Delft  
Dr. J.L.A. Dubbeldam, TU Delft

## Abstract

The theory of morpho-poroelasticity is applied to tumor growth. This allows for the modeling of permanent deformation in the tissue as a result of the presence of tumor cells. The work done in this project can be divided into two parts. In the first part we review existing models for elasticity and build corresponding finite element models. The aim of the first part is to gain understanding in the behaviour of morphoelasticity and poroelasticity. In morphoelasticity the deformation tensor is decomposed into an elastic and plastic component, allowing for permanent deformations to occur in the tissue. The theory of morphoelasticity describes the interplay between a fluid in a porous tissue and the tissue's elastic properties. A common problem in poroelastic models is the occurrence of non-physical oscillations in the pressure. The second part of this project contains novel work. First, a rigorous mathematical derivation is presented of the tuning parameter found in the diffusive stabilization method for poroelastic systems. Secondly, we present a finite element model for the novel combination of morpho- and poroelasticity. The derivation concerning the diffusive stabilization is also applicable to the morpho-poroelastic finite element model. We obtain a promising tumor growth model by combining a biochemical tumor growth model with morpho-poroelasticity. This biochemical model is based on a nutrient transport equation and a tumor cell density evolution equation. Some example output of the tumor growth finite element model is shown. It can be used as a starting point for more elaborate tumor growth models.

# Contents

List of Symbols . . . . .	iv
<b>1 Introduction</b>	<b>1</b>
1.1 Thesis Outline . . . . .	2
1.2 A Note on Implementation . . . . .	2
<b>Part I</b>	<b>4</b>
<b>2 A Sequence of Models for Elasticity</b>	<b>5</b>
2.1 Stress and Strain . . . . .	5
2.2 Pure Elasticity . . . . .	6
2.2.1 Weak Formulation . . . . .	7
2.2.2 Galerkin Equations . . . . .	7
2.2.3 Numerical Experiments . . . . .	8
2.3 Viscoelasticity . . . . .	10
2.3.1 Eulerian Versus Lagrangian Coordinates . . . . .	10
2.3.2 Complete Viscoelastic Model Equations . . . . .	11
2.3.3 Weak Formulation . . . . .	12
2.3.4 Galerkin Equations . . . . .	14
2.3.5 Time Discretization . . . . .	15
2.3.6 Numerical Experiments . . . . .	17
2.4 Morphoelasticity . . . . .	18
2.4.1 Weak Formulation . . . . .	19
2.4.2 Galerkin Equations . . . . .	21
2.4.3 Time Discretization . . . . .	23
2.4.4 Numerical Experiments . . . . .	24
2.5 Comparison Between Elastic Models . . . . .	26
<b>3 The Mechanics of Porous Media</b>	<b>29</b>
3.1 Porous Media . . . . .	29
3.2 Terzaghi Problem . . . . .	30
3.2.1 Linear-Linear Elements . . . . .	31
3.2.2 Quadratic-Linear Elements . . . . .	32
3.2.3 Stabilization . . . . .	34
3.3 Poroelasticity on a Fixed Domain . . . . .	35
3.3.1 Linear-Linear Elements . . . . .	37
3.3.2 Stabilization with MINI Element . . . . .	38
3.3.3 Diffusive Stabilization . . . . .	41
3.3.4 Comparison Between Bubble- and Diffusive Stabilization . . . . .	41
3.4 Visco-Poroelasticity on a Deforming Domain . . . . .	43
3.4.1 Derivation of System of Equations . . . . .	43
3.4.2 Stabilization . . . . .	44
3.4.3 Numerical Experiments . . . . .	44
3.5 Comparison Between Viscoelastic and Visco-Poroelastic Models . . . . .	46

<b>Part II</b>	<b>48</b>
<b>4 Morpho-Poroelasticity and Analysis of Stabilized Finite Element Methods</b>	<b>49</b>
4.1 Stabilizing Poroelastic Systems . . . . .	49
4.2 Analysis of One-Dimensional Visco-Poroelastic System . . . . .	50
4.3 Approximating the Inverse of the Discrete Reaction-Laplacian Operator . . . . .	51
4.3.1 Existence and Uniqueness of Weak Solution to Reaction-Laplacian Equation . . . . .	51
4.3.2 Error Analysis . . . . .	53
4.3.3 Computing the Pseudo-Inverse . . . . .	57
4.4 Determining the Optimal Tuning Parameter . . . . .	59
4.5 Morpho-Poroelasticity . . . . .	61
4.5.1 Relation to Morphoelastic System . . . . .	62
4.5.2 Derivation of System of Equations . . . . .	62
4.5.3 Stabilization . . . . .	64
4.6 Comparison Between Morphoelastic and Morpho-Poroelastic Models . . . . .	65
<b>5 Tumor Growth</b>	<b>66</b>
5.1 Mathematical Model . . . . .	66
5.1.1 Complete Model Equations . . . . .	68
5.1.2 Testing Configurations . . . . .	68
5.2 Analysis of Zero-Dimensional Nutrient Decay & Tissue Growth . . . . .	70
5.2.1 Analysis of Nutrient Absorption . . . . .	70
5.2.2 Coupling Nutrient Absorption and Volume Fraction Growth . . . . .	72
5.3 Finite Element Approach . . . . .	72
5.3.1 Nutrient Transport Equation . . . . .	72
5.3.2 Volume Fraction Evolution . . . . .	74
5.3.3 Updating New Tissue Volume . . . . .	75
5.3.4 Summary of Finite Element Model . . . . .	76
5.4 Numerical Experiments . . . . .	76
5.4.1 Configuration (i) . . . . .	77
5.4.2 Configuration (ii) . . . . .	80
5.4.3 Configuration (iii) . . . . .	80
<b>6 Conclusion</b>	<b>84</b>
6.1 Suggestions for Further Work . . . . .	85
6.2 Suggestions for Improving Model Performance . . . . .	86
<b>Bibliography</b>	<b>87</b>
<b>A Element Integrals</b>	<b>88</b>
A.1 Pure Elasticity . . . . .	89
A.2 Morphoelasticity . . . . .	89
A.3 Bubble Functions . . . . .	91
A.4 Tumor Growth Model . . . . .	92
<b>B Matrix Inversion by Analyzing Delta Problems</b>	<b>93</b>
B.1 Inverting the Laplace Matrix in 1D . . . . .	93
B.2 Connection to Green's Functions . . . . .	94
B.3 Generalization to Non-Uniform Mesh . . . . .	96
B.4 Higher Order Problems . . . . .	96
B.5 Mixed Boundary Conditions . . . . .	97

## List of Symbols

In this work we develop two-dimensional models. These models are assumed slices of three-dimensional objects in which there is no change in the  $z$ -direction. As such, the units we use are the standard three-dimensional SI units.

Table 1: List of symbols and their corresponding SI units.

Symbol	Description	Units
$\mathbf{x}$	Eulerian position	m
$\mathbf{X}$	Lagrangian position	m
$\mathbf{u}$	Displacement	m
$\mathbf{v}$	Displacement velocity	$\text{m s}^{-1}$
$t$	Time	s
$\mathbf{g}$	Body force density	$\text{N m}^{-3}$
$\boldsymbol{\tau}$	Traction vector	Pa
$\sigma_{ij}$	Stress	Pa
$\mathbf{n}$	Unit normal vector	1
$\varepsilon_{ij}$	Strain	1
$\mu_\varepsilon, \lambda_\varepsilon$	Lamé's first and second parameter	Pa
$E$	Young's modulus	Pa
$\nu$	Poisson ratio	1
$\rho$	Density	$\text{kg m}^{-3}$
$\mu_v, \lambda_v$	Viscosities of viscoelastic medium	Pa s
$p$	Poro-fluid pressure	Pa
$\Phi$	Fraction of space occupied by solid material	1
$k_{\text{perm}}$	Permeability of porous medium	$\text{m}^2$
$\mu_{\text{dyn}}$	Dynamic viscosity of poro-fluid	Pa s
$c$	Oxygen (nutrient) concentration	$\text{Mol m}^{-3}$
$\eta$	Ratio of created tissue per existing tissue	1
$\lambda_N$	Nutrient diffusion coefficient	$\text{m}^2 \text{s}^{-1}$
$N_G$	Tumor growth rate	$\text{s}^{-1}$
$G_{\text{max}}$	Maximum growth rate	$\text{s}^{-1}$
$c_G$	Oxygen concentration at which growth rate is half of $G_{\text{max}}$	$\text{Mol m}^{-3}$
$r_d$	Tumor death rate	$\text{s}^{-1}$
$N_A$	Nutrient absorption growth rate	$\text{s}^{-1}$
$A_{\text{max}}$	Maximum nutrient absorption rate	$\text{s}^{-1}$
$c_A$	Oxygen concentration at which absorption rate is half of $A_{\text{max}}$	$\text{Mol m}^{-3}$



# Introduction

The development of a cancerous tumor can be divided into three broad stages [19]. First, there is the initiation or transformation stage. This is a process in which normal cells are altered such that they acquire the potential for autonomous growth and are able to form tumors. Initiation is caused by mutations in the genetic material, DNA. These mutations can occur spontaneously, or they can be induced by environmental agents such as radiation, viruses or chemicals. It has also been shown that people with so-called germline mutations in specific genes leading to hereditary predisposition are more likely to develop cancer. The second stage of tumor development is growth. This stage consists of the two processes, promotion and progression. In promotion, initiated cells can expand and resist mechanisms that would otherwise curtail further development of a tumor. Then, during the progression process, the tumor starts growing rapidly. Tumor growth is dependent on the presence of oxygen and glucose. The final stage of cancer development is called metastasis, which occurs when cancer cells break away from the initial tumor and enter the circulatory system. New tumors can then be formed at different locations in the body.

In this work, we focus on the mechanical aspects of the growth stage of tumor development, specifically the tumor progression process. To this end, a tumor growth model is developed based on the novel combination of morpho- and poroelasticity. Poroelasticity provides us with a mathematical framework that models the interaction of fluid and solid phases in a porous medium. An intuitive way to think out this is to consider a dry sponge that is submerged into a body of water. Due to the water entering the sponge, it will expand. In the field of biomechanics, poroelasticity has been used extensively. For example, to model the flow of blood through arteries [3], bone deformations [5], and of course tumor growth [7, 14]. The addition of morphoelasticity a novel approach. Morphoelasticity is the theory of elastic growth, hence it allows us to take into account permanent deformations of tissue.

In the work done by Eline Kleimann [10] and Daan Smits [16] in their master graduation projects, morphoelasticity is applied to the healing of burn wounds. Their work is based on the PhD thesis of Daniël Koppenol [11], in which a number of different biomedical models are developed to study the healing and treatment of dermal wounds. Although we will apply morphoelasticity to the different field of tumor growth, the underlying mathematics are largely the same. As a tumor grows, it can push the surrounding tissue away. This can lead to complications. For example, a brain tumor exerts force on the surrounding brain tissue, possibly causing permanent damage. Even after a tumor has been surgically removed, the surrounding tissue might be permanently deformed or damaged. The motivation for using a morphoelastic approach is to model these permanent deformations.

In addition to the combination of poro- and morphoelasticity, which we will call the mechanical part of the tumor growth model, we must also consider some biological aspects of tumor growth. For example, a tumor needs nutrients to grow. The biochemical part of the model is based on the work of Tiina Roose et al. [14], in which poroelasticity is coupled with nutrient transport to model the growth of spheroid tumor cells.

This project has a number of goals. First, we want to investigate whether it is possible to construct a finite element model that combines morpo- and poroelasticity. Since this combination is a novel approach, it is unclear whether a standard finite element approach will work. The second goal is to combine this

mechanical model with a biochemical model based on [14].

## 1.1 Thesis Outline

The outline of this thesis is briefly discussed in the following section. It is divided into two parts. In the first part we develop finite element schemes for existing elastic models. The complexity of these models increases, leading up to a morphoelastic model at the end of Chapter 2, and a visco- poroelastic model at the end of Chapter 3. These models form the building blocks of the eventual tumor growth model, in which morpho- and poroelasticity are combined. The second part of this thesis contains new work. First we show the derivation of the stabilization parameter for a visco- poroelastic system. Next a finite element model for a morpho-poroelastic model is developed. In Chapter 5 morpho-poroelasticity is used as the backbone of the tumor growth model. We give a short summary of the chapters.

- **Part I**

- **Chapter 2:** A finite element model is developed for increasingly complex elastic models. We start with a regular elasticity equation, then move on to viscoelasticity in which the inertial forces must be considered. Finally, we arrive at morphoelasticity. The model we develop here is effectively the same as the one used in [16]. Each section has a similar setup. First the weak form is derived, then the Galerkin equations are constructed and finally some numerical experiments are run to test the models. In the derivation of the weak forms, a number of theorems are used. We give the proof for some of these theorems, whereas for others we provide a short motivation and cite an external source that contains the proof.
- **Chapter 3:** The focus is shifted towards poroelasticity. First, a one-dimensional system called the Terzaghi problem is analyzed. It is well known that even for this simple problem the solution can contain non-physical oscillations. We reproduce these oscillations and investigate their severity when a different type of element is used (Taylor-Hood). Next, two-dimensional poroelastic models are developed. Oscillations also occur in these cases. Different methods of stabilization are compared.

- **Part II**

- **Chapter 4:** We investigate in more detail one of the methods used to stabilize a poroelastic system: diffusive stabilization. In this method a tuning parameter must be chosen, which has a big effect on the performance. In literature such as [1, 13] an optimal value of this parameter is shown for the one-dimensional Terzaghi problem. We prove that this parameter is also (approximately) optimal for a visco- and morpho-poroelastic system. To this end, we use a novel technique which can also be used to find the exact expression of the inverse of a Laplace matrix. This tangential result is explored in Appendix B. In the latter part of this chapter, the viscoelastic and morphoelastic models from Sections 2.3 and 2.4 are combined with poroelasticity. The results on stabilization can also be applied to the morpho-poroelastic model.
- **Chapter 5:** The morpho-poroelastic model is combined with a biochemical tumor growth model from [14]. First the behaviour of the newly introduced biochemical unknowns is analyzed by considering zero-dimensional problems. We subsequently build the complete finite element scheme for the tumor growth model, and use it to perform numerical experiments.

- **Appendices**

- **Appendix A:** This chapter contains the element matrices and vectors for any integral that occurs in the Galerkin equations.
- **Appendix B:** We use the same approach as in Chapter 4 to find an exact expression for the inverse of the Laplace matrix in one dimension. Some additional results are shown, such as the connection to Green's functions and higher order problems.

## 1.2 A Note on Implementation

The finite element models are implemented in an object oriented fashion using Python 3.6. We make heavy use of the library NumPy, which allows us to work efficiently with vectors. In addition to NumPy, the sparse matrix library from SciPy is used. Since all finite element coefficient matrices are constructed



---

by looping over all elements, the construction of a sparse matrix object will not be faster than a regular matrix. However, solving the resulting linear system is much faster for certain classes of sparse matrices. We use the CSR (compressed row storage) class. All figures shown in this work are created using the Matplotlib visualization library.

For generating the finite element meshes the library ‘dmsh’ is used [15]. It supports automatically generated meshes for all kinds of geometries in two dimensions. The meshes it produces are of high quality, in the sense that every triangular element is close to an equilateral triangle. Moreover, dmsh has a user-friendly interface due to its pure Python source code. Unfortunately, mesh generation is quite slow and requires a lot of memory. For these reasons, we do not make use of remeshing. Thus only at the start of a simulation dmsh is used to generate an initial mesh. During the simulation grid nodes are able to move, causing elements to deform. It could be beneficial to create a new mesh with respect to the updated domain.

# Part I

# 2

## A Sequence of Models for Elasticity

In this chapter three finite element models are constructed for various elastic systems in two dimensions. We first briefly discuss the basic mathematical framework used to model stress and strain. We subsequently consider the elasticity equation, in which the reaction to a force is instantaneous. A viscous term is then added to expand the model to a viscoelastic system. Finally, we consider the morphoelastic system. Here the strain tensor is made time-dependent and must satisfy a partial differential equation. The right-hand side of this equation is the growth tensor, which allows us to model permanent deformations of the material. Each section has a similar structure. We begin by discussing the mathematical model and present the corresponding equations. Then the weak form is constructed and the corresponding Galerking equations are derived, leading up to a complete algebraic system of equations. We end each section by performing some numerical experiments to test the finite element model. The last section of this chapter contains a comparison between the three models.

### 2.1 Stress and Strain

We first define the notions of stress and strain, and show the derivation of the force balance equation. Stress is a quantity that measures the internal forces in a material. At a given point in the material, the component  $\sigma_{ij}$  denotes the stress along the  $i$ 'th coordinate direction acting on a surface of which the normal points in the  $j$ 'th coordinate direction. Hence, in three dimensions, the stress in a point is completely defined by nine components: there are three normal directions with each three stress components. The nine components are collected in the stress tensor:

$$\boldsymbol{\sigma} = \begin{pmatrix} \sigma_{11} & \sigma_{12} & \sigma_{13} \\ \sigma_{21} & \sigma_{22} & \sigma_{23} \\ \sigma_{31} & \sigma_{32} & \sigma_{33} \end{pmatrix}. \quad (2.1)$$

In two dimensions, there are two normal directions, hence the stress tensor becomes a  $2 \times 2$  array:

$$\boldsymbol{\sigma} = \begin{pmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{pmatrix}. \quad (2.2)$$

The components  $\sigma_{ii}$  are called normal stresses, they result from a force applied perpendicular to a surface. The other components  $\sigma_{ij}$  for  $i \neq j$  are called the shear stresses. They result from a force applied parallel to a surface. Given a normalized direction vector  $\mathbf{n}$ , the traction vector  $\boldsymbol{\tau}^{(\mathbf{n})}$  acting on the surface normal to  $\mathbf{n}$  can be computed using the stress tensor:

$$\boldsymbol{\tau}^{(\mathbf{n})} = \mathbf{n} \cdot \boldsymbol{\sigma} = \boldsymbol{\sigma} \cdot \mathbf{n}. \quad (2.3)$$

We will use the two notations  $\mathbf{n} \cdot \boldsymbol{\sigma}$  and  $\boldsymbol{\sigma} \cdot \mathbf{n}$  interchangeably, they both denote the standard matrix-vector product of  $\boldsymbol{\sigma}$  and  $\mathbf{n}$ . Thus:

$$\tau_i^{(\mathbf{n})} = \sum_j \sigma_{ij} n_j. \quad (2.4)$$

Consider a closed volume  $\Omega$  with piecewise smooth boundary  $\Gamma$ . Let  $\mathbf{n}$  denote the outward pointing unit normal vector to  $\Gamma$ . Assume that  $\Omega$  is subject to a body force  $\mathbf{g}$ . By Newton's third law, the action of

the body force will result in an opposite traction force on the boundary:

$$-\int_{\Gamma} \mathbf{n} \cdot \boldsymbol{\sigma} \, d\Gamma = \int_{\Omega} \mathbf{g} \, d\Omega. \quad (2.5)$$

We apply Gauss' divergence theorem to obtain

$$-\int_{\Omega} \nabla \cdot \boldsymbol{\sigma} \, d\Omega = \int_{\Omega} \mathbf{g} \, d\Omega. \quad (2.6)$$

If this holds for every volume, then the above integral turns into a pointwise equality:

$$-\nabla \cdot \boldsymbol{\sigma} = \mathbf{g}. \quad (2.7)$$

This equation is the momentum balance equation. All models for elasticity are based on this equation. Their differences lie in the choice of the stress tensor  $\boldsymbol{\sigma}$ . In the purely elastic model,  $\boldsymbol{\sigma}$  only contains an elastic component, whereas in the visco- and morphoelastic models, an inertial component involving the displacement velocity is also included. The elastic component of the stress tensor is given in terms of the strain tensor. While there are many different definitions of strain, we will use the Eulerian strain, given by:

$$\boldsymbol{\varepsilon} = \frac{1}{2} (\nabla \mathbf{u} + \nabla \mathbf{u}^{\top}). \quad (2.8)$$

Here  $\mathbf{u} = \mathbf{x} - \mathbf{X}$  is the displacement vector, it relates the initial position of a particle  $\mathbf{X}$  to its current position  $\mathbf{x}$ . The gradient  $\nabla$  is relative to the current coordinates. In Section 2.3.1 we will give more details regarding the different coordinate systems. It is important to note that  $\boldsymbol{\varepsilon}$  is linear with respect to the partial derivatives of  $\mathbf{u}$ .

## 2.2 Pure Elasticity

We will first consider pure elasticity on a fixed domain  $\Omega$ . The equations that we consider are stationary, meaning that there will be no timestepping required. The momentum balance equation on a fixed domain  $\Omega$  is given by:

$$-\nabla \cdot \boldsymbol{\sigma}(\mathbf{u}) = \mathbf{g}, \quad (2.9)$$

where  $\mathbf{g}$  is a body force, and the stress tensor is given by:

$$\begin{aligned} \boldsymbol{\sigma}(\mathbf{u}) &= \frac{1}{2} \mu_{\varepsilon} (\nabla \mathbf{u} + \nabla \mathbf{u}^{\top}) + \lambda_{\varepsilon} (\nabla \cdot \mathbf{u}) \mathbf{I} \\ &= \mu_{\varepsilon} \boldsymbol{\varepsilon}(\mathbf{u}) + \lambda_{\varepsilon} \text{tr}(\boldsymbol{\varepsilon}(\mathbf{u})) \mathbf{I} \end{aligned} \quad (2.10)$$

Here the strain tensor  $\boldsymbol{\varepsilon}$  is defined in (2.8). The constants  $\mu_{\varepsilon}$  and  $\lambda_{\varepsilon}$  are called the Lamé constants. They relate to the Young's modulus  $E$  and Poisson ratio  $\nu$  in the following way:

$$\mu_{\varepsilon} = \frac{E}{1 + \nu}, \quad \lambda_{\varepsilon} = \frac{E\nu}{(1 + \nu)(1 - 2\nu)}. \quad (2.11)$$

We consider a Neumann boundary condition of the form

$$\boldsymbol{\sigma} \cdot \mathbf{n} = \boldsymbol{\tau}, \quad (2.12)$$

which represents a shear force on the boundary. A homogeneous Dirichlet boundary conditions of the form

$$\mathbf{u} = \mathbf{0} \quad (2.13)$$

can also be imposed on part of the boundary. This corresponds to fixing this part of the boundary, such that it can not move as a result of the applied force. To summarize, we will solve the following boundary value problem:

$$-\nabla \cdot \boldsymbol{\sigma}(\mathbf{u}) = \mathbf{g}, \quad \text{in } \Omega, \quad (2.14)$$

$$\boldsymbol{\sigma} \cdot \mathbf{n} = \boldsymbol{\tau}, \quad \text{on } \Gamma_1, \quad (2.15)$$

$$\mathbf{u} = \mathbf{0}, \quad \text{on } \Gamma_2. \quad (2.16)$$

The boundary  $\Gamma$  of  $\Omega$  consists of a 'Neumann'-part  $\Gamma_1$  and a 'Dirichlet'-part  $\Gamma_2$ , such that  $\Gamma_1 \cap \Gamma_2 = \emptyset$ .

### 2.2.1 Weak Formulation

We derive the weak form of the above equation. Let  $\mathbf{H}^1(\Omega)$  be the space of all functions  $\mathbf{u} : \Omega \rightarrow \mathbb{R}^2$  such that  $\mathbf{u} \in \mathbf{L}^2(\Omega)$  and all its partial derivatives are also in  $\mathbf{L}^2(\Omega)$ . Let  $\mathbf{w}$  be a test function in  $\mathbf{H}^1(\Omega)$  such that  $\mathbf{w}|_{\Gamma_2} = \mathbf{0}$ . Take the inner product of  $\mathbf{w}$  with the equation and integrate over  $\Omega$ :

$$-\int_{\Omega} \mathbf{w} \cdot (\nabla \cdot \boldsymbol{\sigma}) \, d\Omega = \int_{\Omega} \mathbf{w} \cdot \mathbf{g} \, d\Omega. \quad (2.17)$$

The left hand-side can be expanded using Green's formula and Gauss' divergence theorem:

$$\begin{aligned} -\int_{\Omega} \mathbf{w} \cdot (\nabla \cdot \boldsymbol{\sigma}) \, d\Omega &= -\int_{\Omega} \nabla \cdot (\mathbf{w} \cdot \boldsymbol{\sigma}) - \nabla \mathbf{w} : \boldsymbol{\sigma} \, d\Omega \\ &= \int_{\Omega} \nabla \mathbf{w} : \boldsymbol{\sigma} \, d\Omega - \int_{\Gamma} (\mathbf{w} \cdot \boldsymbol{\sigma}) \cdot \mathbf{n} \, d\Gamma \\ &= \int_{\Omega} \nabla \mathbf{w} : \boldsymbol{\sigma} \, d\Omega - \int_{\Gamma_1} \mathbf{w} \cdot \boldsymbol{\tau} \, d\Gamma \end{aligned} \quad (2.18)$$

In the last step, we used the fact that  $\mathbf{w}|_{\Gamma_2} = \mathbf{0}$  and that  $(\mathbf{w} \cdot \boldsymbol{\sigma}) \cdot \mathbf{n} = \mathbf{w} \cdot (\boldsymbol{\sigma} \cdot \mathbf{n}) = \mathbf{w} \cdot \boldsymbol{\tau}$  on  $\Gamma_1$ . We obtain the following weak formulation:

Find  $\mathbf{u} \in \mathbf{H}_{\Gamma_2}^1(\Omega)$  such that for all  $\mathbf{w} \in \mathbf{H}_{\Gamma_2}^1(\Omega)$  we have:

$$\int_{\Omega} \boldsymbol{\varepsilon}(\mathbf{w}) : \boldsymbol{\sigma}(\mathbf{u}) \, d\Omega = \int_{\Omega} \mathbf{w} \cdot \mathbf{g} \, d\Omega + \int_{\Gamma_1} \mathbf{w} \cdot \boldsymbol{\tau} \, d\Gamma. \quad (2.19)$$

Here  $\mathbf{H}_{\Gamma_2}^1(\Omega)$  is the linear subspace of  $\mathbf{H}^1(\Omega)$  containing functions that vanish on  $\Gamma_2$ .

We write  $\boldsymbol{\sigma} = \boldsymbol{\sigma}(\mathbf{u})$  to stress<sup>1</sup> that  $\boldsymbol{\sigma}$  is a functional acting on  $\mathbf{u}$ . Also note that in the left-hand side,  $\nabla \mathbf{w}$  has turned into  $\boldsymbol{\varepsilon}(\mathbf{w})$ . Due to the stress tensor  $\boldsymbol{\sigma}$  being symmetric, the double dot products are equal:

$$\begin{aligned} \boldsymbol{\varepsilon}(\mathbf{w}) : \boldsymbol{\sigma}(\mathbf{u}) &= \frac{1}{2} (\nabla \mathbf{w} + \nabla \mathbf{w}^{\top}) : \boldsymbol{\sigma}(\mathbf{u}) \\ &= \frac{1}{2} \nabla \mathbf{w} : \boldsymbol{\sigma}(\mathbf{u}) + \frac{1}{2} \nabla \mathbf{w}^{\top} : \boldsymbol{\sigma}(\mathbf{u}) \\ &= \frac{1}{2} \nabla \mathbf{w} : \boldsymbol{\sigma}(\mathbf{u}) + \frac{1}{2} \nabla \mathbf{w} : \boldsymbol{\sigma}(\mathbf{u})^{\top} \\ &= \nabla \mathbf{w} : \boldsymbol{\sigma}(\mathbf{u}) \end{aligned} \quad (2.20)$$

The element matrices and vectors corresponding to the integrals in (2.19) are given in Section A.1.

### 2.2.2 Galerkin Equations

We derive the Galerkin equations corresponding to the weak form (2.19), which is equivalent to

$$\int_{\Omega} \nabla \mathbf{w} : \boldsymbol{\sigma}(\mathbf{u}) \, d\Omega = \int_{\Omega} \mathbf{w} \cdot \mathbf{g} \, d\Omega + \int_{\Gamma_1} \mathbf{w} \cdot \boldsymbol{\tau} \, d\Gamma. \quad (2.21)$$

From a computational perspective, it is easier to use  $\nabla \mathbf{w}$  instead of  $\boldsymbol{\varepsilon}(\mathbf{w})$ . Let  $\mathbf{u}_h$  be the finite element solution, and let  $\varphi_j$  be a scalar-valued basis function in mesh point  $j$ . We write

$$\mathbf{u}_h(\mathbf{x}) = \sum_{j=1}^n (u_j^1 \varphi_j^1(\mathbf{x}) + u_j^2 \varphi_j^2(\mathbf{x})), \quad (2.22)$$

where the indices  $j = 1, \dots, n$  correspond to grid nodes that are not on the fixed part of the boundary  $\Gamma_2$ . The vector valued basis functions  $\varphi_j^1$  and  $\varphi_j^2$  are given by:

$$\varphi_j^1 = \begin{pmatrix} \varphi_j \\ 0 \end{pmatrix}, \quad \text{and} \quad \varphi_j^2 = \begin{pmatrix} 0 \\ \varphi_j \end{pmatrix}. \quad (2.23)$$

<sup>1</sup>Pun not intended.

Substitute  $\mathbf{u}_h$  into (2.21) to obtain

$$\sum_{j=1}^n \left[ u_j^1 \int_{\Omega} \nabla \mathbf{w} : \boldsymbol{\sigma}(\boldsymbol{\varphi}_j^1) \, d\Omega + u_j^2 \int_{\Omega} \nabla \mathbf{w} : \boldsymbol{\sigma}(\boldsymbol{\varphi}_j^2) \, d\Omega \right] = \int_{\Omega} \mathbf{w} \cdot \mathbf{g} \, d\Omega + \int_{\Gamma_1} \mathbf{w} \cdot \boldsymbol{\tau} \, d\Gamma. \quad (2.24)$$

Setting  $\mathbf{w}$  equal to  $\boldsymbol{\varphi}_i^1$  and  $\boldsymbol{\varphi}_i^2$  for  $i = 1, \dots, n$  yields  $2n$  linear equations for the coefficients  $u_j^1$  and  $u_j^2$ . These equations are of the form

$$\begin{cases} S^{11} \mathbf{u}^1 + S^{12} \mathbf{u}^2 &= \mathbf{g}^1 + \boldsymbol{\tau}^1 \\ S^{21} \mathbf{u}^1 + S^{22} \mathbf{u}^2 &= \mathbf{g}^2 + \boldsymbol{\tau}^2 \end{cases}, \quad \text{or} \quad \begin{pmatrix} S^{11} & S^{12} \\ S^{21} & S^{22} \end{pmatrix} \begin{pmatrix} \mathbf{u}^1 \\ \mathbf{u}^2 \end{pmatrix} = \begin{pmatrix} \mathbf{g}^1 + \boldsymbol{\tau}^1 \\ \mathbf{g}^2 + \boldsymbol{\tau}^2 \end{pmatrix}, \quad (2.25)$$

where the latter representation is in block-matrix form. The coefficient vectors  $\mathbf{u}^1$  and  $\mathbf{u}^2$  contain the unknowns:

$$\mathbf{u}^1 = (u_1^1, u_2^1, \dots, u_n^1)^\top, \quad \mathbf{u}^2 = (u_1^2, u_2^2, \dots, u_n^2)^\top. \quad (2.26)$$

Moreover, the vectors  $\mathbf{g}^1$  and  $\mathbf{g}^2$  contain the integrals involving the body force  $\mathbf{g}$ , whereas the vectors  $\boldsymbol{\tau}^1$  and  $\boldsymbol{\tau}^2$  contain integrals involving the shear force  $\boldsymbol{\tau}$ . The system can be written in the even shorter notation  $S\mathbf{u} = \mathbf{g} + \boldsymbol{\tau}$ , where  $\mathbf{u} = (\mathbf{u}^1, \mathbf{u}^2)^\top$ , etc. To compute the entries of the blocks  $S^{11}$ ,  $S^{12}$ ,  $S^{21}$  and  $S^{22}$ , we must first compute the gradients of the  $\boldsymbol{\varphi}_j^1$  and  $\boldsymbol{\varphi}_j^2$ :

$$\nabla \boldsymbol{\varphi}_j^1 = \begin{pmatrix} \frac{\partial \varphi_j}{\partial x} & \frac{\partial \varphi_j}{\partial y} \\ 0 & 0 \end{pmatrix}, \quad \nabla \boldsymbol{\varphi}_j^2 = \begin{pmatrix} 0 & 0 \\ \frac{\partial \varphi_j}{\partial x} & \frac{\partial \varphi_j}{\partial y} \end{pmatrix}. \quad (2.27)$$

From which it follows that

$$\text{sym}(\nabla \boldsymbol{\varphi}_j^1) = \begin{pmatrix} \frac{\partial \varphi_j}{\partial x} & \frac{1}{2} \frac{\partial \varphi_j}{\partial y} \\ \frac{1}{2} \frac{\partial \varphi_j}{\partial y} & 0 \end{pmatrix}, \quad \text{sym}(\nabla \boldsymbol{\varphi}_j^2) = \begin{pmatrix} 0 & \frac{1}{2} \frac{\partial \varphi_j}{\partial x} \\ \frac{1}{2} \frac{\partial \varphi_j}{\partial x} & \frac{\partial \varphi_j}{\partial y} \end{pmatrix}. \quad (2.28)$$

Also note:

$$\nabla \cdot \boldsymbol{\varphi}_j^1 = \frac{\partial \varphi_j}{\partial x}, \quad \nabla \cdot \boldsymbol{\varphi}_j^2 = \frac{\partial \varphi_j}{\partial y}. \quad (2.29)$$

The elements of the  $S$ -blocks are then given by:

$$S_{ij}^{11} = \int_{\Omega} \nabla \boldsymbol{\varphi}_i^1 : \boldsymbol{\sigma}(\boldsymbol{\varphi}_j^1) \, d\Omega = \int_{\Omega} (\mu_\varepsilon + \lambda_\varepsilon) \frac{\partial \varphi_i}{\partial x} \frac{\partial \varphi_j}{\partial x} + \frac{1}{2} \mu_\varepsilon \frac{\partial \varphi_i}{\partial y} \frac{\partial \varphi_j}{\partial y} \, d\Omega, \quad (2.30)$$

$$S_{ij}^{12} = \int_{\Omega} \nabla \boldsymbol{\varphi}_i^1 : \boldsymbol{\sigma}(\boldsymbol{\varphi}_j^2) \, d\Omega = \int_{\Omega} \lambda_\varepsilon \frac{\partial \varphi_i}{\partial x} \frac{\partial \varphi_j}{\partial y} + \frac{1}{2} \mu_\varepsilon \frac{\partial \varphi_i}{\partial y} \frac{\partial \varphi_j}{\partial x} \, d\Omega, \quad (2.31)$$

$$S_{ij}^{21} = \int_{\Omega} \nabla \boldsymbol{\varphi}_i^2 : \boldsymbol{\sigma}(\boldsymbol{\varphi}_j^1) \, d\Omega = \int_{\Omega} \frac{1}{2} \mu_\varepsilon \frac{\partial \varphi_i}{\partial x} \frac{\partial \varphi_j}{\partial y} + \lambda_\varepsilon \frac{\partial \varphi_i}{\partial y} \frac{\partial \varphi_j}{\partial x} \, d\Omega, \quad (2.32)$$

$$S_{ij}^{22} = \int_{\Omega} \nabla \boldsymbol{\varphi}_i^2 : \boldsymbol{\sigma}(\boldsymbol{\varphi}_j^2) \, d\Omega = \int_{\Omega} \frac{1}{2} \mu_\varepsilon \frac{\partial \varphi_i}{\partial x} \frac{\partial \varphi_j}{\partial x} + (\mu_\varepsilon + \lambda_\varepsilon) \frac{\partial \varphi_i}{\partial y} \frac{\partial \varphi_j}{\partial y} \, d\Omega. \quad (2.33)$$

The elements of the right-hand side vectors are given by:

$$g_i^1 = \int_{\Omega} \varphi_i g^1 \, d\Omega, \quad g_i^2 = \int_{\Omega} \varphi_i g^2 \, d\Omega \quad (2.34)$$

$$\tau_i^1 = \int_{\Gamma_1} \varphi_i \tau^1 \, d\Gamma, \quad \tau_i^2 = \int_{\Gamma_1} \varphi_i \cdot \tau^2 \, d\Gamma \quad (2.35)$$

### 2.2.3 Numerical Experiments

Let us test the finite element implementation on a simple situation. We consider a cross section of a beam of length  $L$  and height 1 that is fixed to a rigid wall at its left boundary. We assume the beam has uniform density and is subject to gravity. We model gravity using the body force

$$\mathbf{g}(x, y) = \begin{pmatrix} 0 \\ -F \end{pmatrix},$$

for some constant  $F$ . At the left boundary ( $x = 0$ ), the homogeneous Dirichlet boundary condition  $\mathbf{u} = \mathbf{0}$  is imposed. At all other boundaries, we assume there is zero stress in the normal direction, yielding the boundary condition  $\boldsymbol{\sigma} \cdot \mathbf{n} = \mathbf{0}$ . Having fully specified the problem, we apply the finite element model constructed in the previous section. For generating the following figures, we have fixed  $\mu_\varepsilon = 0.5$ ,  $\lambda_\varepsilon = 0$  and vary  $L$  and  $F$ . Figures 2.1 and 2.2 show the solution for  $L = 2$  and  $F = 0.01$ . A mesh size of  $h_{\max} = 0.1$  has been used.

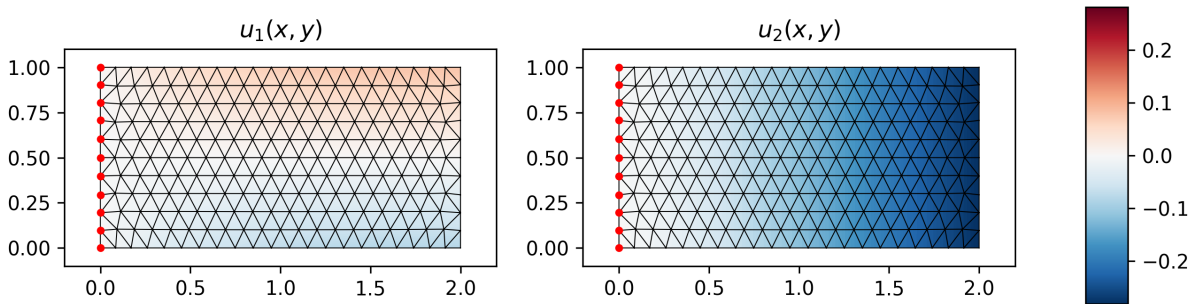


Figure 2.1: Contour plots of both components of the finite element solution for  $L = 2$  and  $F = 0.01$ .

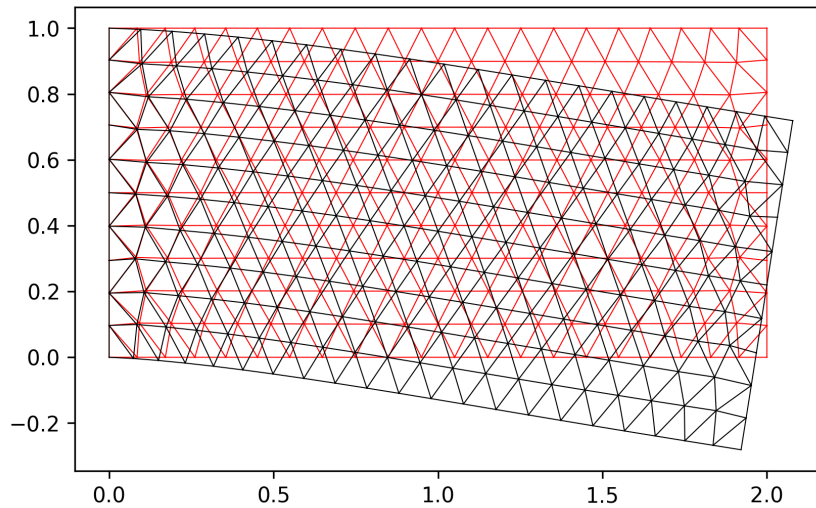


Figure 2.2: Mesh displacement corresponding to the solution in Figure 2.1.

We now double the magnitude of the body force, while keeping  $L = 2$  fixed.

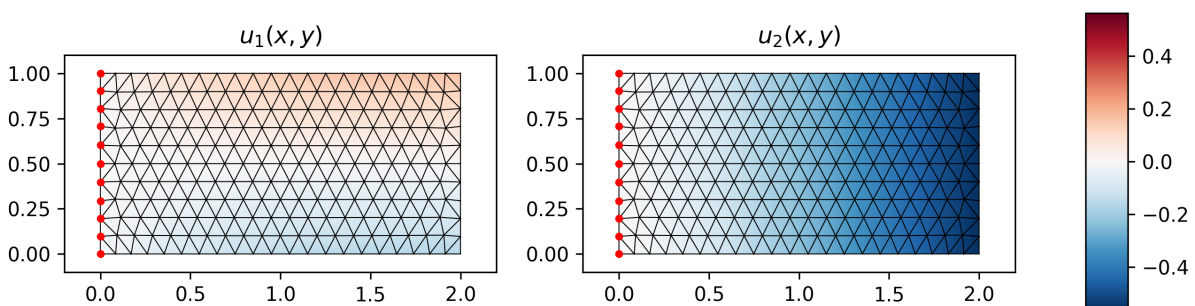


Figure 2.3: Contour plots of both components of the finite element solution for  $L = 2$  and  $F = 0.02$ .

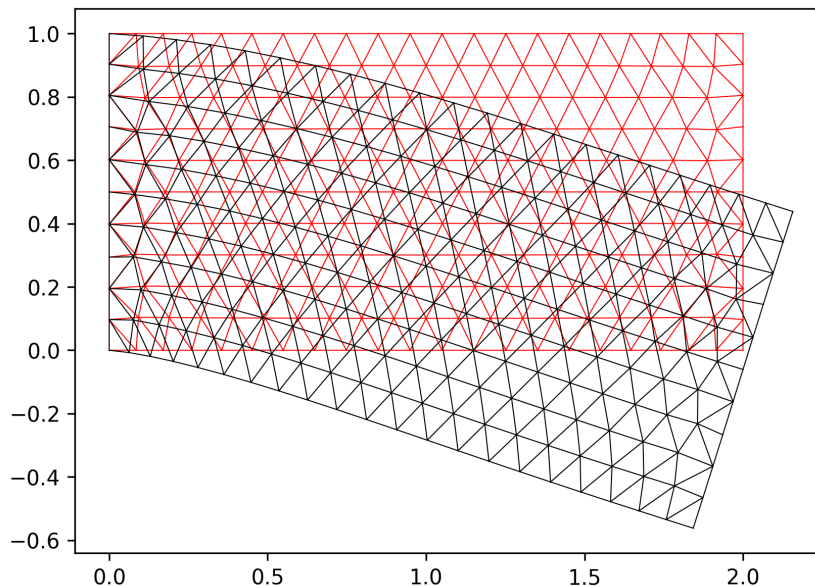


Figure 2.4: Mesh displacement corresponding to the solution in Figure 2.3.

## 2.3 Viscoelasticity

We add an inertial- or viscous term to the stationary elasticity equation, which will make the equations time-dependent. The viscoelastic momentum-balance equation is given by

$$\rho \left( \frac{D\mathbf{v}}{Dt} + (\nabla \cdot \mathbf{v})\mathbf{v} \right) - \nabla \cdot \boldsymbol{\sigma}(\mathbf{u}, \mathbf{v}) = \mathbf{g}, \quad (2.36)$$

where  $\mathbf{v}$  is the displacement velocity as defined in (2.44). Furthermore,  $\rho$  is the material density, and  $\mathbf{g}$  is again a body force. Due to the inclusion of inertial forces, the stress tensor must also be adapted. It now consists of a purely elastic part and a viscous part:

$$\boldsymbol{\sigma}(\mathbf{u}, \mathbf{v}) = \boldsymbol{\sigma}_{\text{el}}(\mathbf{u}) + \boldsymbol{\sigma}_{\text{vis}}(\mathbf{v}), \quad (2.37)$$

where

$$\boldsymbol{\sigma}_{\text{el}}(\mathbf{u}) = \frac{1}{2}\mu_{\varepsilon} (\nabla \mathbf{u} + \nabla \mathbf{u}^{\top}) + \lambda_{\varepsilon} (\nabla \cdot \mathbf{u})\mathbf{I} = \mu_{\varepsilon} \boldsymbol{\varepsilon}(\mathbf{u}) + \lambda_{\varepsilon} \text{tr}(\boldsymbol{\varepsilon}(\mathbf{u}))\mathbf{I}, \quad (2.38)$$

$$\boldsymbol{\sigma}_{\text{vis}}(\mathbf{v}) = \frac{1}{2}\mu_v (\nabla \mathbf{v} + \nabla \mathbf{v}^{\top}) + \lambda_v (\nabla \cdot \mathbf{v})\mathbf{I}. \quad (2.39)$$

Equation (2.36) is solved on a moving domain  $\Omega(t)$ , it is defined as:

$$\Omega(t) = \{ \mathbf{X} + \mathbf{u}(\mathbf{x}(\mathbf{X}, t), t) \mid \mathbf{X} \in \Omega(0) \}. \quad (2.40)$$

Furthermore, the gradient  $\nabla$  is relative to the current coordinates. In developing the finite element model we therefore have two choices. As a first option we can use a fixed mesh. In this setting, to compute the gradient with respect to the deformed domain, we must use a mapping that relates current coordinates to the initial coordinates in the fixed mesh. This is called an Eulerian approach. The second option would be to make the mesh dynamic: the grid nodes follow the flow of the domain. In this case we do not need a coordinate mapping, but we need to update the mesh in each timestep. The second method is called the Lagrangian approach, which is the approach that we choose to use. In Section 2.3.1 we give a brief overview of the mathematics behind the distinction between Eulerian and Lagrangian approaches.

### 2.3.1 Eulerian Versus Lagrangian Coordinates

The finite element model developed in this work is based on the Lagrangian framework. In this section, we will briefly discuss the difference between a Lagrangian and Eulerian specification. We can define the flow of a particle by considering its initial position  $\mathbf{X}$  and its position at time  $t$ , which we denote by  $\mathbf{x}(\mathbf{X}, t)$ . Note that we can also turn it around; assuming that a particle has position  $\mathbf{x}$  at some time  $t$ ,



we can find its initial position  $\mathbf{X}(\mathbf{x}, t)$ .

Let  $\varphi$  be a physical quantity that depends on space and time. In a Eulerian framework we consider the evolution of  $\varphi$  in a fixed point and write  $\varphi = \varphi(\mathbf{x}, t)$ . In the Lagrangian framework, the value of  $\varphi$  is considered along the trajectory of a particle. Thus  $\varphi$  effectively becomes a function of the initial position of the particle that is being followed:  $\tilde{\varphi}(\mathbf{X}, t) = \varphi(\mathbf{x}(\mathbf{X}, t), t)$ . We use  $\tilde{\varphi}$  instead of  $\varphi$  to emphasize it being a different function in the mathematical sense. Care must be taken in evaluating the rate of change of  $\varphi$ . The standard partial time-derivative  $\partial\varphi/\partial t$  measures the rate of change of  $\varphi$  at a fixed position. The full time-derivative takes into account the movement of particles:

$$\frac{d\varphi}{dt}(\mathbf{x}, t) = \frac{\partial\varphi}{\partial t} + \frac{\partial\varphi}{\partial\mathbf{x}} \cdot \frac{d\mathbf{x}}{dt}. \quad (2.41)$$

The above expression is called the *material derivative*. By introducing the notation  $d\mathbf{x}/dt = \mathbf{v}$  and writing  $\partial\varphi/\partial\mathbf{x} = \nabla_{\mathbf{x}}\varphi$ , the material derivative is denoted by

$$\frac{D\varphi}{Dt} = \frac{\partial\varphi}{\partial t} + \mathbf{v} \cdot \nabla_{\mathbf{x}}\varphi. \quad (2.42)$$

We write  $\nabla_{\mathbf{x}}$  to emphasize we are dealing with the Eulerian gradient: the spatial derivative are taken with respect to the current particle locations  $\mathbf{x}$ . In the rest of this work, the subscript is omitted and we simply use  $\nabla$  to denote the Eulerian gradient. The Eulerian and Lagrangian coordinates are related through the displacement vector  $\mathbf{u}$ , which is defined as the difference between the current and initial position of a particle. Thus, in the Lagrangian framework we have:

$$\mathbf{u}(\mathbf{X}, t) = \mathbf{x}(\mathbf{X}, t) - \mathbf{X}. \quad (2.43)$$

From this definition it follows that

$$\frac{D\mathbf{u}}{Dt} = \frac{d\mathbf{x}}{dt} = \mathbf{v}, \quad (2.44)$$

since  $d\mathbf{X}/dt = \mathbf{0}$ . Therefore,  $\mathbf{v}$  is called the *displacement velocity*. The relation between the initial position  $\mathbf{X}$  and the current position  $\mathbf{x}$  is encoded in the *deformation tensor*, which is given by

$$\mathbf{F} = \frac{\partial\mathbf{x}}{\partial\mathbf{X}}. \quad (2.45)$$

If  $\mathbf{F}$  is close to the identity matrix, then its inverse is too. In this case, the deformations are small, and we have

$$\begin{aligned} \frac{D\mathbf{u}}{Dt} &= \frac{\partial\mathbf{u}}{\partial t} + \nabla\mathbf{u} \cdot \mathbf{v} \\ &= \frac{\partial\mathbf{u}}{\partial t} + \nabla(\mathbf{x} - \mathbf{X}) \cdot \mathbf{v} \\ &= \frac{\partial\mathbf{u}}{\partial t} + (\mathbf{I} - \mathbf{F}^{-1}) \cdot \mathbf{v} \\ &\approx \frac{\partial\mathbf{u}}{\partial t} \end{aligned} \quad (2.46)$$

Thus, the material derivative of  $\mathbf{u}$  is closely approximated by the partial time-derivative. In this situation we are dealing with *infinitesimal deformations*. In most of this work, finite deformations are considered, meaning that  $\|\mathbf{F} - \mathbf{I}\| > 0$ . Only in the first two sections of Chapter 3 do we consider infinitesimal deformations, which is done to investigate the behaviour of the pressure.

As stated, our finite element model will be developed using a Lagrangian framework. As a result, the mesh nodes of the finite element mesh will follow the displacement velocity field.

### 2.3.2 Complete Viscoelastic Model Equations

To fully specify the viscoelastic system we need to consider the boundary conditions. We denote the boundary of  $\Omega(t)$  by  $\Gamma(t)$ . The boundary is again split into two disjoint parts:  $\Gamma(t) = \Gamma_1(t) \cup \Gamma_2(t)$  such that  $\Gamma_1(t) \cap \Gamma_2(t) = \emptyset$  for all  $t \geq 0$ . Since the domain is be moving, it is not straight-forward to determine where one part begins and the other ends. To solve this problem,  $\Gamma_1(t)$  is fixed in place by imposed the

boundary condition  $\mathbf{u} = \mathbf{0}$ . The other part  $\Gamma_2(t)$  is free to move, and may be subject to a shear force  $\tau$ . In summary, we obtain the following viscoelastic system:

$$\rho \left( \frac{D\mathbf{v}}{Dt} + (\nabla \cdot \mathbf{v})\mathbf{v} \right) - \nabla \cdot \boldsymbol{\sigma}(\mathbf{u}, \mathbf{v}) = \mathbf{g}, \quad \text{in } \Omega(t), \quad (2.47)$$

$$\frac{D\mathbf{u}}{Dt} = \mathbf{v}, \quad \text{in } \Omega(t), \quad (2.48)$$

with boundary conditions

$$\mathbf{u} = \mathbf{0}, \quad \text{on } \Gamma_1(t), \quad (2.49)$$

$$\boldsymbol{\sigma} \cdot \mathbf{n} = \boldsymbol{\tau}, \quad \text{on } \Gamma_2(t). \quad (2.50)$$

In subsequent sections, we derive the weak form and Galerkin equations corresponding to the system. To achieve this, an important result summarized in Theorem 3 is required. This theorem is proved using a number of well-known results, such as Reynold's Transport Theorem and the product rule for the material derivative. Theorem 3 is used numerous times throughout this work, as it allows us to greatly simplify convert complex integrals involving material derivatives.

### 2.3.3 Weak Formulation

We derive the weak form corresponding to equation (2.47). Let  $\varphi$  be a basis function. Take the inner product of  $\varphi$  with the momentum-balance equation and integrate over  $\Omega(t)$  to obtain the following weak form:

$$\rho \left[ \int_{\Omega(t)} \frac{D\mathbf{v}}{Dt} \cdot \varphi \, d\Omega + \int_{\Omega(t)} (\nabla \cdot \mathbf{v})\mathbf{v} \cdot \varphi \, d\Omega \right] - \int_{\Omega(t)} (\nabla \cdot \boldsymbol{\sigma}) \cdot \varphi \, d\Omega = \int_{\Omega(t)} \mathbf{f} \cdot \varphi \, d\Omega. \quad (2.51)$$

We shall use the following theorems and lemma to manipulate the weak form. First, we state Reynolds Transport Theorem:

**Theorem 1** (Reynold's Transport Theorem). *Let  $f(\mathbf{x}, t)$  be a sufficiently smooth scalar-valued function. Let  $\Omega(t)$  be a time-dependent domain with piecewise smooth boundary  $\Gamma(t)$ , and  $\mathbf{v}(\mathbf{x}, t)$  the velocity of a point  $\mathbf{x} \in \Omega(t)$ , then:*

$$\frac{d}{dt} \int_{\Omega(t)} f(\mathbf{x}, t) \, d\Omega = \int_{\Omega(t)} \frac{\partial f}{\partial t}(\mathbf{x}, t) \, d\Omega + \int_{\Gamma(t)} f(\mathbf{x}, t) \mathbf{v} \cdot \mathbf{n} \, d\Gamma,$$

where  $\mathbf{n}$  is the outward-pointing unit normal vector to  $\Gamma(t)$ .

We will not give the proof of this well-known theorem. In addition, we need the following result:

**Theorem 2** (Dziuk & Elliott). *If  $\varphi$  is a Lagrangian basis function, then*

$$\frac{D\varphi}{Dt} = 0.$$

The proof can be found in [6]. A Lagrangian basis function is a Lagrangian interpolation polynomial on an element that interpolates between the vertices. We will demonstrate the validity of the theorem by considering a linear interpolatory function in one dimension. Let  $x_1(t)$ ,  $x_2(t)$  be two moving grid nodes. We define

$$\varphi(x; t) = \frac{x_2(t) - x}{x_2(t) - x_1(t)} \varphi_1 + \frac{x - x_1(t)}{x_2(t) - x_1(t)} \varphi_2, \quad (2.52)$$

where  $\varphi_1$  and  $\varphi_2$  are constant in time. Thus,  $\varphi$  linearly interpolates the values  $\varphi_1$  and  $\varphi_2$  on the moving interval  $[x_1(t), x_2(t)]$ . It follows that:

$$\frac{\partial \varphi}{\partial x}(x; t) = \frac{\varphi_2 - \varphi_1}{x_2(t) - x_1(t)}, \quad (2.53)$$

and

$$\frac{\partial \varphi}{\partial t}(x; t) = -\frac{\varphi_2 - \varphi_1}{x_2(t) - x_1(t)} \left[ \frac{x_2(t) - x}{x_2(t) - x_1(t)} \frac{dx_1}{dt}(t) + \frac{x - x_1(t)}{x_2(t) - x_1(t)} \frac{dx_2}{dt}(t) \right]. \quad (2.54)$$

We can find the velocity at the point  $x$  by considering the velocities of the endpoints:

$$v(x, t) = \frac{x_2(t) - x}{x_2(t) - x_1(t)} \frac{dx_1}{dt}(t) + \frac{x - x_1(t)}{x_2(t) - x_1(t)} \frac{dx_2}{dt}(t). \quad (2.55)$$

From which it follows that

$$\frac{\partial \varphi}{\partial t}(x; t) = -\frac{\varphi_2 - \varphi_1}{x_2(t) - x_1(t)}v(x, t). \quad (2.56)$$

By the definition of the material derivative and equations (2.53) and (2.56), we get

$$\begin{aligned} \frac{D\varphi}{Dt} &= \frac{\partial \varphi}{\partial t} + v(x, t) \frac{\partial \varphi}{\partial x} \\ &= -\frac{\varphi_2 - \varphi_1}{x_2(t) - x_1(t)}v(x, t) + \frac{\varphi_2 - \varphi_1}{x_2(t) - x_1(t)}v(x, t) = 0, \end{aligned} \quad (2.57)$$

Showing that Theorem 2 indeed holds in this case. It follows immediately from Theorem 2 that for a vector valued Lagrangian basis function  $\boldsymbol{\varphi}$  we also have  $D\boldsymbol{\varphi}/Dt = \mathbf{0}$ . Indeed, we can write  $\boldsymbol{\varphi} = (\varphi_1, \dots, \varphi_d)^\top$ , thus

$$\frac{D\boldsymbol{\varphi}}{Dt} = \begin{pmatrix} \frac{D\varphi_1}{Dt} \\ \vdots \\ \frac{D\varphi_d}{Dt} \end{pmatrix} = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}. \quad (2.58)$$

The last ingredient for the proof of Theorem 3 is the product rule for the material derivative.

**Lemma 1** (Product Rule). *Let  $\Omega(t) \subset \mathbb{R}^d$  be a domain subject to a velocity  $\mathbf{v}(\mathbf{x}, t)$ . Let  $f(\mathbf{x}, t)$  and  $g(\mathbf{x}, t)$  be two differentiable functions on  $\Omega(t)$ . Then we have:*

$$\frac{D}{Dt}(\mathbf{f} \cdot \mathbf{g}) = \frac{D\mathbf{f}}{Dt} \cdot \mathbf{g} + \mathbf{f} \cdot \frac{D\mathbf{g}}{Dt}.$$

*Proof.* The statement is proved using the definition of the material derivative:

$$\begin{aligned} \frac{D}{Dt}(\mathbf{f} \cdot \mathbf{g}) &= \frac{\partial}{\partial t}(\mathbf{f} \cdot \mathbf{g}) + \mathbf{v} \cdot \nabla(\mathbf{f} \cdot \mathbf{g}) \\ &= \frac{\partial}{\partial t} \left( \sum_{i=1}^d f_i g_i \right) + \mathbf{v} \cdot \left( \sum_{i=1}^d \nabla(f_i g_i) \right) \\ &= \sum_{i=1}^d \left( \frac{\partial f_i}{\partial t} g_i + f_i \frac{\partial g_i}{\partial t} \right) + \sum_{k=1}^d v_k \left( \sum_{i=1}^d \frac{\partial}{\partial x_k}(f_i g_i) \right) \\ &= \sum_{i=1}^d \left( \frac{\partial f_i}{\partial t} g_i + f_i \frac{\partial g_i}{\partial t} \right) + \sum_{i=1}^d \sum_{k=1}^d v_k \left( \frac{\partial f_i}{\partial x_k} g_i + f_i \frac{\partial g_i}{\partial x_k} \right) \\ &= \sum_{i=1}^d \left( \frac{\partial f_i}{\partial t} g_i + f_i \frac{\partial g_i}{\partial t} \right) + \sum_{i=1}^d (g_i(\mathbf{v} \cdot \nabla f_i) + f_i(\mathbf{v} \cdot \nabla g_i)) \\ &= \sum_{i=1}^d \left( g_i \left( \frac{\partial f_i}{\partial t} + \mathbf{v} \cdot \nabla f_i \right) + f_i \left( \frac{\partial g_i}{\partial t} + \mathbf{v} \cdot \nabla g_i \right) \right) \\ &= \frac{D\mathbf{f}}{Dt} \cdot \mathbf{g} + \mathbf{f} \cdot \frac{D\mathbf{g}}{Dt}. \end{aligned}$$

□

Theorems 1, 2 and Lemma 1 are used to prove the following theorem. It is particularly useful for reducing complicated expressions involving the material derivative to easy to discretize terms using Finite Elements.

**Theorem 3.** *Let  $\mathbf{f}(\mathbf{x}, t)$  be a sufficiently smooth vector function. Furthermore, let  $\boldsymbol{\varphi}$  be a Lagrangian basis function, and let  $\mathbf{v}(\mathbf{x}, t)$  be the velocity of the point  $\mathbf{x} \in \Omega(t)$ . Then:*

$$\int_{\Omega(t)} \left( \frac{D\mathbf{f}}{Dt} + (\nabla \cdot \mathbf{v})\mathbf{f} \right) \cdot \boldsymbol{\varphi} \, d\Omega = \frac{d}{dt} \int_{\Omega(t)} \mathbf{f} \cdot \boldsymbol{\varphi} \, d\Omega.$$

*Proof.* Using Lemma 1 and Theorem 2, we find:

$$\int_{\Omega(t)} \frac{D\mathbf{f}}{Dt} \cdot \boldsymbol{\varphi} \, d\Omega = \int_{\Omega(t)} \frac{D}{Dt}(\mathbf{f} \cdot \boldsymbol{\varphi}) - \mathbf{f} \cdot \frac{D\boldsymbol{\varphi}}{Dt} \, d\Omega = \int_{\Omega(t)} \frac{D}{Dt}(\mathbf{f} \cdot \boldsymbol{\varphi}) \, d\Omega$$

Define  $\xi = \mathbf{f} \cdot \boldsymbol{\varphi}$ , then:

$$\begin{aligned} \int_{\Omega(t)} \frac{D}{Dt} (\mathbf{f} \cdot \boldsymbol{\varphi}) \, d\Omega &= \int_{\Omega(t)} \frac{D\xi}{Dt} \, d\Omega \\ &= \int_{\Omega(t)} \frac{\partial \xi}{\partial t} + \mathbf{v} \cdot \nabla \xi \, d\Omega \\ &= \int_{\Omega(t)} \frac{\partial \xi}{\partial t} + \nabla \cdot (\xi \mathbf{v}) - (\nabla \cdot \mathbf{v}) \xi \, d\Omega \\ &= \int_{\Omega(t)} \frac{\partial \xi}{\partial t} - (\nabla \cdot \mathbf{v}) \xi \, d\Omega + \int_{\Gamma(t)} \xi (\mathbf{v} \cdot \mathbf{n}) \, d\Gamma \end{aligned}$$

Using Theorem 1 we find:

$$\begin{aligned} \int_{\Omega(t)} \left( \frac{D\mathbf{f}}{Dt} + (\nabla \cdot \mathbf{v}) \mathbf{f} \right) \cdot \boldsymbol{\varphi} \, d\Omega &= \int_{\Omega(t)} \frac{D\xi}{Dt} + (\nabla \cdot \mathbf{v}) \xi \, d\Omega \\ &= \int_{\Omega(t)} \frac{\partial \xi}{\partial t} - (\nabla \cdot \mathbf{v}) \xi + (\nabla \cdot \mathbf{v}) \xi \, d\Omega + \int_{\Gamma(t)} \xi (\mathbf{v} \cdot \mathbf{n}) \, d\Gamma \\ &= \int_{\Omega(t)} \frac{\partial \xi}{\partial t} \, d\Omega + \int_{\Gamma(t)} \xi (\mathbf{v} \cdot \mathbf{n}) \, d\Gamma \\ &= \frac{d}{dt} \int_{\Omega(t)} \xi \, d\Omega = \frac{d}{dt} \int_{\Omega(t)} \mathbf{f} \cdot \boldsymbol{\varphi} \, d\Omega \end{aligned}$$

□

Thus, if we assume that  $\boldsymbol{\varphi}$  is a Lagrangian basis function, then by Theorem 3 the weak form (2.51) is equivalent to:

$$\rho \frac{d}{dt} \int_{\Omega(t)} \mathbf{v} \cdot \boldsymbol{\varphi} \, d\Omega - \int_{\Omega(t)} (\nabla \cdot \boldsymbol{\sigma}) \cdot \boldsymbol{\varphi} \, d\Omega = \int_{\Omega(t)} \mathbf{g} \cdot \boldsymbol{\varphi} \, d\Omega. \quad (2.59)$$

Using the same derivation as in (2.18), the weak form can be further expanded into

$$\rho \frac{d}{dt} \int_{\Omega(t)} \mathbf{v} \cdot \boldsymbol{\varphi} \, d\Omega + \int_{\Omega(t)} \boldsymbol{\sigma}(\mathbf{u}, \mathbf{v}) : \boldsymbol{\varepsilon}(\boldsymbol{\varphi}) \, d\Omega = \int_{\Omega(t)} \mathbf{g} \cdot \boldsymbol{\varphi} \, d\Omega + \int_{\Gamma_2(t)} \boldsymbol{\tau} \cdot \boldsymbol{\varphi} \, d\Gamma. \quad (2.60)$$

The stress tensor is given by (2.37). This splitting is useful, because the two tensors  $\boldsymbol{\sigma}_{\text{el}}$  and  $\boldsymbol{\sigma}_{\text{vis}}$  are very similar to the stress tensor in (2.10). Thus, the element matrix entries are also very similar to those derived in section 2.2.2.

### 2.3.4 Galerkin Equations

Let  $\mathbf{u}_h$  and  $\mathbf{v}_h$  be the Finite Element solutions, we write

$$\mathbf{u}_h(\mathbf{x}, t) = \sum_{j=1}^n u_j^1(t) \boldsymbol{\varphi}_j^1(\mathbf{x}; t) + u_j^2(t) \boldsymbol{\varphi}_j^2(\mathbf{x}; t), \quad (2.61)$$

$$\mathbf{v}_h(\mathbf{x}, t) = \sum_{j=1}^n v_j^1(t) \boldsymbol{\varphi}_j^1(\mathbf{x}; t) + v_j^2(t) \boldsymbol{\varphi}_j^2(\mathbf{x}; t). \quad (2.62)$$

where  $\boldsymbol{\varphi}_{j1}^1$  and  $\boldsymbol{\varphi}_j^2$  again contain the scalar-valued basis function  $\varphi_j$  in respectively their first and second position, and zero in the other. We use linear basis functions for  $\varphi_j(\mathbf{x}; t)$ , centered in the vertex  $\mathbf{x}_j$ . Note that the basis functions depend implicitly on  $t$ , since the vertices  $\mathbf{x}_j$  depend on time. Moreover, the basis functions are Lagrangian, making the weak form (2.60) valid. Substitute (2.61) and (2.62) into (2.60), and set  $\boldsymbol{\varphi}$  equal to  $\boldsymbol{\varphi}_i^1$  and  $\boldsymbol{\varphi}_i^2$  for  $i = 1, \dots, n$ . Due to the splitting of  $\boldsymbol{\sigma}$  into an elastic and viscous part, we obtain the following system of equations:

$$\frac{d}{dt} \begin{pmatrix} \rho M \mathbf{v}^1 \\ \rho M \mathbf{v}^2 \end{pmatrix} + \begin{pmatrix} S_{\text{el}}^{11} & S_{\text{el}}^{12} \\ S_{\text{el}}^{21} & S_{\text{el}}^{22} \end{pmatrix} \begin{pmatrix} \mathbf{u}^1 \\ \mathbf{u}^2 \end{pmatrix} + \begin{pmatrix} S_{\text{vis}}^{11} & S_{\text{vis}}^{12} \\ S_{\text{vis}}^{21} & S_{\text{vis}}^{22} \end{pmatrix} \begin{pmatrix} \mathbf{v}^1 \\ \mathbf{v}^2 \end{pmatrix} = \begin{pmatrix} \mathbf{g}^1 \\ \mathbf{g}^2 \end{pmatrix}, \quad (2.63)$$

where the vectors  $\mathbf{u}^1$ ,  $\mathbf{u}^2$ ,  $\mathbf{v}^1$  and  $\mathbf{v}^2$  contain the coefficients  $u_j^1$ ,  $u_j^2$ ,  $v_j^1$  and  $v_j^2$  respectively. The mass matrix  $M$  is given by

$$M_{ij} = \int_{\Omega(t)} \varphi_j \varphi_i \, d\Omega. \quad (2.64)$$

The right-hand side vectors of (2.63) now include both the body force and shear force, they are given by:

$$(\mathbf{g}^1)_i = \int_{\Omega(t)} \mathbf{g} \cdot \boldsymbol{\varphi}_i^1 \, d\Omega + \int_{\Gamma_2(t)} \boldsymbol{\tau} \cdot \boldsymbol{\varphi}_i^1 \, d\Gamma = \int_{\Omega(t)} g^1 \varphi_i \, d\Omega + \int_{\Gamma_2(t)} \tau^1 \varphi_i \, d\Gamma, \quad (2.65)$$

$$(\mathbf{g}^2)_i = \int_{\Omega(t)} \mathbf{g} \cdot \boldsymbol{\varphi}_i^2 \, d\Omega + \int_{\Gamma_2(t)} \boldsymbol{\tau} \cdot \boldsymbol{\varphi}_i^2 \, d\Gamma = \int_{\Omega(t)} g^2 \varphi_i \, d\Omega + \int_{\Gamma_2(t)} \tau^2 \varphi_i \, d\Gamma. \quad (2.66)$$

Moreover, the elements of the  $S_{\text{el}}$  and  $S_{\text{vis}}$  blocks are very similar to (2.30)-(2.33):

$$(S_{\text{el}}^{11})_{ij} = \int_{\Omega(t)} \nabla \boldsymbol{\varphi}_i^1 : \boldsymbol{\sigma}_{\text{el}}(\boldsymbol{\varphi}_j^1) \, d\Omega = \int_{\Omega(t)} (\mu_\varepsilon + \lambda_\varepsilon) \frac{\partial \varphi_i}{\partial x} \frac{\partial \varphi_j}{\partial x} + \frac{1}{2} \mu_\varepsilon \frac{\partial \varphi_i}{\partial y} \frac{\partial \varphi_j}{\partial y} \, d\Omega, \quad (2.67)$$

$$(S_{\text{el}}^{12})_{ij} = \int_{\Omega(t)} \nabla \boldsymbol{\varphi}_i^1 : \boldsymbol{\sigma}_{\text{el}}(\boldsymbol{\varphi}_j^2) \, d\Omega = \int_{\Omega(t)} \lambda_\varepsilon \frac{\partial \varphi_i}{\partial x} \frac{\partial \varphi_j}{\partial y} + \frac{1}{2} \mu_\varepsilon \frac{\partial \varphi_i}{\partial y} \frac{\partial \varphi_j}{\partial x} \, d\Omega, \quad (2.68)$$

$$(S_{\text{el}}^{21})_{ij} = \int_{\Omega(t)} \nabla \boldsymbol{\varphi}_i^2 : \boldsymbol{\sigma}_{\text{el}}(\boldsymbol{\varphi}_j^1) \, d\Omega = \int_{\Omega(t)} \frac{1}{2} \mu_\varepsilon \frac{\partial \varphi_i}{\partial x} \frac{\partial \varphi_j}{\partial y} + \lambda_\varepsilon \frac{\partial \varphi_i}{\partial y} \frac{\partial \varphi_j}{\partial x} \, d\Omega, \quad (2.69)$$

$$(S_{\text{el}}^{22})_{ij} = \int_{\Omega(t)} \nabla \boldsymbol{\varphi}_i^2 : \boldsymbol{\sigma}_{\text{el}}(\boldsymbol{\varphi}_j^2) \, d\Omega = \int_{\Omega(t)} \frac{1}{2} \mu_\varepsilon \frac{\partial \varphi_i}{\partial x} \frac{\partial \varphi_j}{\partial x} + (\mu_\varepsilon + \lambda_\varepsilon) \frac{\partial \varphi_i}{\partial y} \frac{\partial \varphi_j}{\partial y} \, d\Omega. \quad (2.70)$$

And:

$$(S_{\text{vis}}^{11})_{ij} = \int_{\Omega(t)} \nabla \boldsymbol{\varphi}_i^1 : \boldsymbol{\sigma}_v(\boldsymbol{\varphi}_j^1) \, d\Omega = \int_{\Omega(t)} (\mu_v + \lambda_v) \frac{\partial \varphi_i}{\partial x} \frac{\partial \varphi_j}{\partial x} + \frac{1}{2} \mu_v \frac{\partial \varphi_i}{\partial y} \frac{\partial \varphi_j}{\partial y} \, d\Omega, \quad (2.71)$$

$$(S_{\text{vis}}^{12})_{ij} = \int_{\Omega(t)} \nabla \boldsymbol{\varphi}_i^1 : \boldsymbol{\sigma}_v(\boldsymbol{\varphi}_j^2) \, d\Omega = \int_{\Omega(t)} \lambda_v \frac{\partial \varphi_i}{\partial x} \frac{\partial \varphi_j}{\partial y} + \frac{1}{2} \mu_v \frac{\partial \varphi_i}{\partial y} \frac{\partial \varphi_j}{\partial x} \, d\Omega, \quad (2.72)$$

$$(S_{\text{vis}}^{21})_{ij} = \int_{\Omega(t)} \nabla \boldsymbol{\varphi}_i^2 : \boldsymbol{\sigma}_v(\boldsymbol{\varphi}_j^1) \, d\Omega = \int_{\Omega(t)} \frac{1}{2} \mu_v \frac{\partial \varphi_i}{\partial x} \frac{\partial \varphi_j}{\partial y} + \lambda_v \frac{\partial \varphi_i}{\partial y} \frac{\partial \varphi_j}{\partial x} \, d\Omega, \quad (2.73)$$

$$(S_{\text{vis}}^{22})_{ij} = \int_{\Omega(t)} \nabla \boldsymbol{\varphi}_i^2 : \boldsymbol{\sigma}_v(\boldsymbol{\varphi}_j^2) \, d\Omega = \int_{\Omega(t)} \frac{1}{2} \mu_v \frac{\partial \varphi_i}{\partial x} \frac{\partial \varphi_j}{\partial x} + (\mu_v + \lambda_v) \frac{\partial \varphi_i}{\partial y} \frac{\partial \varphi_j}{\partial y} \, d\Omega. \quad (2.74)$$

Note that we have changed the strain tensor  $\boldsymbol{\varepsilon}(\boldsymbol{\varphi})$  in the weak form (2.60) into  $\nabla \boldsymbol{\varphi}$ . This is valid as per the derivation in (2.20). We can write the system in (2.63) in the even shorter notation

$$\rho \frac{d}{dt} (\widetilde{M} \mathbf{v}) + S_{\text{el}} \mathbf{u} + S_{\text{vis}} \mathbf{v} = \mathbf{g}, \quad (2.75)$$

where

$$\mathbf{u} = \begin{pmatrix} \mathbf{u}^1 \\ \mathbf{u}^2 \end{pmatrix}, \quad \mathbf{v} = \begin{pmatrix} \mathbf{v}^1 \\ \mathbf{v}^2 \end{pmatrix}, \quad \mathbf{g} = \begin{pmatrix} \mathbf{g}^1 \\ \mathbf{g}^2 \end{pmatrix}, \quad (2.76)$$

and

$$\widetilde{M} = \begin{pmatrix} M & 0 \\ 0 & M \end{pmatrix}, \quad S_{\text{el}} = \begin{pmatrix} S_{\text{el}}^{11} & S_{\text{el}}^{12} \\ S_{\text{el}}^{21} & S_{\text{el}}^{22} \end{pmatrix}, \quad S_{\text{vis}} = \begin{pmatrix} S_{\text{vis}}^{11} & S_{\text{vis}}^{12} \\ S_{\text{vis}}^{21} & S_{\text{vis}}^{22} \end{pmatrix}. \quad (2.77)$$

Note that we are ‘overloading’ the symbols  $\mathbf{u}$ ,  $\mathbf{v}$  and  $\mathbf{g}$ . They denote both a physical quantity and a numerical vector containing unknowns. Their meaning should be clear from the context.

### 2.3.5 Time Discretization

We use Implicit Euler to discretize (2.75) in time. It is important to note that the matrices  $M$ ,  $S_{\text{el}}$ ,  $S_{\text{vis}}$  and the vector  $\mathbf{g}$  all depend on  $t$ , since they all involve integrations over the time-dependent domain  $\Omega(t)$  of the time-dependent basis functions. Applying Implicit Euler yields

$$\left( \widetilde{M}^m + \Delta t S_{\text{vis}}^m \right) \mathbf{v}^m + \Delta t S_{\text{el}}^m \mathbf{u}^m = \widetilde{M}^{m-1} \mathbf{v}^{m-1} + \Delta t \mathbf{g}^m. \quad (2.78)$$

The superscript denotes a time level. To solve for  $\mathbf{v}^m$  and  $\mathbf{u}^m$ , we need another equation. This equation is obtained by discretizing the second equation (2.48). In the framework of a moving grid, this turns out to be rather straight-forward.

**Lemma 2.** *Let  $\mathbf{u}_i(t)$  be shorthand for  $\mathbf{u}(\mathbf{x}(\mathbf{X}_i, t), t)$ , then*

$$\frac{D\mathbf{u}_i}{Dt} = \frac{\mathbf{u}_i(t + \Delta t) - \mathbf{u}_i(t)}{\Delta t} + \mathcal{O}(\Delta t).$$

While the statement in Lemma 2 might seem obvious, it is more nuanced than a simple first order approximation of a time-derivative. The movement of the grid is encoded in the displacement vector  $\mathbf{u}$ .

*Proof.* Use a first order Taylor expansion to expand  $\mathbf{u}_i(t + \Delta t)$  around  $\mathbf{u}_i(t)$ :

$$\mathbf{u}_i(t + \Delta t) = \mathbf{u}_i(t) + \Delta t \frac{d}{dt} \mathbf{u}_i(t) + \mathcal{O}(\Delta t^2).$$

Let us calculate the time derivative of  $\mathbf{u}_i$ :

$$\begin{aligned} \frac{d}{dt} \mathbf{u}_i(t) &= \frac{d}{dt} \mathbf{u}(\mathbf{x}(\mathbf{X}_i, t), t) \\ &= \frac{\partial \mathbf{u}}{\partial t}(\mathbf{x}(\mathbf{X}_i, t), t) + \nabla \mathbf{u}(\mathbf{x}(\mathbf{X}_i, t), t) \cdot \frac{\partial \mathbf{x}}{\partial t}(\mathbf{X}_i, t) \\ &= \frac{D\mathbf{u}}{Dt}(\mathbf{x}(\mathbf{X}_i, t), t) = \frac{D\mathbf{u}_i}{Dt}. \end{aligned}$$

By rearranging terms and dividing by  $\Delta t$  one finds the desired expression.  $\square$

Using Lemma 2, we again discretize (2.48) by

$$\mathbf{u}^m - \mathbf{u}^{m-1} = \Delta t \mathbf{v}^m, \quad (2.79)$$

where we again use Implicit Euler. Joining equations (2.78) and (2.79) into a single system yields

$$\begin{pmatrix} \widetilde{M}^m + \Delta t S_{\text{vis}}^m & \Delta t S_{\text{el}}^m \\ -\Delta t I & I \end{pmatrix} \begin{pmatrix} \mathbf{v}^m \\ \mathbf{u}^m \end{pmatrix} = \begin{pmatrix} \widetilde{M}^{m-1} & 0 \\ 0 & I \end{pmatrix} \begin{pmatrix} \mathbf{v}^{m-1} \\ \mathbf{u}^{m-1} \end{pmatrix} + \Delta t \begin{pmatrix} \mathbf{g}^m \\ \mathbf{0} \end{pmatrix}. \quad (2.80)$$

Since the grid at time level  $m$  depends on  $\mathbf{u}^m$ , we cannot solve (2.80) directly for  $\mathbf{u}^m$  and  $\mathbf{v}^m$ . Instead, we use an iterative scheme. Note that the system in (2.80) can be abstractly written as

$$A(z)z = c + b(z), \quad (2.81)$$

where  $z$  is the unknown vector,  $A(z)$  is the left-hand side matrix,  $c$  is a constant vector corresponding to the solution at time  $t$ , and  $b(z)$  is the right-hand side vector containing  $\mathbf{g}^m$ . We approximate the solution to (2.81) by generating a sequence  $(z_k)$  using the following update rule:

$$A(z_k)z_{k+1} = c + b(z_k). \quad (2.82)$$

Where  $z_0$  is equal to the unknowns at time level  $m - 1$ . If we translate this scheme back to the original variables, we get

$$\begin{pmatrix} \widetilde{M}_k + \Delta t (S_{\text{vis}})_k & \Delta t (S_{\text{el}})_k \\ -\Delta t I & I \end{pmatrix} \begin{pmatrix} \mathbf{v}_{k+1} \\ \mathbf{u}_{k+1} \end{pmatrix} = \begin{pmatrix} \widetilde{M}^{m-1} & 0 \\ 0 & I \end{pmatrix} \begin{pmatrix} \mathbf{v}^{m-1} \\ \mathbf{u}^{m-1} \end{pmatrix} + \Delta t \begin{pmatrix} \mathbf{g}^m \\ \mathbf{0} \end{pmatrix}. \quad (2.83)$$

Where the subscript  $k$  means the matrix or vector has been calculated on the grid based on  $\mathbf{u}_k$ . In words, one iteration of the scheme amounts to:

- Create new grid based on  $\mathbf{u}_k$ .
- Calculate matrices  $\widetilde{M}_k$ ,  $(S_{\text{vis}})_k$ ,  $(S_{\text{el}})_k$  and vector  $\mathbf{g}_k$  on the new grid.
- Solve (2.83) for  $\mathbf{u}_{k+1}$  (and  $\mathbf{v}_{k+1}$ ).

If the iterative process converges, we have found a fixed point of (2.83), and thus a solution to (2.80). In practice, we say that the scheme has converged if the relative difference between two consecutive iterates is sufficiently small. That is, if

$$\frac{\|\mathbf{u}_{k+1} - \mathbf{u}_k\|}{\|\mathbf{u}_k\|} + \frac{\|\mathbf{v}_{k+1} - \mathbf{v}_k\|}{\|\mathbf{v}_k\|} < \delta, \quad (2.84)$$

for some small  $\delta > 0$ .

### 2.3.6 Numerical Experiments

We run the Finite Element model using the parameters  $\rho = 1$ ,  $\mu_\varepsilon = \lambda_\varepsilon = \mu_v = \lambda_v = 1$ . The initial domain is the unit square:  $\Omega(0) = (0, 1)^2$ . We apply the body force  $\mathbf{g} = (0, -0.1)^\top$ . On the left border ( $x = 0$ ) a homogeneous Dirichlet boundary condition is imposed; this border is fixed. On the remaining boundaries, a homogeneous Neumann boundary condition is imposed, corresponding to  $\boldsymbol{\tau} = \mathbf{0}$ . The solution is expected to converge to the stationary solution, which is shown in Figure 2.5. Note that the stationary solution has been computed using the elastic Finite Element model from section 2.2.

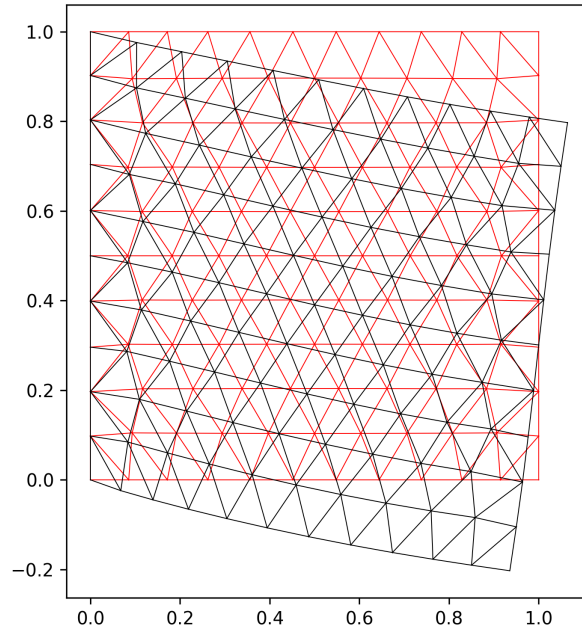


Figure 2.5: Stationary solution. A meshwidth of  $h = 0.1$  has been used.

We solve the viscoelastic system from  $t = 0$  to  $t = 8$  using a stepsize of  $\Delta t = 0.05$ . The same (initial) meshwidth of  $h = 0.1$  is used. Figures 2.6a - 2.6d show the intermediate finite element solutions. Note that in Figure 2.6c the mesh has passed the equilibrium state. Due to the elastic nature of the material, it will again bounce back, as can be seen in Figure 2.6d. This is a result incorporating the viscous term in the force-balance equation.

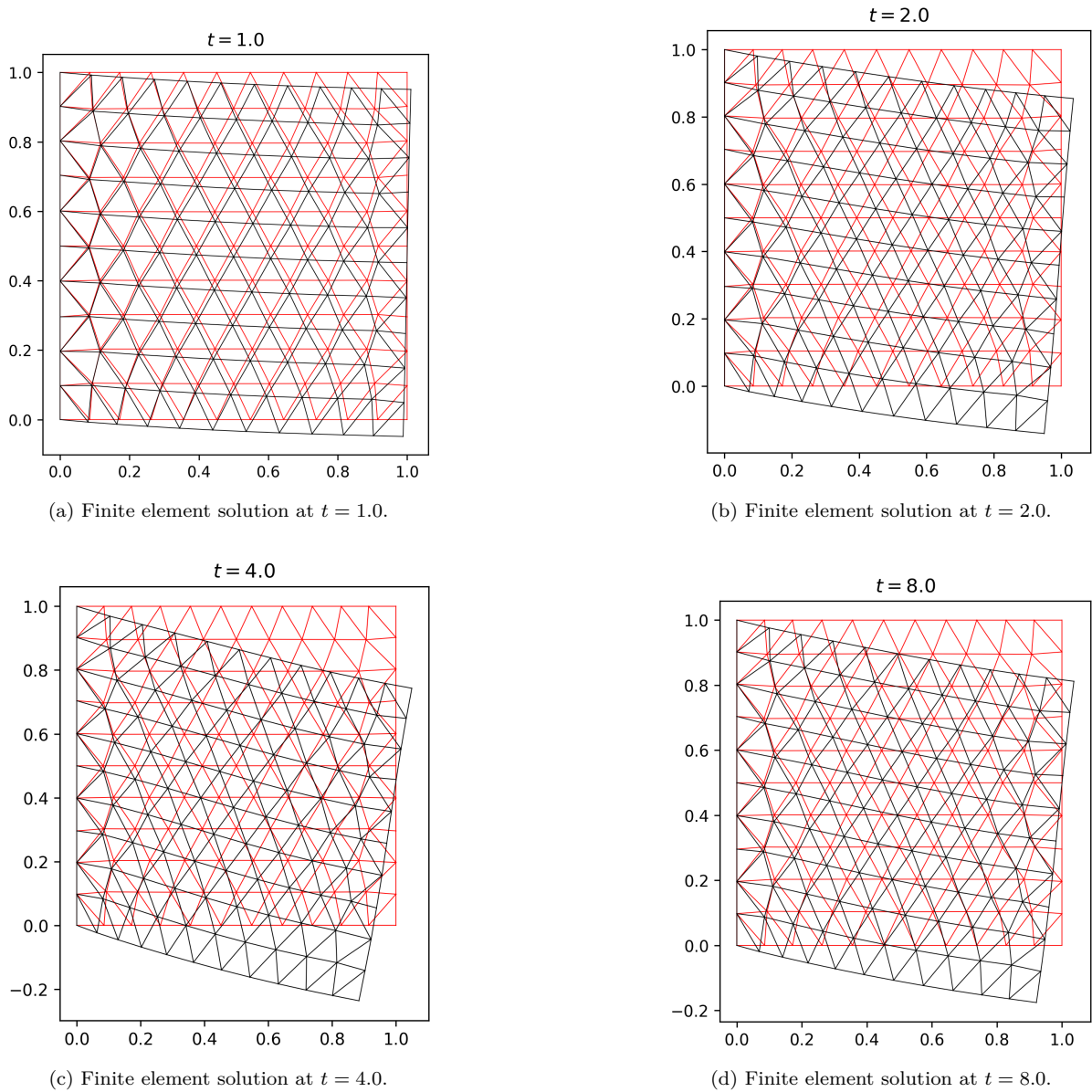


Figure 2.6: Plots showing the finite element solution at various points in time.

## 2.4 Morphoelasticity

Using morphoelasticity allows us to model elastic growth. The idea as introduced in [12] is to decompose the deformation tensor  $\mathbf{F}$  as defined in (2.45) into an elastic part  $\mathbf{F}^e$  and plastic part  $\mathbf{F}^p$ :

$$\mathbf{F} = \mathbf{F}^e \mathbf{F}^p. \quad (2.85)$$

An intuitive way to look at this decomposition is to consider the corresponding mappings that map coordinates onto each other. Between the initial state  $\mathbf{X}$  and current state  $\mathbf{x}$ , a stress-free state  $\mathbf{z}$  is introduced. In the stress free state, only the plastic deformation is applied. The decomposition can then be represented as

$$\mathbf{F} : \mathbf{X} \xrightarrow{\mathbf{F}^p} \mathbf{z} \xrightarrow{\mathbf{F}^e} \mathbf{x}. \quad (2.86)$$

The stress-free state is a virtual configuration; the mapping  $\mathbf{F}^p$  need not even be continuous. In [8], the theory of morphoelasticity is used to obtain an evolution equation for the strain tensor  $\boldsymbol{\varepsilon}$ . It is important to note that the structure of the problem now becomes completely different compared to viscoelasticity. In the latter framework,  $\boldsymbol{\varepsilon}$  directly related to the displacement  $\mathbf{u}$ , allowing us to construct a model that computes  $\mathbf{u}$  directly. In the framework of morphoelasticity,  $\mathbf{u}$  will no longer be explicitly featured in any



of the equations. Instead, the model only allows us to compute  $\boldsymbol{\varepsilon}$  and the displacement velocity  $\mathbf{v}$ , from which  $\mathbf{u}$  can then be determined in a post-processing step. We will not reproduce the derivation of the strain evolution equation from [8], as it involves many pages of algebraic manipulation. Together with (2.36), we obtain the following morphoelastic system of equations:

$$\rho \left( \frac{D\mathbf{v}}{Dt} + (\nabla \cdot \mathbf{v})\mathbf{v} \right) - \nabla \cdot \boldsymbol{\sigma}(\mathbf{v}, \boldsymbol{\varepsilon}) = \mathbf{g} \quad (2.87)$$

$$\frac{D\boldsymbol{\varepsilon}}{Dt} + \boldsymbol{\varepsilon} \text{skw}(\nabla \mathbf{v}) - \text{skw}(\nabla \mathbf{v})\boldsymbol{\varepsilon} + (\text{tr}(\boldsymbol{\varepsilon}) - 1)\text{sym}(\nabla \mathbf{v}) = -\mathbf{G}. \quad (2.88)$$

The stress tensor is again given by an elastic and viscous part such as in (2.37), but the elastic part now depends directly on  $\boldsymbol{\varepsilon}$  instead of  $\mathbf{u}$ :

$$\boldsymbol{\sigma}_{\text{el}}(\boldsymbol{\varepsilon}) = \mu_{\varepsilon}\boldsymbol{\varepsilon} + \lambda_{\varepsilon}\text{tr}(\boldsymbol{\varepsilon})\mathbf{I}. \quad (2.89)$$

Note that to obtain the updated locations of the grid nodes, one must still compute  $\mathbf{u}$ . It can be extracted from  $\mathbf{v}$ . Like in the previous section, part of the boundary is fixed by setting  $\mathbf{v} = \mathbf{0}$  on  $\Gamma_1(t)$ , and on the other part the shear force is prescribed:

$$\boldsymbol{\sigma} \cdot \mathbf{n} = \boldsymbol{\tau}, \quad \text{on } \Gamma_2(t). \quad (2.90)$$

To reduce the number of unknowns by one, we shall first prove the following lemma.

**Lemma 3** (Daan Smits [16]). *If  $\mathbf{G}$  is symmetric for all  $t \geq 0$  and  $\boldsymbol{\varepsilon}$  is symmetric at  $t = 0$ , then  $\boldsymbol{\varepsilon}$  is symmetric for all  $t \geq 0$ .*

*Proof.* Consider the transpose of equation (2.88):

$$\frac{D\boldsymbol{\varepsilon}^{\top}}{Dt} + \text{skw}(\nabla \mathbf{v})^{\top} \boldsymbol{\varepsilon}^{\top} - \boldsymbol{\varepsilon}^{\top} \text{skw}(\nabla \mathbf{v})^{\top} + (\text{tr}(\boldsymbol{\varepsilon}) - 1)\text{sym}(\nabla \mathbf{v})^{\top} = -\mathbf{G}^{\top}.$$

Using the facts that  $\mathbf{G}^{\top} = \mathbf{G}$ ,  $\text{skw}(\nabla \mathbf{v})^{\top} = -\text{skw}(\nabla \mathbf{v})$  and  $\text{sym}(\nabla \mathbf{v})^{\top} = \text{sym}(\nabla \mathbf{v})$ , we get

$$\frac{D\boldsymbol{\varepsilon}^{\top}}{Dt} + \boldsymbol{\varepsilon}^{\top} \text{skw}(\nabla \mathbf{v}) - \text{skw}(\nabla \mathbf{v})\boldsymbol{\varepsilon}^{\top} + (\text{tr}(\boldsymbol{\varepsilon}) - 1)\text{sym}(\nabla \mathbf{v}) = -\mathbf{G}.$$

Subtract the above equation from (2.88) to obtain

$$\frac{D(\boldsymbol{\varepsilon} - \boldsymbol{\varepsilon}^{\top})}{Dt} + (\boldsymbol{\varepsilon} - \boldsymbol{\varepsilon}^{\top})\text{skw}(\nabla \mathbf{v}) - \text{skw}(\nabla \mathbf{v})(\boldsymbol{\varepsilon} - \boldsymbol{\varepsilon}^{\top}) = \mathbf{0}.$$

Thus, whenever  $\boldsymbol{\varepsilon}$  is symmetric, we have

$$\frac{D(\boldsymbol{\varepsilon} - \boldsymbol{\varepsilon}^{\top})}{Dt} = \mathbf{0}.$$

Since  $\boldsymbol{\varepsilon}$  is symmetric at  $t = 0$ , it will be symmetric for all  $t \geq 0$ .  $\square$

In the numerical experiments, we will set  $\mathbf{G} = \alpha\boldsymbol{\varepsilon}$ . In this case Lemma 3 holds whenever  $\boldsymbol{\varepsilon}$  is symmetric at  $t = 0$ . We always set  $\boldsymbol{\varepsilon}|_{t=0} = \mathbf{0}$ . From now on, we will assume the conditions of Lemma 3 are met. Since  $\varepsilon^{12} = \varepsilon^{21}$ , we are left with the 5 unknowns  $v^1, v^2$  and  $\varepsilon^{11}, \varepsilon^{12}, \varepsilon^{22}$ .

## 2.4.1 Weak Formulation

### Velocity Evolution Equation

We first consider the weak form of (2.87). Note that we can repeat the steps in section 2.3.3 up until equation (2.60). We must now however split the stress-tensor into a viscous part and strain part to obtain the following weak form:

$$\rho \frac{d}{dt} \int_{\Omega(t)} \mathbf{v} \cdot \boldsymbol{\varphi} \, d\Omega + \int_{\Omega(t)} (\boldsymbol{\sigma}_{\text{vis}}(\mathbf{v}) + \boldsymbol{\sigma}_{\text{el}}(\boldsymbol{\varepsilon})) : \nabla \boldsymbol{\varphi} \, d\Omega = \int_{\Omega(t)} \mathbf{g} \cdot \boldsymbol{\varphi} \, d\Omega + \int_{\Gamma(t)} \boldsymbol{\tau} \cdot \boldsymbol{\varphi} \, d\Omega. \quad (2.91)$$

The integral of  $\boldsymbol{\sigma}_{\text{vis}}(\mathbf{v}) : \nabla \boldsymbol{\varphi}$  will lead to the same Galerkin equations as in section 2.3.4. Therefore we will focus on the part involving  $\boldsymbol{\varepsilon}$ . Let  $\boldsymbol{\varphi}$  be a scalar-valued Lagrangian basis function. If we first let  $\boldsymbol{\varphi} = (\boldsymbol{\varphi}, 0)^{\top}$ , then

$$\boldsymbol{\sigma}_{\text{el}}(\boldsymbol{\varepsilon}) : \nabla \boldsymbol{\varphi} = ((\mu_{\varepsilon} + \lambda_{\varepsilon})\varepsilon^{11} + \lambda_{\varepsilon}\varepsilon^{22}) \frac{\partial \boldsymbol{\varphi}}{\partial x} + \mu_{\varepsilon}\varepsilon^{12} \frac{\partial \boldsymbol{\varphi}}{\partial y}. \quad (2.92)$$

Now let  $\boldsymbol{\varphi} = (0, \varphi)^\top$ , then:

$$\boldsymbol{\sigma}_{\text{el}}(\boldsymbol{\varepsilon}) : \nabla \boldsymbol{\varphi} = \mu_\varepsilon \varepsilon^{12} \frac{\partial \varphi}{\partial x} + (\lambda_\varepsilon \varepsilon^{11} + (\mu_\varepsilon + \lambda_\varepsilon) \varepsilon^{22}) \frac{\partial \varphi}{\partial y}. \quad (2.93)$$

By considering these expressions and using (2.71) - (2.74), we get the following weak form for the equation for  $v_1$ :

$$\begin{aligned} \rho \frac{d}{dt} \int_{\Omega(t)} v^1 \varphi \, d\Omega + \int_{\Omega(t)} \left( (\mu_v + \lambda_v) \frac{\partial v^1}{\partial x} + \lambda_v \frac{\partial v^2}{\partial y} + (\mu_\varepsilon + \lambda_\varepsilon) \varepsilon^{11} + \lambda_\varepsilon \varepsilon^{22} \right) \frac{\partial \varphi}{\partial x} \\ + \left( \frac{1}{2} \mu_v \frac{\partial v^1}{\partial y} + \frac{1}{2} \mu_v \frac{\partial v^2}{\partial x} + \mu_\varepsilon \varepsilon^{12} \right) \frac{\partial \varphi}{\partial y} \, d\Omega = \int_{\Omega(t)} g^1 \varphi \, d\Omega + \int_{\Gamma(t)} \tau^1 \varphi \, d\Gamma, \end{aligned} \quad (2.94)$$

and for the equation for  $v_2$ :

$$\begin{aligned} \rho \frac{d}{dt} \int_{\Omega(t)} v^2 \varphi \, d\Omega + \int_{\Omega(t)} \left( \frac{1}{2} \mu_v \frac{\partial v^1}{\partial y} + \frac{1}{2} \mu_v \frac{\partial v^2}{\partial x} + \mu_\varepsilon \varepsilon^{12} \right) \frac{\partial \varphi}{\partial x} \\ + \left( \lambda_v \frac{\partial v^1}{\partial x} + (\mu_v + \lambda_v) \frac{\partial v^2}{\partial y} + \lambda_\varepsilon \varepsilon^{11} + (\mu_\varepsilon + \lambda_\varepsilon) \varepsilon^{22} \right) \frac{\partial \varphi}{\partial y} \, d\Omega = \int_{\Omega(t)} g^2 \varphi \, d\Omega + \int_{\Gamma(t)} \tau^2 \varphi \, d\Gamma, \end{aligned} \quad (2.95)$$

### Strain Evolution Equation

We write out the strain evolution equation (2.88) component-wise:

$$\frac{D\varepsilon^{11}}{Dt} + \varepsilon^{21} \left( \frac{\partial v^2}{\partial x} - \frac{\partial v^1}{\partial y} \right) + (\varepsilon^{11} + \varepsilon^{22} - 1) \frac{\partial v^1}{\partial x} = -G^{11}, \quad (2.96)$$

$$\frac{D\varepsilon^{12}}{Dt} + \frac{1}{2} (\varepsilon^{11} - \varepsilon^{22}) \left( \frac{\partial v^1}{\partial y} - \frac{\partial v^2}{\partial x} \right) + \frac{1}{2} (\varepsilon^{11} + \varepsilon^{22} - 1) \left( \frac{\partial v^2}{\partial x} + \frac{\partial v^1}{\partial y} \right) = -G^{12}, \quad (2.97)$$

$$\frac{D\varepsilon^{22}}{Dt} + \varepsilon^{21} \left( \frac{\partial v^1}{\partial y} - \frac{\partial v^2}{\partial x} \right) + (\varepsilon^{11} + \varepsilon^{22} - 1) \frac{\partial v^2}{\partial y} = -G^{22}. \quad (2.98)$$

Before we multiply these equations by a test function and integrate them over  $\Omega(t)$ , we add a term involving  $\nabla \cdot \mathbf{v}$  to both sides. This is done in order to use Theorem 3. After rearranging some terms, we find that (2.96) - (2.98) are equivalent to:

$$\left( \frac{D\varepsilon^{11}}{Dt} + (\nabla \cdot \mathbf{v}) \varepsilon^{11} \right) - \varepsilon^{11} \frac{\partial v^2}{\partial y} + \varepsilon^{12} \left( \frac{\partial v^2}{\partial x} - \frac{\partial v^1}{\partial y} \right) + \varepsilon^{22} \frac{\partial v^1}{\partial x} = \frac{\partial v^1}{\partial x} - G^{11}, \quad (2.99)$$

$$\left( \frac{D\varepsilon^{12}}{Dt} + (\nabla \cdot \mathbf{v}) \varepsilon^{12} \right) + \varepsilon^{11} \frac{\partial v^1}{\partial y} - \varepsilon^{12} \left( \frac{\partial v^1}{\partial x} + \frac{\partial v^2}{\partial y} \right) + \varepsilon^{22} \frac{\partial v^2}{\partial x} = \frac{1}{2} \left( \frac{\partial v^2}{\partial x} + \frac{\partial v^1}{\partial y} \right) - G^{12}, \quad (2.100)$$

$$\left( \frac{D\varepsilon^{22}}{Dt} + (\nabla \cdot \mathbf{v}) \varepsilon^{22} \right) + \varepsilon^{11} \frac{\partial v^2}{\partial y} - \varepsilon^{12} \left( \frac{\partial v^2}{\partial x} - \frac{\partial v^1}{\partial y} \right) - \varepsilon^{22} \frac{\partial v^1}{\partial x} = \frac{\partial v^2}{\partial y} - G^{22}. \quad (2.101)$$

Let  $\varphi$  be a Lagrangian test function, then by Theorem 3 we find the following weak forms:

$$\begin{aligned} \frac{d}{dt} \int_{\Omega(t)} \varepsilon^{11} \varphi \, d\Omega + \int_{\Omega(t)} \left( -\varepsilon^{11} \frac{\partial v^2}{\partial y} + \varepsilon^{12} \left( \frac{\partial v^2}{\partial x} - \frac{\partial v^1}{\partial y} \right) + \varepsilon^{22} \frac{\partial v^1}{\partial x} \right) \varphi \, d\Omega \\ = \int_{\Omega(t)} \left( \frac{\partial v^1}{\partial x} - G^{11} \right) \varphi \, d\Omega, \end{aligned} \quad (2.102)$$

$$\begin{aligned} \frac{d}{dt} \int_{\Omega(t)} \varepsilon^{12} \varphi \, d\Omega + \int_{\Omega(t)} \left( \varepsilon^{11} \frac{\partial v^1}{\partial y} - \varepsilon^{12} \left( \frac{\partial v^1}{\partial x} + \frac{\partial v^2}{\partial y} \right) + \varepsilon^{22} \frac{\partial v^2}{\partial x} \right) \varphi \, d\Omega \\ = \int_{\Omega(t)} \left( \frac{1}{2} \left( \frac{\partial v^2}{\partial x} + \frac{\partial v^1}{\partial y} \right) - G^{12} \right) \varphi \, d\Omega, \end{aligned} \quad (2.103)$$

$$\begin{aligned} \frac{d}{dt} \int_{\Omega(t)} \varepsilon^{22} \varphi \, d\Omega + \int_{\Omega(t)} \left( \varepsilon^{11} \frac{\partial v^2}{\partial y} - \varepsilon^{12} \left( \frac{\partial v^2}{\partial x} - \frac{\partial v^1}{\partial y} \right) - \varepsilon^{22} \frac{\partial v^1}{\partial x} \right) \varphi \, d\Omega \\ = \int_{\Omega(t)} \left( \frac{\partial v^2}{\partial y} - G^{22} \right) \varphi \, d\Omega, \end{aligned} \quad (2.104)$$

Another way to obtain these equations is to multiply (2.88) with a symmetric tensor-valued test field  $\zeta$ , and integrate over  $\Omega(t)$ . To use Theorem 3, first add and subtract  $(\nabla \cdot \mathbf{v})\boldsymbol{\varepsilon}$ . We then get the weak form

$$\begin{aligned} & \int_{\Omega(t)} \left( \frac{D\boldsymbol{\varepsilon}}{Dt} + (\nabla \cdot \mathbf{v})\boldsymbol{\varepsilon} \right) : \zeta \, d\Omega \\ & \quad + \int_{\Omega(t)} \left( \boldsymbol{\varepsilon} \text{skw}(\nabla \mathbf{v}) - \text{skw}(\nabla \mathbf{v})\boldsymbol{\varepsilon} + (\text{tr}(\boldsymbol{\varepsilon}) - 1)\text{sym}(\nabla \mathbf{v}) - (\nabla \cdot \mathbf{v})\boldsymbol{\varepsilon} \right) : \zeta \, d\Omega \\ & \hspace{20em} = - \int_{\Omega(t)} \mathbf{G} : \zeta \, d\Omega. \end{aligned} \quad (2.105)$$

Applying Theorem 3 yields

$$\begin{aligned} & \frac{d}{dt} \int_{\Omega(t)} \boldsymbol{\varepsilon} : \zeta \, d\Omega \\ & \quad + \int_{\Omega(t)} \left( \boldsymbol{\varepsilon} \text{skw}(\nabla \mathbf{v}) - \text{skw}(\nabla \mathbf{v})\boldsymbol{\varepsilon} + (\text{tr}(\boldsymbol{\varepsilon}) - 1)\text{sym}(\nabla \mathbf{v}) - (\nabla \cdot \mathbf{v})\boldsymbol{\varepsilon} \right) : \zeta \, d\Omega \\ & \hspace{20em} = - \int_{\Omega(t)} \mathbf{G} : \zeta \, d\Omega. \end{aligned} \quad (2.106)$$

This step is valid provided  $\zeta$  is a Lagrangian test field. Setting  $\zeta$  equal to

$$\begin{pmatrix} \varphi & 0 \\ 0 & 0 \end{pmatrix}, \quad \begin{pmatrix} 0 & \varphi \\ 0 & 0 \end{pmatrix}, \quad \text{and} \quad \begin{pmatrix} 0 & 0 \\ 0 & \varphi \end{pmatrix},$$

yields equations (2.102), (2.103) and (2.104) respectively. To further motivate this, the tensor in the middle term of equation (2.106) is computed component-wise:

$$\begin{aligned} & \boldsymbol{\varepsilon} \text{skw}(\nabla \mathbf{v}) - \text{skw}(\nabla \mathbf{v})\boldsymbol{\varepsilon} + (\text{tr}(\boldsymbol{\varepsilon}) - 1)\text{sym}(\nabla \mathbf{v}) - (\nabla \cdot \mathbf{v})\boldsymbol{\varepsilon} \\ & = \begin{pmatrix} -\varepsilon^{11} \frac{\partial v^2}{\partial y} + \varepsilon^{12} \left( \frac{\partial v^2}{\partial x} - \frac{\partial v^1}{\partial y} \right) + \varepsilon^{22} \frac{\partial v^1}{\partial x} & \varepsilon^{11} \frac{\partial v^1}{\partial y} - \varepsilon^{12} \left( \frac{\partial v_1}{\partial x} + \frac{\partial v^2}{\partial y} \right) + \varepsilon^{22} \frac{\partial v^2}{\partial x} \\ \varepsilon^{11} \frac{\partial v^1}{\partial y} - \varepsilon^{12} \left( \frac{\partial v^1}{\partial x} + \frac{\partial v^2}{\partial y} \right) + \varepsilon^{22} \frac{\partial v^2}{\partial x} & \varepsilon^{11} \frac{\partial v^2}{\partial y} - \varepsilon^{12} \left( \frac{\partial v^2}{\partial x} - \frac{\partial v^1}{\partial y} \right) - \varepsilon^{22} \frac{\partial v^1}{\partial x} \end{pmatrix} \\ & \quad - \begin{pmatrix} \frac{\partial v^1}{\partial x} & \frac{1}{2} \left( \frac{\partial v^1}{\partial y} + \frac{\partial v^2}{\partial x} \right) \\ \frac{1}{2} \left( \frac{\partial v^1}{\partial y} + \frac{\partial v^2}{\partial x} \right) & \frac{\partial v^2}{\partial y} \end{pmatrix}. \end{aligned} \quad (2.107)$$

One can then immediately see that taking the double dot product with one of the test tensors above leads to equations (2.102) - (2.104).

### 2.4.2 Galerkin Equations

We derive the Galerkin equations corresponding to weak forms (2.94), (2.95) and (2.102), (2.103), (2.104). Let  $\varphi_j$  be the linear basis function centered in grid node  $j$ . The finite element solutions are then of the form

$$v^1 = \sum_{j=1}^n v_j^1 \varphi_j, \quad v^2 = \sum_{j=1}^n v_j^2 \varphi_j, \quad \varepsilon^{11} = \sum_{j=1}^n \varepsilon_j^{11} \varphi_j, \quad \varepsilon^{12} = \sum_{j=1}^n \varepsilon_j^{12} \varphi_j, \quad \varepsilon^{22} = \sum_{j=1}^n \varepsilon_j^{22} \varphi_j, \quad (2.108)$$

The coefficients  $v_j^1, \dots, \varepsilon_j^{22}$  depend only on time, whereas the basis functions depend on space and time;  $\varphi_j = \varphi_j(\mathbf{x}; t)$ . Substitute the finite element solutions into the weak forms, and set  $\varphi = \varphi_i$  for  $i = 1, \dots, n$ . Then each weak form will yield  $n$  equations. Note that the equations derived from (2.102), (2.103) and (2.104) will contain a non-linear part, due to the products between  $\varepsilon$ -components and spatial derivatives

of the components of  $\mathbf{v}$ . The resulting non-linear system is given in block matrix form by:

$$\begin{aligned} \frac{d}{dt} \begin{pmatrix} \rho M \mathbf{v}^1 \\ \rho M \mathbf{v}^2 \\ M \boldsymbol{\varepsilon}^{11} \\ M \boldsymbol{\varepsilon}^{12} \\ M \boldsymbol{\varepsilon}^{22} \end{pmatrix} + \begin{pmatrix} S_{\text{vis}}^{11} & S_{\text{vis}}^{12} & (\mu_\varepsilon + \lambda_\varepsilon) B_x & \mu_\varepsilon B_y & \lambda_\varepsilon B_x \\ S_{\text{vis}}^{21} & S_{\text{vis}}^{22} & \lambda_\varepsilon B_y & \mu_\varepsilon B_x & (\mu_\varepsilon + \lambda_\varepsilon) B_y \\ -B_x^\top & 0 & 0 & 0 & 0 \\ -\frac{1}{2} B_y^\top & -\frac{1}{2} B_x^\top & 0 & 0 & 0 \\ 0 & -B_y^\top & 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} \mathbf{v}^1 \\ \mathbf{v}^2 \\ \boldsymbol{\varepsilon}^{11} \\ \boldsymbol{\varepsilon}^{12} \\ \boldsymbol{\varepsilon}^{22} \end{pmatrix} + \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \\ \mathbf{N}^{11}(\mathbf{z}) \\ \mathbf{N}^{12}(\mathbf{z}) \\ -\mathbf{N}^{11}(\mathbf{z}) \end{pmatrix} \\ = \begin{pmatrix} \mathbf{g}^1 + \boldsymbol{\tau}^1 \\ \mathbf{g}^2 + \boldsymbol{\tau}^2 \\ -\mathbf{G}^{11} \\ -\mathbf{G}^{12} \\ -\mathbf{G}^{22} \end{pmatrix}. \end{aligned} \quad (2.109)$$

Here  $\mathbf{v}^1 = (v_1^1, \dots, v_n^1)^\top$  etc, and the vector  $\mathbf{z}$  contains all the unknowns:  $\mathbf{z} = (\mathbf{v}^1, \mathbf{v}^2, \boldsymbol{\varepsilon}^{11}, \boldsymbol{\varepsilon}^{12}, \boldsymbol{\varepsilon}^{22})^\top$ . The components of the matrices  $S_{\text{vis}}^{11}$ ,  $S_{\text{vis}}^{12}$ ,  $S_{\text{vis}}^{21}$  and  $S_{\text{vis}}^{22}$  are given by (2.71), (2.72), (2.73), and (2.74) respectively. Similarly, the right-hand side vectors  $\mathbf{g}^1$  and  $\mathbf{g}^2$  are given by (2.65) and (2.66) respectively. The mass matrix  $M$  is again given by:

$$M_{ij} = \int_{\Omega(t)} \varphi_i \varphi_j \, d\Omega. \quad (2.110)$$

The components of the matrices  $B_x$  and  $B_y$  are given by:

$$(B_x)_{ij} = \int_{\Omega(t)} \varphi_j \frac{\partial \varphi_i}{\partial x} \, d\Omega, \quad (2.111)$$

$$(B_y)_{ij} = \int_{\Omega(t)} \varphi_j \frac{\partial \varphi_i}{\partial y} \, d\Omega. \quad (2.112)$$

The non-linear functionals  $\mathbf{N}^{11}$  and  $\mathbf{N}^{12}$  are defined component-wise by:

$$N_i^{11}(\mathbf{z}) = \sum_{j=1}^n \sum_{k=1}^n \left( -N_{ijk}^y \varepsilon_j^{11} v_k^2 - N_{ijk}^y \varepsilon_j^{12} v_k^1 + N_{ijk}^x \varepsilon_j^{12} v_k^2 + N_{ijk}^x \varepsilon_j^{22} v_k^1 \right), \quad (2.113)$$

$$N_i^{12}(\mathbf{z}) = \sum_{j=1}^n \sum_{k=1}^n \left( N_{ijk}^y \varepsilon_j^{11} v_k^1 - N_{ijk}^x \varepsilon_j^{12} v_k^1 - N_{ijk}^y \varepsilon_j^{12} v_k^2 + N_{ijk}^x \varepsilon_j^{22} v_k^2 \right), \quad (2.114)$$

where

$$N_{ijk}^x = \int_{\Omega(t)} \varphi_i \varphi_j \frac{\partial \varphi_k}{\partial x} \, d\Omega, \quad (2.115)$$

$$N_{ijk}^y = \int_{\Omega(t)} \varphi_i \varphi_j \frac{\partial \varphi_k}{\partial y} \, d\Omega. \quad (2.116)$$

Finally, the right-hand side vectors  $\mathbf{G}^{11}$ ,  $\mathbf{G}^{12}$  and  $\mathbf{G}^{22}$  are given by:

$$G_i^{11} = \int_{\Omega(t)} \varphi_i G^{11} \, d\Omega, \quad G_i^{12} = \int_{\Omega(t)} \varphi_i G^{12} \, d\Omega, \quad G_i^{22} = \int_{\Omega(t)} \varphi_i G^{22} \, d\Omega, \quad i, j, k \in \{1, 2, 3\}. \quad (2.117)$$

The element vectors and matrices of all objects found in system (2.109) are found in Appendix A.2. To model permanent deformations, we let the growth tensor  $\mathbf{G}$  be proportional to the strain tensor  $\boldsymbol{\varepsilon}$ :

$$\mathbf{G} = \alpha \boldsymbol{\varepsilon}, \quad (2.118)$$

for some proportionality constant  $\alpha \geq 0$ . In this case we have  $\mathbf{G}^{11} = \alpha M \boldsymbol{\varepsilon}^{11}$ ,  $\mathbf{G}^{12} = \alpha M \boldsymbol{\varepsilon}^{12}$  and  $\mathbf{G}^{22} = \alpha M \boldsymbol{\varepsilon}^{22}$ . Thus, system (2.109) becomes:

$$\begin{aligned} \frac{d}{dt} \begin{pmatrix} \rho M \mathbf{v}^1 \\ \rho M \mathbf{v}^2 \\ M \boldsymbol{\varepsilon}^{11} \\ M \boldsymbol{\varepsilon}^{12} \\ M \boldsymbol{\varepsilon}^{22} \end{pmatrix} + \begin{pmatrix} S_{\text{vis}}^{11} & S_{\text{vis}}^{12} & (\mu_\varepsilon + \lambda_\varepsilon) B_x & \mu_\varepsilon B_y & \lambda_\varepsilon B_x \\ S_{\text{vis}}^{21} & S_{\text{vis}}^{22} & \lambda_\varepsilon B_y & \mu_\varepsilon B_x & (\mu_\varepsilon + \lambda_\varepsilon) B_y \\ -B_x^\top & 0 & \alpha M & 0 & 0 \\ -\frac{1}{2} B_y^\top & -\frac{1}{2} B_x^\top & 0 & \alpha M & 0 \\ 0 & -B_y^\top & 0 & 0 & \alpha M \end{pmatrix} \begin{pmatrix} \mathbf{v}^1 \\ \mathbf{v}^2 \\ \boldsymbol{\varepsilon}^{11} \\ \boldsymbol{\varepsilon}^{12} \\ \boldsymbol{\varepsilon}^{22} \end{pmatrix} + \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \\ \mathbf{N}^{11}(\mathbf{z}) \\ \mathbf{N}^{12}(\mathbf{z}) \\ -\mathbf{N}^{11}(\mathbf{z}) \end{pmatrix} \\ = \begin{pmatrix} \mathbf{g}^1 + \boldsymbol{\tau}^1 \\ \mathbf{g}^2 + \boldsymbol{\tau}^2 \\ -\mathbf{G}^{11} \\ -\mathbf{G}^{12} \\ -\mathbf{G}^{22} \end{pmatrix}. \end{aligned} \quad (2.119)$$

It will become useful to write the above system in a more compact form. To this end, we define the block vectors  $\mathbf{v} = (\mathbf{v}^1, \mathbf{v}^2)^\top$  and  $\boldsymbol{\varepsilon} = (\boldsymbol{\varepsilon}^{11}, \boldsymbol{\varepsilon}^{21}, \boldsymbol{\varepsilon}^{22})^\top$ . Then system (2.119) can be written as

$$\frac{d}{dt} \begin{pmatrix} \rho M_v \mathbf{v} \\ M_\varepsilon \boldsymbol{\varepsilon} \end{pmatrix} + \begin{pmatrix} S_{\text{vis}} & S_{\text{el}} \\ -B^\top & \alpha M_\varepsilon \end{pmatrix} \begin{pmatrix} \mathbf{v} \\ \boldsymbol{\varepsilon} \end{pmatrix} + \begin{pmatrix} \mathbf{0} \\ \mathbf{N}(\mathbf{v}, \boldsymbol{\varepsilon}) \end{pmatrix} = \begin{pmatrix} \mathbf{g} + \boldsymbol{\tau} \end{pmatrix}. \quad (2.120)$$

Here  $M_v$  and  $M_\varepsilon$  the block matrices with respectively 2 and 3 times the mass matrix  $M$  on its block-diagonal. The other block matrices and vectors should be clear from the structure of the system when compared to (2.119). The per-element contributions, or element matrices/vectors, of all objects found in this section are computed in A.2.

### 2.4.3 Time Discretization

We use the following short notation to denote (2.109) or (2.119):

$$\frac{d}{dt} (\widetilde{M} \mathbf{z}) + A \mathbf{z} + \mathbf{N}(\mathbf{z}) = \mathbf{f}. \quad (2.121)$$

Note that  $\widetilde{M}$ ,  $A$ ,  $\mathbf{N}(\cdot)$  and  $\mathbf{f}$  all depend on  $\mathbf{z}$  implicitly, because they involve integrations over the current domain  $\Omega(t)$ . We use the implicit Euler timestepping method. Applying this method to (2.121) yields the following update equation:

$$(\widetilde{M}^m + \Delta t A^m) \mathbf{z}^m + \Delta t \mathbf{N}^m(\mathbf{z}^m) = \widetilde{M}^{m-1} \mathbf{z}^{m-1} + \Delta t \mathbf{f}^m. \quad (2.122)$$

The superscript denotes a time-level. In the case of the coefficient matrices and right-hand side vector, the time-level is shown explicitly to emphasize that these objects depend on the grid. Note that the grid at time-level  $m$  can be computed using  $\mathbf{v}^m$  like in Section 2.3.5. Given  $\mathbf{z}^{m-1}$ , to obtain  $\mathbf{z}^m$ , we must solve an equation of the form

$$Q[\zeta] \zeta + q[\zeta](\zeta) = \widetilde{f}[\zeta] \quad (2.123)$$

for  $\zeta = \mathbf{z}^m$ . This equation is an even more abstract form of (2.122). Due to the non-linearities (concerning both the grid update and the non-linear part), it must be solved using a Picard-type iterative method. We consider the following update rule:

$$Q[\zeta^\ell] \zeta^{\ell+1} + q[\zeta^\ell](\zeta^\ell) = \widetilde{f}[\zeta^\ell], \quad (2.124)$$

where  $\zeta^0 = \mathbf{z}^{m-1}$ . Note that  $\zeta^{\ell+1}$  can be obtained by solving a linear system. Using this method, we only need a single Picard-loop in each timestep. The algorithm is presented in pseudocode in Algorithm 1.

---

#### Algorithm 1 Single Loop Timestep

---

```

set  $\zeta^0 = \mathbf{z}^t$ 
for  $\ell = 0, 1, \dots$  do
   $\zeta^{\ell+1} = Q[\zeta^\ell]^{-1}(\widetilde{f}[\zeta^\ell] - q[\zeta^\ell](\zeta^\ell))$ 
  if  $\|\zeta^{\ell+1} - \zeta^\ell\| < \delta$  then
    set  $\mathbf{z}^{t+\Delta t} = \zeta^{\ell+1}$ 
    break for-loop
  end if
end for

```

---

Another option would be to use a different update rule, such as

$$Q[\zeta^\ell]\zeta^{\ell+1} + q[\zeta^\ell](\zeta^{\ell+1}) = \tilde{f}[\zeta^\ell]. \quad (2.125)$$

Now, in order to calculate  $\zeta^{\ell+1}$ , one must invert the non-linear operator  $q[\zeta^\ell](\cdot)$ . As this is not feasible in practice, a nested iterative scheme must be used to obtain  $\zeta^{\ell+1}$ . Such a nested iterative process is presented in Algorithm 2.

---

**Algorithm 2** Nested Loop Timestep
 

---

```

set  $\zeta^0 = \mathbf{z}^t$ 
for  $\ell = 0, 1, \dots$  do
  set  $\zeta^{\ell+1,0} = \zeta^\ell$ 
  for  $m = 0, 1, \dots$  do
     $\zeta^{\ell+1,m+1} = Q[\zeta^\ell]^{-1}(\tilde{f}[\zeta^\ell] - q[\zeta^\ell](\zeta^{\ell+1,m}))$ 
    if  $\|\zeta^{\ell+1,m+1} - \zeta^{\ell+1,m}\| < \delta_1$  then
      set  $\zeta^{\ell+1} = \zeta^{\ell+1,m+1}$ 
      break for-loop
    end if
  end for
  if  $\|\zeta^{\ell+1} - \zeta^\ell\| < \delta_2$  then
    set  $\mathbf{z}^{t+\Delta t} = \zeta^{\ell+1}$ 
    break for-loop
  end if
end for

```

---

In practice we see that Algorithm 1 performs well. It should be noted that although an unconditionally stable method is used, we still have a restriction on the timestep. If  $\Delta t$  is too big, then the Picard iterations will not converge.

#### 2.4.4 Numerical Experiments

Let the initial domain be the unit square,  $\Omega(0) = (0, 1)^2$ . The left border ( $x = 0$ ) is fixed in place. A homogeneous Neumann boundary condition is imposed on the other boundaries; thus  $\boldsymbol{\sigma} \cdot \mathbf{n} = \mathbf{0}$ . The body force  $\mathbf{f}$  is set to

$$\mathbf{g}(\mathbf{x}, t) = \begin{pmatrix} F \\ 0 \end{pmatrix} \mathbb{1}_{[0,1]}(t),$$

where  $F = 1/2$ . We also set

$$\rho = 1, \quad \mu_\varepsilon = \lambda_\varepsilon = \mu_v = \lambda_v = 1.$$

Note that this corresponds to setting the Poisson ratio to  $\nu = 1/3$ , hence we expect some amount of transaxial compression as a consequence of axial stretching. We set  $\mathbf{G} = \alpha \boldsymbol{\varepsilon}$  for various values of  $\alpha$ . When  $\alpha = 0$ , we expect there to be no permanent deformations; the system is purely elastic. In this case, the domain should ‘bounce back’ to its original state after  $t = 1$ . For  $\alpha > 0$ , we expect there to be some permanent deformations. Figures 2.7, 2.8 and 2.9 show the finite element solution at  $t = 10$  for  $\alpha = 0$ ,  $\alpha = 1$  and  $\alpha = 10$  respectively.

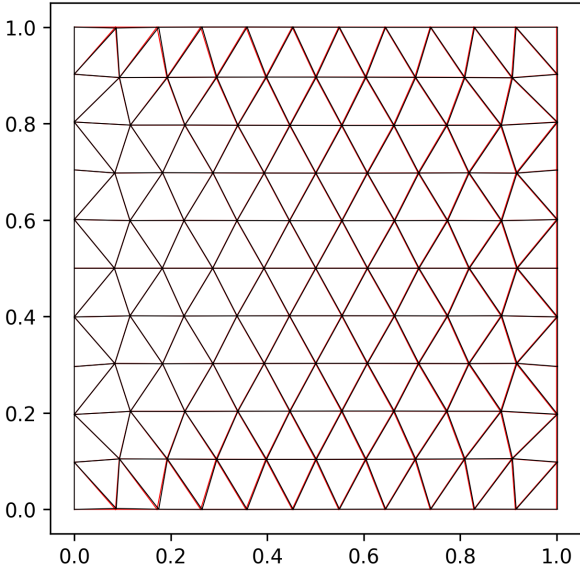


Figure 2.7: finite element solution at  $t = 10$  for  $\alpha = 0$ .

Note that as expected, the domain has bounced back to its original state.

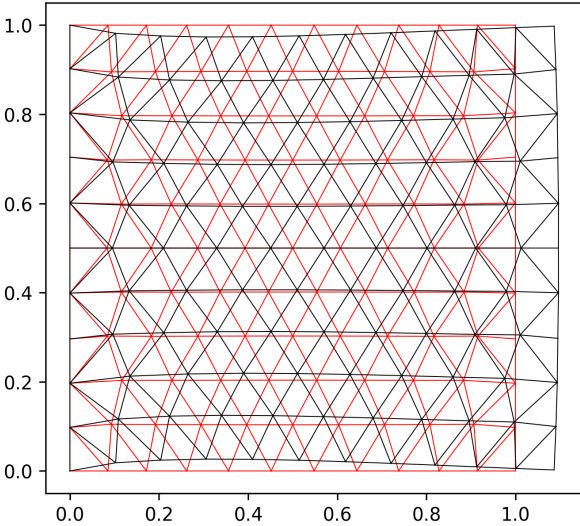
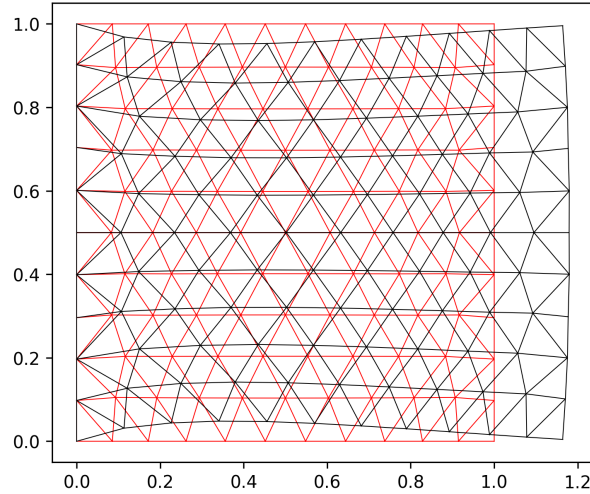


Figure 2.8: finite element solution at  $t = 10$  for  $\alpha = 1$ .

Figure 2.9: finite element solution at  $t = 10$  for  $\alpha = 10$ .

## 2.5 Comparison Between Elastic Models

In this section, we show the qualitative differences between the three elastic models. To this end, a time-dependent force is applied to the domain. In the case of the purely elastic model, equation (2.14) is made time-dependent. That is, since the body force is a function of space and time;  $\mathbf{g} = \mathbf{g}(\mathbf{x}, t)$ , so is the displacement:  $\mathbf{u} = \mathbf{u}(\mathbf{x}, t)$ . For the first comparison, we consider the unit square, and set the body force equal to

$$\mathbf{g}(\mathbf{x}, t) = \begin{pmatrix} 0 \\ -0.1 \sin(\pi t) \mathbb{1}_{[0,1]}(t) \end{pmatrix} \quad (2.126)$$

Thus for  $t > 1$  no force is applied. The left boundary of the domain ( $x = 0$ ) is fixed in place. The shear force on the rest of the boundary is set equal to zero. To compare the different models, we track the  $y$ -coordinate of the bottom right corner of the square. Moreover, a uniform initial grid with maximum meshwidth  $h = 0.05\sqrt{2}$  is used. Figure 2.10 shows this grid, and the position of the tracked point.

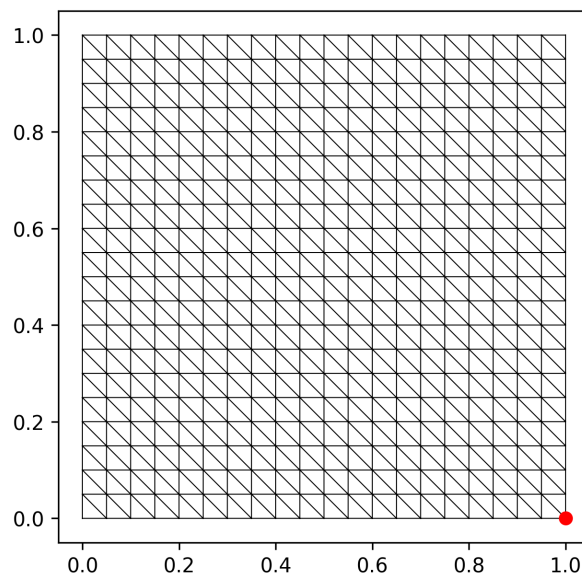


Figure 2.10: Finite element mesh with tracked grid node shown in red.

As the body force is applied to the domain, we expect it to move in a similar fashion to Figure 2.6. In the elastic case, the displacement will react instantaneously to the force, whereas in the visco- and morphoelastic models, the reaction will be delayed due to the inertial forces. Furthermore, in the morphoelastic case we expect some permanent deformations, hence the domain will not return to its original shape. We



run the simulations until  $t = 5$  using a timestep of size  $\Delta t = 0.05$ . The mechanical parameters are set to

$$\mu_\varepsilon = \lambda_\varepsilon = 1, \quad \mu_v = \lambda_v = 0.5, \quad \rho = 0.1. \quad (2.127)$$

Figure 2.11 shows the  $y$ -coordinate of the red point as a function of time. Note that we have considered two different morphoelastic models, one with  $\alpha = 0.5$  and the other with  $\alpha = 1$ . We see that in the latter case, there is more permanent deformation.

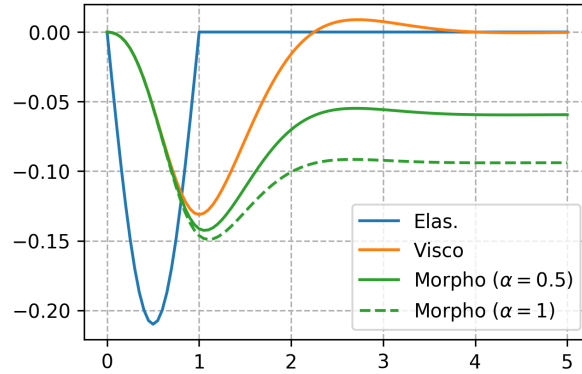


Figure 2.11:  $y$ -coordinate of bottom right corner as a function time for various elastic models.

Also observe that in the purely elastic model, the reaction the force is indeed instantaneous. After the body force vanishes, the domain will immediately return to its initial state. In the viscoelastic model, the domain also returns to the initial state, but this takes more time. Moreover, in both the viscoelastic and morphoelastic models, there is a small amount of ‘overshoot’. Due to the inclusion of inertial forces, the domain slightly oscillates around the equilibrium position.

For the second comparison, we use the same initial domain and boundary condition. The body force is set to

$$\mathbf{g}(\mathbf{x}, t) = \begin{pmatrix} 0.2 \sin(\pi t) \mathbb{1}_{[0,1]}(t) \\ 0 \end{pmatrix}. \quad (2.128)$$

Thus, we now expect the domain to stretch in the  $x$ -direction. The following mechanical parameters are used:

$$\mu_\varepsilon = \lambda_\varepsilon = 1, \quad \mu_v = \lambda_v = 0.5, \quad \rho = 0.5. \quad (2.129)$$

Thus, due to  $\lambda_\varepsilon$  being positive, we also expect some transaxial compression. To measure the stretching and compressing of the domain, we track the maximum domain width (in the  $x$ -direction), and the minimum domain height (in the  $y$ -direction). Figures 2.12a and 2.12b show these two quantities as a function of time.

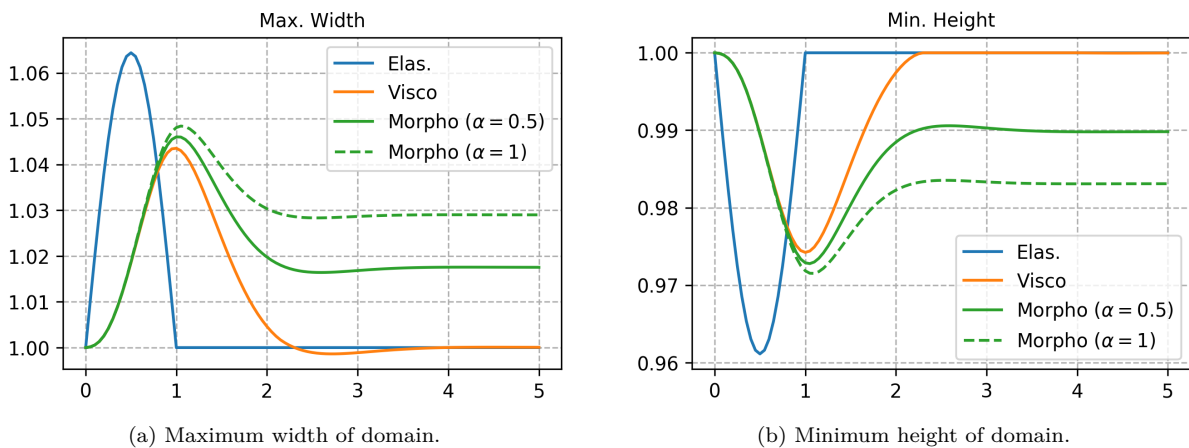


Figure 2.12: Maximum domain width and minimum domain height as a function of time for various elastic models.

Note that in Figure 2.12b the line corresponding to the viscoelastic simulation contains a non-differentiable point. This is simply the result of fixing the left border. Hence there will always be a height of length 1 in the domain.

An interesting property of the morphoelastic model is that although the domain does not converge to its original state, and thus the displacement does not converge to zero, the strain energy *does* converge to zero. In Figure 2.13 the strain energy is shown for the morphoelastic simulation with  $\alpha = 1$ , thus it corresponds to the dashed green lines in Figures 2.12a and 2.12b. We define the elastic strain energy by

$$\begin{aligned} E_{\text{strain}}(t) &= \frac{1}{2} \int_{\Omega(t)} \boldsymbol{\sigma}_{\text{el}}(\boldsymbol{\varepsilon}) : \boldsymbol{\varepsilon} \, d\Omega \\ &= \frac{1}{2} \int_{\Omega(t)} \mu_{\varepsilon} \boldsymbol{\varepsilon} : \boldsymbol{\varepsilon} + \lambda_{\varepsilon} \text{tr}(\boldsymbol{\varepsilon})^2 \, d\Omega. \end{aligned} \quad (2.130)$$

This shows us that the classical definition of the strain tensor (2.8) does not make sense in a morphoelastic setting. The strain tensor is no longer directly related to the displacement.

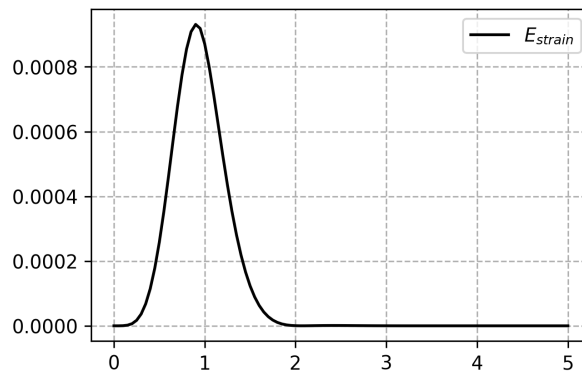


Figure 2.13: Strain energy as a function of time corresponding to the morphoelastic simulation with  $\alpha = 1$ .

# 3

## The Mechanics of Porous Media

In this chapter a number of finite element models are developed for different poroelastic systems. We start with the one-dimensional Terzaghi problem, which is a simplified poroelastic model. It is solved using linear-linear elements, meaning that we use linear basis functions for both  $\mathbf{u}$  and  $p$ , and quadratic-linear elements, where we use quadratic basis functions for  $\mathbf{u}$  and linear for  $p$ . In both cases we encounter non-physical oscillations in the pressure solution for small values of  $k$ , or small timesteps. It is shown that these oscillations can be fixed by using diffusive stabilization.

We then move on to two-dimensional poroelasticity. First we consider infinitesimal deformations corresponding to high values of the Lamé constants. This poroelastic system is solved on a fixed domain, since the deformation tensor is assumed to be close to unity. We subsequently consider visco- poroelasticity on a moving domain. This model essentially combines poroelasticity with the viscoelastic model from Section 2.3. The chapter is concluded by a comparison between the viscoelastic and visco-poroelastic models.

### 3.1 Porous Media

In this section we will discuss the basic physical assumptions on which poroelastic systems are based. A porous medium is a material that contains small, connected voids. These voids, called pores, are completely filled with one or more fluids. Fluid flows through porous media are used in many different fields. For example in the field of geomechanics, the interaction between water and oil in a porous rock is often modeled using a two-phase fluid flow in a porous medium. In the context of this project the porous medium represents a tissue in which tumor cells are present. The poro-fluid is the extra-cellular fluid.

The pores in the medium are very small. We consider the porous medium from a macroscopic point of view, hence we are not interested in the precise locations and shapes of the pores. Instead we consider the parameter  $\Phi$ , which represents the fraction of space occupied by solid material. Hence  $1 - \Phi$  is the fraction of space containing voids and is called the porosity. Instead of considering the velocity of individual fluid particles, an average poro-fluid flux called the Darcy velocity is used. It is denoted by  $\mathbf{v}_D$ . To relate the Darcy velocity to the actual poro-fluid velocity we need to take into account the porosity:

$$\mathbf{v}_F = \frac{\mathbf{v}_D}{1 - \Phi}. \quad (3.1)$$

If  $\Phi$  is close to 1 then there is not much space for the fluid to occupy. Thus to achieve a certain Darcy velocity, the fluid velocity must be high. Darcy's law relates the Darcy flux to the pressure gradient:

$$\mathbf{v}_D = -\frac{k_{\text{perm}}}{\mu_{\text{dyn}}} \nabla p. \quad (3.2)$$

Here  $k_{\text{perm}}$  is the permeability of the medium, and  $\mu_{\text{dyn}}$  is the dynamic viscosity of the poro-fluid. Note that  $\mu_{\text{dyn}}$  might be different from the viscoelastic parameter  $\mu_v$ . The reason for this difference is that the medium might contain other fluids or materials, leading to a different viscosity of the medium compared to the poro-fluid. From this point, we will write  $k = k_{\text{perm}}/\mu_{\text{dyn}}$ . The parameter  $k$  is assumed to be constant in the remainder of this work. It will become important in Chapter 3: small values can lead to

problems with the solution. Darcy's law tells us that the Darcy velocity is proportional to the pressure gradient, and flows from high-pressure locations to low-pressure locations.

To relate  $\mathbf{v}_D$  to the displacement velocity, we define the net velocity  $\mathbf{v}_{\text{net}} = \mathbf{v} + \mathbf{v}_D$ , where  $\mathbf{v}$  represents the displacement velocity. The net velocity is assumed to be incompressible:  $\nabla \cdot \mathbf{v}_{\text{net}} = 0$ . This assumption follows from conservation of matter. By applying Darcy's law we find the equation

$$\nabla \cdot \mathbf{v} - \nabla \cdot (k \nabla p) = 0. \quad (3.3)$$

In the poroelastic systems from Chapter 3, we use the slightly adapted equation

$$\nabla \cdot \mathbf{v} - \nabla \cdot (k \nabla p) = f. \quad (3.4)$$

Here  $f$  represent a mass source (or sink) function. In a poroelastic medium the fluid pressure must also be taken into account in the force balance equation. For this reason the term  $-p\mathbf{I}$  is added to the stress tensor.

### 3.2 Terzaghi Problem

We develop several different finite element models for the Terzaghi problem in one dimension. This problem is analyzed in [1] in order to show that non-physical oscillations occur for small timesteps. We consider the domain  $\Omega = (0, 1)$ , assume the medium is porous and purely elastic, and consider only infinitesimal deformations. The corresponding stress tensor in one dimension is given by:

$$\sigma = (\mu_\varepsilon + \lambda_\varepsilon) \frac{\partial u}{\partial x} - p. \quad (3.5)$$

The force balance equation then becomes

$$-\frac{\partial \sigma}{\partial x} = 0, \quad \text{hence} \quad -(\mu_\varepsilon + \lambda_\varepsilon) \frac{\partial u}{\partial x} + \frac{\partial p}{\partial x} = 0 \quad (3.6)$$

Using the fact that for infinitesimal deformations we have  $v = \partial u / \partial t$ , equation (3.3) becomes

$$\frac{\partial}{\partial t} \left( \frac{\partial u}{\partial x} \right) - k \frac{\partial^2 p}{\partial x^2} = 0. \quad (3.7)$$

The Terzaghi problem is obtained by setting  $\mu_\varepsilon + \lambda_\varepsilon = 1$  and  $k = 1$ . We then get the following system of partial differential equations:

$$-\frac{\partial^2 u}{\partial x^2} + \frac{\partial p}{\partial x} = 0, \quad x \in (0, 1), \quad t > 0, \quad (3.8)$$

$$\frac{\partial}{\partial t} \left( \frac{\partial u}{\partial x} \right) - \frac{\partial^2 p}{\partial x^2} = 0, \quad x \in (0, 1), \quad t > 0. \quad (3.9)$$

The following boundary conditions are imposed:

$$\frac{\partial u}{\partial x}(0, t) = -1, \quad p(0, t) = 0, \quad (3.10)$$

$$u(1, t) = 0, \quad \frac{\partial p}{\partial x}(1, t) = 0, \quad (3.11)$$

as well as the initial condition  $\partial u / \partial x(x, 0) = 0$ . Define the following function spaces:

$$\mathcal{V}_0 = \{\varphi \in H^1(0, 1) : \varphi(1) = 0\}, \quad (3.12)$$

$$\mathcal{W}_0 = \{\psi \in H^1(0, 1) : \psi(0) = 0\}. \quad (3.13)$$

The weak form of the Terzaghi problem is given by:

For every  $t > 0$  find  $u(\cdot, t) \in \mathcal{V}_0$  and  $p(\cdot, t) \in \mathcal{W}_0$  such that

$$\begin{aligned} \int_0^1 \frac{\partial u}{\partial x} \frac{d\varphi}{dx} + \frac{\partial p}{\partial x} \varphi \, dx &= \varphi(0), \quad \forall \varphi \in \mathcal{V}_0, \\ \int_0^1 \frac{\partial}{\partial t} \left( \frac{\partial u}{\partial x} \right) \psi + \frac{\partial p}{\partial x} \frac{d\psi}{dx} \, dx &= 0, \quad \forall \psi \in \mathcal{W}_0. \end{aligned} \quad (3.14)$$

Note that the  $\varphi(0)$  in the right-hand side of the first equation is a result of boundary condition (3.10). Let  $u_h$  and  $p_h$  be the finite element solutions. We write

$$u_h(x, t) = \sum_{j=1}^{n_u} u_j(t) \varphi_j(x), \quad p_h(x, t) = \sum_{j=1}^{n_p} p_j(t) \psi_j(x), \quad (3.15)$$

where the  $\varphi_j$  span a finite dimensional subspace of  $\mathcal{V}_0$  and the  $\psi_j$  span a finite dimensional subspace of  $\mathcal{W}_0$ . The Galerkin equations corresponding to the weak form (3.14) are given by

$$\begin{aligned} \sum_{j=1}^{n_u} u_j \int_0^1 \frac{d\varphi_i}{dx} \frac{d\varphi_j}{dx} dx + \sum_{j=1}^{n_p} p_j \int_0^1 \varphi_i \frac{d\psi_j}{dx} dx &= \varphi_i(0), \quad i \in \{1, \dots, n_u\}, \\ \sum_{j=1}^{n_u} \frac{du_j}{dt} \int_0^1 \psi_i \frac{d\varphi_j}{dx} dx + \sum_{j=1}^{n_p} p_j \int_0^1 \frac{d\psi_i}{dx} \frac{d\psi_j}{dx} dx &= 0, \quad i \in \{1, \dots, n_p\}. \end{aligned} \quad (3.16)$$

In block matrix form we have

$$\begin{pmatrix} 0 & 0 \\ D & 0 \end{pmatrix} \begin{pmatrix} \dot{\mathbf{u}} \\ \dot{\mathbf{p}} \end{pmatrix} + \begin{pmatrix} S & G \\ 0 & L \end{pmatrix} \begin{pmatrix} \mathbf{u} \\ \mathbf{p} \end{pmatrix} = \begin{pmatrix} \mathbf{b} \\ \mathbf{0} \end{pmatrix}. \quad (3.17)$$

Here  $\mathbf{u}$  and  $\mathbf{p}$  are the vectors

$$\mathbf{u} = (u_1(t), u_2(t), \dots, u_{n_u}(t))^{\top}, \quad (3.18)$$

$$\mathbf{p} = (p_1(t), p_2(t), \dots, p_{n_p}(t))^{\top}. \quad (3.19)$$

The blocks are given component-wise by:

$$D_{ij} = \int_0^1 \psi_i \frac{d\varphi_j}{dx} dx, \quad i \in \{1, \dots, n_p\}, j \in \{1, \dots, n_u\}, \quad (3.20)$$

$$G_{ij} = \int_0^1 \varphi_i \frac{d\psi_j}{dx} dx, \quad i \in \{1, \dots, n_u\}, j \in \{1, \dots, n_p\}, \quad (3.21)$$

$$S_{ij} = \int_0^1 \frac{d\varphi_i}{dx} \frac{d\varphi_j}{dx} dx, \quad i, j \in \{1, \dots, n_u\}, \quad (3.22)$$

$$L_{ij} = \int_0^1 \frac{d\psi_i}{dx} \frac{d\psi_j}{dx} dx, \quad i, j \in \{1, \dots, n_p\}. \quad (3.23)$$

The right-hand side vector  $\mathbf{b}$  is given component-wise by:

$$b_i = \varphi_i(0). \quad (3.24)$$

In practice we will use basis functions that equal 1 at exactly one grid point and zero everywhere else. In this case  $\mathbf{b}$  will have a 1 at its first position, and contain zeroes everywhere else. From (3.18) and (3.19) it follows that  $D$  is a  $n_p \times n_u$  matrix,  $G$  is  $n_u \times n_p$ ,  $S$  is  $n_u \times n_u$ , and  $L$  is  $n_p \times n_p$ . We apply implicit Euler to the system in (3.17). This leads to the following system:

$$\begin{pmatrix} S & G \\ D & \Delta t L \end{pmatrix} \begin{pmatrix} \mathbf{u}^m \\ \mathbf{p}^m \end{pmatrix} = \begin{pmatrix} 0 & 0 \\ D & 0 \end{pmatrix} \begin{pmatrix} \mathbf{u}^{m-1} \\ \mathbf{p}^{m-1} \end{pmatrix} + \begin{pmatrix} \mathbf{b} \\ \mathbf{0} \end{pmatrix}. \quad (3.25)$$

The superscript denotes a time level.

### 3.2.1 Linear-Linear Elements

We first use linear basis functions for  $u$  and  $p$ . In this we have  $n_p = n_u$ . Let  $\{0 = x_1 < x_2 < \dots < x_{n+1} = 1\}$  be a uniform partition of  $[0, 1]$  with meshwidth  $h$ . We denote the linear basis function centered in  $x_k$  by  $\lambda_k$ . Let  $e_k$  be the element  $[x_k, x_{k+1}]$  for  $k \in \{1, \dots, n\}$ . On  $e_k$ , only the functions  $\lambda_k$  and  $\lambda_{k+1}$  are nonzero. Furthermore, on  $e_k$  we have

$$\frac{d\lambda_k}{dx} = -\frac{1}{h}, \quad \frac{d\lambda_{k+1}}{dx} = \frac{1}{h}. \quad (3.26)$$

The element matrices corresponding to (3.20)-(3.23) are given by:

$$D^{e_k} = G^{e_k} = \frac{1}{2} \begin{bmatrix} -1 & 1 \\ -1 & 1 \end{bmatrix}, \quad (3.27)$$

$$S^{e_k} = L^{e_k} = \frac{1}{h} \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}. \quad (3.28)$$

In [1] it is shown that in order to obtain an oscillation-free solution, we must have

$$\Delta t > \frac{h^2}{4}. \quad (3.29)$$

Figure 3.1 and 3.2 show the finite element solution after one timestep, using a meshwidth of  $h = 1/10$  and a timestep of  $\Delta t = 1/200$  and  $\Delta t = 1/800$  respectively. Note that in Figure 3.1 the bound in (3.29) holds, whereas in Figure 3.2 it does not.

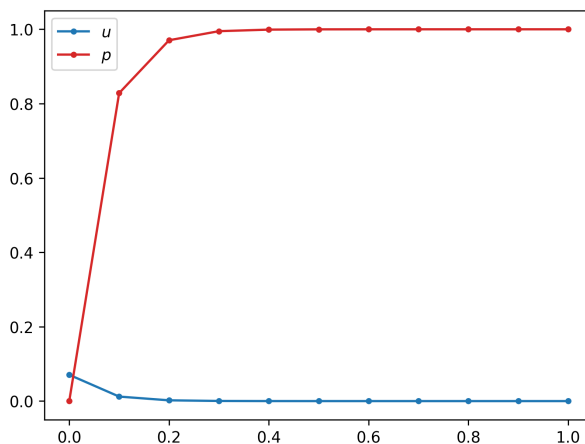


Figure 3.1: FE-solution after one timestep using  $h = 1/10$  and  $\Delta t = 1/200$ .

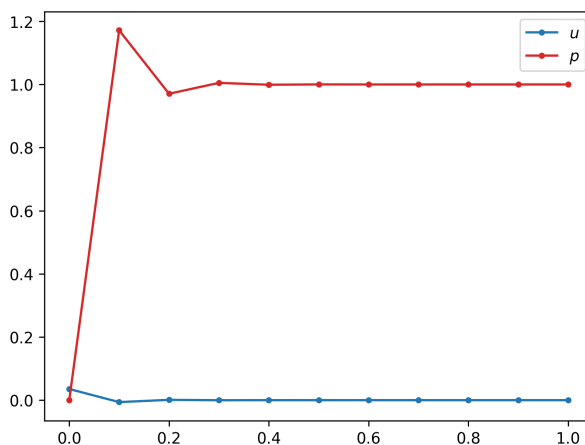


Figure 3.2: FE-solution after one timestep using  $h = 1/10$  and  $\Delta t = 1/800$ .

### 3.2.2 Quadratic-Linear Elements

We now use quadratic basis functions for  $u$  and linear basis functions for  $p$ . Such elements are called Taylor-Hood elements. Since quadratic basis functions need an additional control vertex, we have  $n_u = 2n_p - 1$ . We again use a uniform grid with meshwidth  $h$ . Let  $x_{k+\frac{1}{2}}$  be the midpoint of the element  $e_k = [x_k, x_{k+1}]$ . The quadratic basis functions  $\varphi_k, \varphi_{k+\frac{1}{2}}, \varphi_{k+1}$  are defined by the fact that they equal 1 in their respective grid point and zero on the other two points. They can be expressed in terms of the

linear basis functions:

$$\varphi_k = \lambda_k(2\lambda_k - 1), \quad (3.30)$$

$$\varphi_{k+\frac{1}{2}} = 4\lambda_k\lambda_{k+1}, \quad (3.31)$$

$$\varphi_{k+1} = \lambda_{k+1}(2\lambda_{k+1} - 1). \quad (3.32)$$

Figure 3.3 shows the quadratic basis functions on the element  $e_k$ .

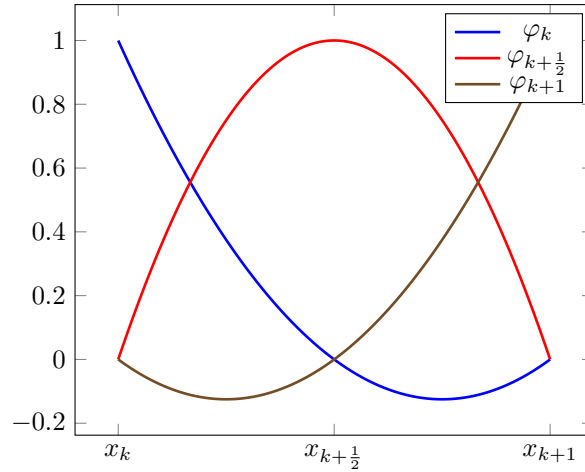


Figure 3.3: Quadratic basis functions on the element  $e_k$ .

From these expressions one can easily compute the derivatives of the quadratic basis functions:

$$\frac{d\varphi_k}{dx} = -\frac{1}{h}(4\lambda_k - 1), \quad (3.33)$$

$$\frac{d\varphi_{k+\frac{1}{2}}}{dx} = \frac{4}{h}(\lambda_k - \lambda_{k+1}), \quad (3.34)$$

$$\frac{d\varphi_{k+1}}{dx} = \frac{1}{h}(4\lambda_{k+1} - 1). \quad (3.35)$$

Now Holand & Bell's theorem [17] can be used to compute the element matrices. Note that most element matrices are no longer  $2 \times 2$ . For example,  $D^{e_k}$  is now of the form

$$D^{e_k} = \begin{bmatrix} D_{k,k}^{e_k} & D_{k,k+\frac{1}{2}}^{e_k} & D_{k,k+1}^{e_k} \\ D_{k+1,k}^{e_k} & D_{k+1,k+\frac{1}{2}}^{e_k} & D_{k+1,k+1}^{e_k} \end{bmatrix}, \quad (3.36)$$

since  $D$  is multiplied with  $\mathbf{u}$ , which now also contains coefficients corresponding to the midpoints of each element. The element matrices corresponding to (3.20)-(3.23) in the case of quadratic basis functions for  $u$  and linear for  $p$  are given by:

$$D^{e_k} = \frac{1}{6} \begin{bmatrix} -5 & 4 & 1 \\ -1 & -4 & 5 \end{bmatrix}, \quad (3.37)$$

$$G^{e_k} = \frac{1}{6} \begin{bmatrix} -1 & 1 \\ -4 & 4 \\ -1 & 1 \end{bmatrix}, \quad (3.38)$$

$$S^{e_k} = \frac{1}{3h} \begin{bmatrix} 7 & -8 & 1 \\ -8 & 16 & -8 \\ 1 & -8 & 7 \end{bmatrix}, \quad (3.39)$$

$$L^{e_k} = \frac{1}{h} \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}. \quad (3.40)$$

In [1] it is shown that the lower bound for  $\Delta t$  in the case of quadratic-linear elements is given by

$$\Delta t > \frac{h^2}{6}. \quad (3.41)$$

While this bound allows for smaller timesteps compared to the linear-linear case, it is only reduced by a factor  $2/3$ . Figures 3.4 and 3.5 again show the finite element solutions after one timestep, using the same parameters as in Figures 3.1 and 3.2.

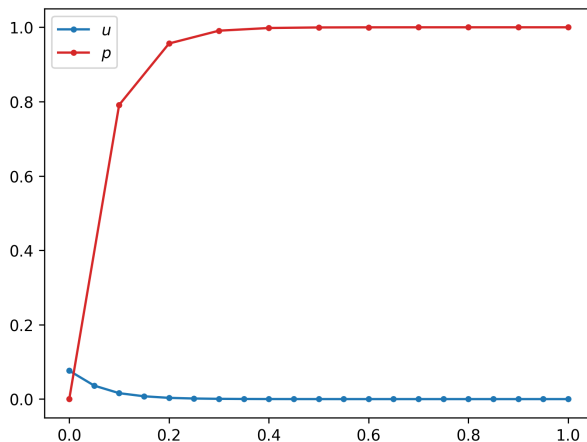


Figure 3.4: FE-solution after one timestep using  $h = 1/10$  and  $\Delta t = 1/200$ .

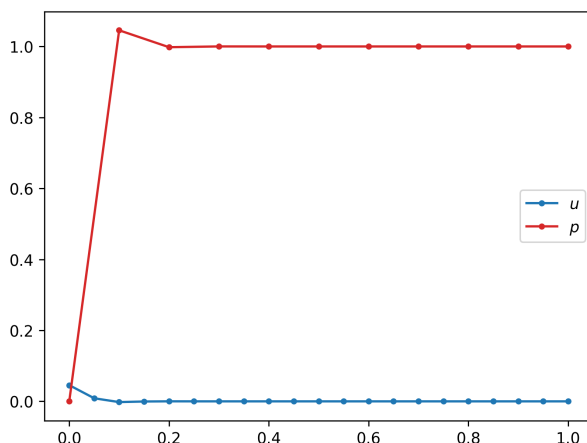


Figure 3.5: FE-solution after one timestep using  $h = 1/10$  and  $\Delta t = 1/800$ .

Note that oscillations still occur for  $\Delta t = 1/800$ , but they are smaller in magnitude compared to the linear-linear case.

### 3.2.3 Stabilization

In [13] it is shown that the non-physical oscillations can be removed by introducing a stabilization term to equation (3.9). In the case of a uniform grid, equation (3.9) is replaced by:

$$\frac{\partial}{\partial t} \left( \frac{\partial u}{\partial x} \right) - \frac{\partial^2 p}{\partial x^2} - \beta h^2 \frac{\partial}{\partial t} \left( \frac{\partial^2 p}{\partial x^2} \right) = 0, \quad (3.42)$$

where  $h$  is the meshwidth and  $\beta$  is a constant that depends on the spatial discretization. For linear-linear elements we have  $\beta = 1/4$  and for Taylor-Hood elements we have  $\beta = 1/6$ . The corresponding Galerkin equations are then given by:

$$\begin{pmatrix} 0 & 0 \\ D & L_s \end{pmatrix} \begin{pmatrix} \dot{\mathbf{u}} \\ \dot{\mathbf{p}} \end{pmatrix} + \begin{pmatrix} S & G \\ 0 & L \end{pmatrix} \begin{pmatrix} \mathbf{u} \\ \mathbf{p} \end{pmatrix} = \begin{pmatrix} \mathbf{b} \\ \mathbf{0} \end{pmatrix}. \quad (3.43)$$

Notice that the difference with (3.17) is the  $L_s$  block in the first matrix. It is given by:

$$L_s = \beta h^2 L. \quad (3.44)$$



We note once again that these expressions are only valid in the case of a uniform grid. Applying backwards Euler to (3.43) yields

$$\begin{pmatrix} S & G \\ D & L_s + \Delta t L \end{pmatrix} \begin{pmatrix} \mathbf{u}^m \\ \mathbf{p}^m \end{pmatrix} = \begin{pmatrix} 0 & 0 \\ D & L_s \end{pmatrix} \begin{pmatrix} \mathbf{u}^{m-1} \\ \mathbf{p}^{m-1} \end{pmatrix} + \begin{pmatrix} \mathbf{b} \\ \mathbf{0} \end{pmatrix}. \quad (3.45)$$

Figures 3.6 and 3.7 show the finite element solutions for linear-linear and quadratic-linear elements after one very small timestep. In both cases, the stabilized solutions show no oscillations.

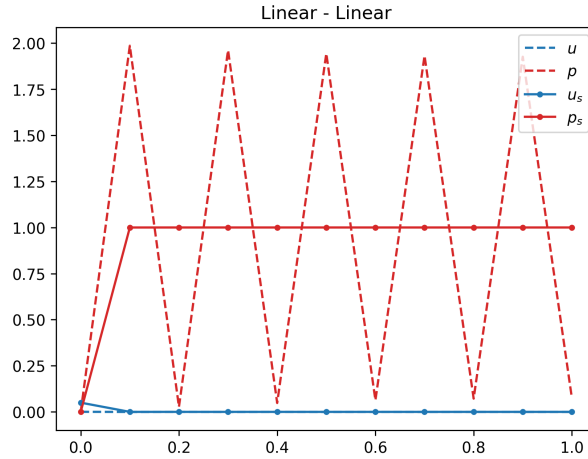


Figure 3.6: Stabilized and non-stabilized FE-solutions using linear-linear elements, after one timestep using  $h = 1/10$  and  $\Delta t = 10^{-6}$ . The solid lines represent the stabilized solutions and the dashed lines represent the non-stabilized solutions

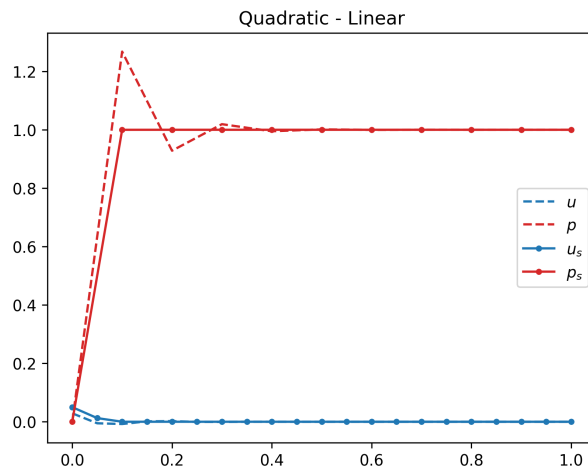


Figure 3.7: Stabilized and non-stabilized FE-solutions using quadratic-linear elements, after one timestep using  $h = 1/10$  and  $\Delta t = 10^{-6}$ . The solid lines represent the stabilized solutions and the dashed lines represent the non-stabilized solutions

### 3.3 Poroelasticity on a Fixed Domain

Now consider a two-dimensional poroelastic system in the framework of infinitesimal deformations: it is assumed that  $\|\nabla \mathbf{u}\| \ll 1$ . Therefore the displacement velocity simplifies to

$$\mathbf{v} \approx \frac{\partial \mathbf{u}}{\partial t}, \quad (3.46)$$

by the derivation in (2.46). For this approximation to be accurate, high values of the Lamé constants  $\mu_\varepsilon$  and  $\lambda_\varepsilon$  are used. Combining the purely elastic model from Section (2.2) with the simplified version of

equation (3.4) yields:

$$-\nabla \cdot \boldsymbol{\sigma}(\mathbf{u}, p) = \mathbf{g}, \quad \text{on } \Omega, \quad (3.47)$$

$$\nabla \cdot \frac{\partial \mathbf{u}}{\partial t} - \nabla \cdot (k \nabla p) = f, \quad \text{on } \Omega. \quad (3.48)$$

The stress tensor is given by:

$$\boldsymbol{\sigma}(\mathbf{u}, p) = \boldsymbol{\sigma}_{\text{el}}(\mathbf{u}) - p \mathbf{I}, \quad (3.49)$$

where the elasticity tensor  $\boldsymbol{\sigma}_{\text{el}}$  is given in (2.10). We impose the following boundary conditions:

$$\mathbf{u} = \mathbf{0}, \quad k \nabla p \cdot \mathbf{n} = 0, \quad \text{on } \Gamma_1, \quad (3.50)$$

$$\boldsymbol{\sigma} \cdot \mathbf{n} = \boldsymbol{\tau}, \quad p = 0, \quad \text{on } \Gamma_2, \quad (3.51)$$

We use Implicit Euler to discretize in time:

$$-\nabla \cdot \boldsymbol{\sigma}(\mathbf{u}^m, p^m) = \mathbf{g}^m, \quad (3.52)$$

$$\nabla \cdot \mathbf{u}^m - \Delta t \nabla \cdot (k \nabla p^m) = \Delta t f^m + \nabla \cdot \mathbf{u}^{m-1}. \quad (3.53)$$

Here  $\Delta t$  denotes the timestep. Using the same notation as in [13] the weak form is given by:

For each  $m \geq 1$  find  $\mathbf{u}^m \in [H_{\Gamma_1}^1(\Omega)]^2$  and  $p^m \in H_{\Gamma_2}^1(\Omega)$  such that

$$a(\mathbf{u}^m, \mathbf{v}) - (p^m, \nabla \cdot \mathbf{v}) = (\mathbf{g}^m, \mathbf{v}) + (\boldsymbol{\tau}^m, \mathbf{v})_{\Gamma_2}, \quad \forall \mathbf{v} \in [H_{\Gamma_1}^1(\Omega)]^2, \quad (3.54)$$

$$(\nabla \cdot \mathbf{u}^m, q) + \Delta t b(p^m, q) = (\Delta t f^m + \nabla \cdot \mathbf{u}^{m-1}, q), \quad \forall q \in H_{\Gamma_2}^1(\Omega). \quad (3.55)$$

The function spaces  $H_{\Gamma_1}^1(\Omega)$  and  $H_{\Gamma_2}^1(\Omega)$  are the Sobolev spaces of degree 1 with functions that vanish on  $\Gamma_1$  and  $\Gamma_2$  respectively. Moreover, the bilinear forms  $a(\cdot, \cdot)$  and  $b(\cdot, \cdot)$  are given by:

$$a(\mathbf{u}, \mathbf{v}) = \mu_\varepsilon \int_{\Omega} \boldsymbol{\varepsilon}(\mathbf{u}) : \boldsymbol{\varepsilon}(\mathbf{v}) \, d\Omega + \lambda_\varepsilon \int_{\Omega} (\nabla \cdot \mathbf{u})(\nabla \cdot \mathbf{v}) \, d\Omega, \quad (3.56)$$

$$b(p, q) = \int_{\Omega} k \nabla p \cdot \nabla q \, d\Omega. \quad (3.57)$$

Finally,  $(\cdot, \cdot)$  denotes the standard inner product on its respective Sobolev space, and

$$(\boldsymbol{\tau}, \mathbf{v})_{\Gamma_2} = \int_{\Gamma_2} \boldsymbol{\tau} \cdot \mathbf{v} \, d\Gamma. \quad (3.58)$$

To obtain a finite element formulation, the infinite dimensional Sobolev spaces  $[H_{\Gamma_1}^1(\Omega)]^2$  and  $H_{\Gamma_2}^1(\Omega)$  are replaced by the finite dimensional ones  $V_h$  and  $Q_h$  respectively. We get the following formulation:

For each  $m \geq 1$  find  $\mathbf{u}_h^m \in V_h$  and  $p_h^m \in Q_h$  such that

$$a(\mathbf{u}_h^m, \mathbf{v}) - (p_h^m, \nabla \cdot \mathbf{v}) = (\mathbf{g}^m, \mathbf{v}) + (\boldsymbol{\tau}^m, \mathbf{v})_{\Gamma_2}, \quad \forall \mathbf{v} \in V_h, \quad (3.59)$$

$$(\nabla \cdot \mathbf{u}_h^m, q) + \Delta t b(p_h^m, q) = (\Delta t f^m + \nabla \cdot \mathbf{u}_h^{m-1}, q), \quad \forall q \in Q_h. \quad (3.60)$$

In general, the domain  $\Omega$  will be replaced by an approximate domain  $\Omega_h$ . All integrations in (3.59) and (3.60) will be done over  $\Omega_h$ . We derive the Galerkin equations for general basis functions  $\varphi_i$  and  $\psi_i$  for  $\mathbf{u}$  and  $p$  respectively. Let  $n_u, n_p$  be the number of grid nodes for the displacement and pressure respectively. We define the finite element solutions by

$$(u^1)_h^m(\mathbf{x}) = \sum_{j=1}^{n_u} (u_j^1)^m \varphi_j(\mathbf{x}), \quad (u^2)_h^m(\mathbf{x}) = \sum_{j=1}^{n_u} (u_j^2)^m \varphi_j(\mathbf{x}), \quad p_h^m(\mathbf{x}) = \sum_{j=1}^{n_p} p_j^m \psi_j(\mathbf{x}). \quad (3.61)$$

The basis functions  $\varphi_j$  and  $\psi_j$  form a basis of  $V_h$  and  $Q_h$  respectively, such that they equal 1 in their corresponding grid node  $\mathbf{x}_j$ . To obtain the Galerkin equations, substitute the finite element solutions

into (3.54) and (3.55). Furthermore set  $\mathbf{v} = (\varphi_i, 0)^\top$  and  $\mathbf{v} = (0, \varphi_i)^\top$  for  $i = 1, \dots, n_u$  in (3.54) and set  $q = \psi_i$  for  $i = 1, \dots, n_p$  in (3.55). We then obtain a linear system of the form

$$\begin{pmatrix} S_{\text{el}} & -D^\top \\ D & k\Delta t L \end{pmatrix} \begin{pmatrix} \mathbf{u}^m \\ \mathbf{p}^m \end{pmatrix} = \begin{pmatrix} \mathbf{g}^m + \boldsymbol{\tau}^m \\ \Delta t \mathbf{f}^m \end{pmatrix} + \begin{pmatrix} 0 \\ D\mathbf{u}^{m-1} \end{pmatrix}. \quad (3.62)$$

Where  $\mathbf{u}^m = ((\mathbf{u}^1)^m, (\mathbf{u}^2)^m)^\top$  and

$$\mathbf{u}_1^m = ((u_1^1)^m, (u_2^1)^m, \dots, (u_{n_u}^1)^m)^\top, \quad (3.63)$$

$$\mathbf{u}_2^m = ((u_1^2)^m, (u_2^2)^m, \dots, (u_{n_u}^2)^m)^\top, \quad (3.64)$$

$$\mathbf{p}^m = (p_1^m, p_2^m, \dots, p_{n_p}^m)^\top. \quad (3.65)$$

Moreover, the matrices  $S$  and  $D$  are the block matrices

$$S_{\text{el}} = \begin{pmatrix} S_{\text{el}}^{11} & S_{\text{el}}^{12} \\ S_{\text{el}}^{21} & S_{\text{el}}^{22} \end{pmatrix}, \quad D = (D^1 \quad D^2), \quad (3.66)$$

where the blocks  $S_{\text{el}}^{11}, S_{\text{el}}^{12}, S_{\text{el}}^{21}$  and  $S_{\text{el}}^{22}$  are given element-wise by (2.30)-(2.33). We also write  $\mathbf{g} = (\mathbf{g}^1, \mathbf{g}^2)^\top$  and  $\boldsymbol{\tau} = (\boldsymbol{\tau}^1, \boldsymbol{\tau}^2)^\top$ . The remaining blocks and vectors are given element-wise by:

$$(D^1)_{ij} = \int_{\Omega} \psi_i \frac{\partial \varphi_j}{\partial x} \, d\Omega, \quad (D^2)_{ij} = \int_{\Omega} \psi_i \frac{\partial \varphi_j}{\partial y} \, d\Omega, \quad i \in \{1, \dots, n_p\}, \quad j \in \{1, \dots, n_u\}. \quad (3.67)$$

$$(\mathbf{g}^1)_i = \int_{\Omega} g^1 \varphi_i \, d\Omega, \quad (\mathbf{g}^2)_i = \int_{\Omega} g^2 \varphi_i \, d\Omega, \quad i \in \{1, \dots, n_u\}. \quad (3.68)$$

$$(\boldsymbol{\tau}^1)_i = \int_{\Gamma_2} \tau^1 \varphi_i \, d\Gamma, \quad (\boldsymbol{\tau}^2)_i = \int_{\Gamma_2} \tau^2 \varphi_i \, d\Gamma, \quad i \in \{1, \dots, n_u\}. \quad (3.69)$$

$$(\mathbf{f})_i = \int_{\Omega} f \psi_i \, d\Omega, \quad i \in \{1, \dots, n_p\}. \quad (3.70)$$

The matrix  $L$  is a Laplace matrix:

$$L_{ij} = \int_{\Omega} \nabla \psi_i \cdot \nabla \psi_j \, d\Omega, \quad i, j \in \{1, \dots, n_p\}. \quad (3.71)$$

### 3.3.1 Linear-Linear Elements

We investigate the monotonicity and stability of the finite element solution for linear-linear elements. To this end we introduce some helpful notations:

**Definition 1.** A triangulation  $\mathcal{T}_h$  of  $\Omega$  is a set of triangles such that:

$$\bigcup_{T \in \mathcal{T}_h} T =: \Omega_h \approx \Omega.$$

The parameter  $h$  is called the meshwidth, which is often defined as

$$h = \max_{T \in \mathcal{T}_h} \{\text{diam}(T)\}.$$

Given a triangulation  $\mathcal{T}_h$ , we define the following polynomial function spaces on  $\Omega_h$ :

**Definition 2.** The space of locally polynomial functions of degree  $k$  with respect to the triangulation  $\mathcal{T}_h$  is defined by:

$$\mathbb{P}^k(\mathcal{T}_h) = \{\varphi \in H^1(\Omega_h) : \forall T \in \mathcal{T}_h : \varphi|_T \text{ is a polynomial of degree } k\}. \quad (3.72)$$

Set  $V_h = [\mathbb{P}_{\Gamma_1}^1(\mathcal{T}_h)]^2$  and  $Q_h = \mathbb{P}_{\Gamma_2}^1(\mathcal{T}_h)$ , where the subscripts denote part of the boundary on which the functions vanish. Note that in this case we have  $n_u = n_p$ . In [1], the poroelastic problem (3.47)-(3.48) with boundary conditions (3.50)-(3.51) is solved numerically on  $\Omega = (-4, 4)^2$ , using the right-hand side

functions  $\mathbf{g} = (0, -1)^\top$  and  $f = 0$ . The stress vector  $\boldsymbol{\tau}$  on  $\Gamma_2$  is set to  $\mathbf{0}$ . Furthermore, the boundary components are given by:

$$\Gamma_2 = \{(x, y) \in \Gamma : x = 4 \text{ or } y = 4\}, \quad \Gamma_1 = \Gamma \setminus \Gamma_2.$$

We set  $\mu_\varepsilon = \lambda_\varepsilon = 1000$  and vary the parameter  $k$ . In figures 3.8 and 3.9 the finite element solutions are shown after one timestep using a meshwidth of  $h = 0.2$ , for  $k = 10^{-2}$  and  $k = 10^{-8}$  respectively.

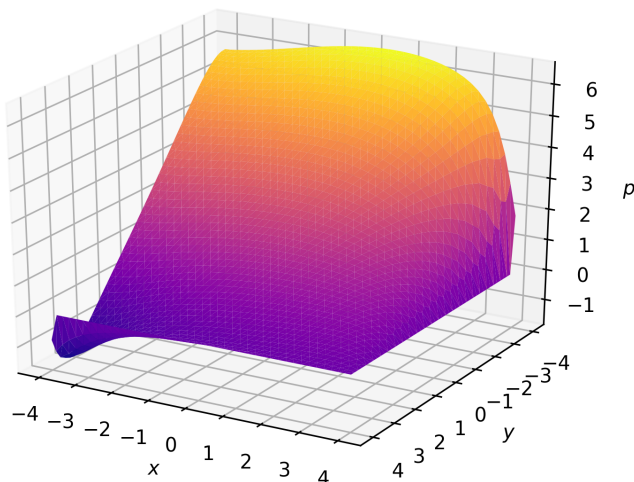


Figure 3.8: 3D-profile of the pressure at  $t = 0.01$  for  $k = 10^{-2}$ .

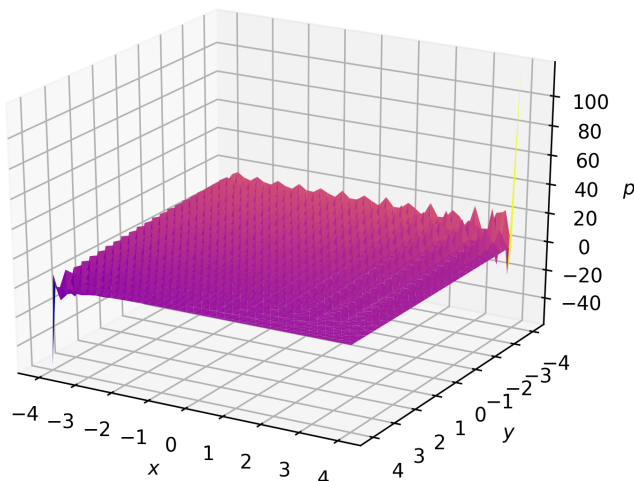
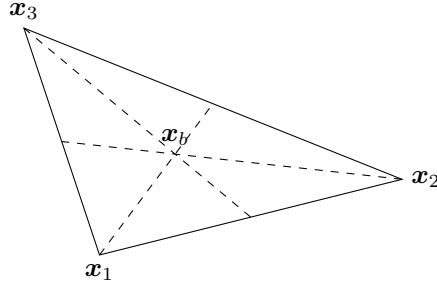


Figure 3.9: 3D-profile of the pressure at  $t = 0.01$  for  $k = 10^{-8}$ .

We observe non-physical oscillations for  $k = 10^{-8}$ . These oscillations also occur for higher values of  $k$  if the timestep is reduced. These results are similar to the 1D case, where non-physical oscillations occur due the left-hand side matrix is no longer being an M-matrix for small timesteps. This phenomenon is a corollary of the fact that no inf-sup condition is satisfied when using linear-linear elements. In [13] an inf-sup condition is proved for the MINI element. This means that for this particular finite element discretization, no non-physical oscillations will occur: the solution is stabilized. In the next section, we will derive this discretization.

### 3.3.2 Stabilization with MINI Element

A MINI element is a triangular element with four nodes: its vertices and its barycenter. A generic MINI element is shown in Figure 3.10.

Figure 3.10: Nodes of a MINI element, where  $\mathbf{x}_b$  is the barycenter of  $T$ 

On this element, there are four basis functions; the linear barycentric coordinate functions that correspond to a vertex, and a cubic ‘bubble function’ that corresponds to the barycenter.

**Definition 3.** Let  $T \in \mathcal{T}_h$  be a triangular element. Let  $\lambda_1, \lambda_2, \lambda_3$  denote the barycentric coordinates of  $T$ , viewed as functions of position  $\mathbf{x} \in T$ . Then the *bubble function*  $\varphi_{b,T}$  is defined as:

$$\varphi_{b,T} = 27\lambda_1\lambda_2\lambda_3. \quad (3.73)$$

The space of bubble functions w.r.t the triangulation  $\mathcal{T}_h$  is given by

$$\mathbb{B}(\mathcal{T}_h) = \text{Span} \{ \varphi_{b,T} : T \in \mathcal{T}_h \}. \quad (3.74)$$

The function  $\varphi_{b,T}$  is called a bubble function because it vanishes on the edges of  $T$  and it is equal to 1 at the barycenter of  $T$ . We use a superscripts to denote function spaces of vector-valued functions. For example, when we write  $\mathbf{u} \in [\mathbb{P}^k(\mathcal{T}_h)]^2$ , then in fact we mean  $\mathbf{u} = (u^1, u^2)$  with  $u^1, u^2 \in \mathbb{P}^k(\mathcal{T}_h)$ .

**Theorem 4.** Let  $a(\cdot, \cdot)$  be bilinear form as defined in (3.56). Furthermore, let  $\mathbf{u}_\ell, \mathbf{v}_\ell \in [\mathbb{P}^1(\mathcal{T}_h)]^2$  and  $\mathbf{u}_b, \mathbf{v}_b \in [\mathbb{B}(\mathcal{T}_h)]^2$ , then

$$a(\mathbf{u}_\ell + \mathbf{u}_b, \mathbf{v}_\ell + \mathbf{v}_b) = a(\mathbf{u}_\ell, \mathbf{v}_\ell) + a(\mathbf{u}_b, \mathbf{v}_b).$$

*Proof.* Note that it suffices to show that  $a(\mathbf{u}_\ell, \mathbf{v}_b) = 0$  and  $a(\mathbf{u}_b, \mathbf{v}_\ell) = 0$ . Since  $a(\cdot, \cdot)$  is symmetric, only one of these equalities has to be proved. From its definition, it is clear that  $a(\mathbf{u}_\ell, \mathbf{v}_b)$  is a weighted sum of integrals of the form

$$\int_{\Omega_h} \frac{\partial u_\ell^i}{\partial x_\alpha} \frac{\partial v_b^j}{\partial x_\beta} d\Omega,$$

where  $i, j, \alpha, \beta \in \{1, 2\}$ . We split this integral into integrals over every triangle  $T \in \mathcal{T}_h$ :

$$\int_{\Omega_h} \frac{\partial u_\ell^i}{\partial x_\alpha} \frac{\partial v_b^j}{\partial x_\beta} d\Omega = \sum_{T \in \mathcal{T}_h} \int_T \frac{\partial u_\ell^i}{\partial x_\alpha} \frac{\partial v_b^j}{\partial x_\beta} d\Omega.$$

We will use the integration by parts formula combined with Gauss’ divergence theorem:

$$\int_T \xi(\nabla \cdot \mathbf{F}) d\Omega = \int_{\partial T} \xi(\mathbf{n} \cdot \mathbf{F}) d\Gamma - \int_T \nabla \xi \cdot \mathbf{F} d\Omega.$$

For  $\beta = 1$  let  $\mathbf{F} = (v_b^j, 0)$ , and for  $\beta = 2$  let  $\mathbf{F} = (0, v_b^j)$ . Moreover, let  $\xi = \partial u_\ell^i / \partial x_\alpha$ . Then by the above formula we get

$$\int_T \frac{\partial u_\ell^i}{\partial x_\alpha} \frac{\partial v_b^j}{\partial x_\beta} d\Omega = \int_{\partial T} \frac{\partial v_b^j}{\partial x_\alpha} n_\beta v_b^j d\Gamma - \int_T \frac{\partial^2 u_\ell^i}{\partial x_\beta \partial x_\alpha} v_b^j d\Omega.$$

Since  $v_b^j$  is a bubble function it vanishes on the boundary  $\partial T$ . Thus the first integral equals zero. Furthermore, since  $u_\ell^i$  is linear on  $T$  all its second derivatives vanish. Therefore the second integral also equals zero. We conclude that on each triangle  $T$  the above integral vanishes, and thus  $a(\mathbf{u}_\ell, \mathbf{v}_b) = 0$ .  $\square$

To stabilize the solution we set  $V_h = [\mathbb{P}^1(\mathcal{T}_h) \oplus \mathbb{B}(\mathcal{T}_h)]^2$  and  $Q_h = \mathbb{P}^1(\mathcal{T}_h)$ . By Theorem 4 we obtain a system of the form:

$$\begin{pmatrix} S_b & 0 & -D_b^\top \\ 0 & S_{\text{el}} & -D^\top \\ D_b & D & k\Delta t L \end{pmatrix} \begin{pmatrix} \mathbf{u}_b^m \\ \mathbf{u}^m \\ \mathbf{p}^m \end{pmatrix} = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ D_b & D & 0 \end{pmatrix} \begin{pmatrix} \mathbf{u}_b^{m-1} \\ \mathbf{u}^{m-1} \\ \mathbf{p}^{m-1} \end{pmatrix} + \begin{pmatrix} \mathbf{g}_b^m \\ \mathbf{g}^m + \boldsymbol{\tau}^m \\ \Delta t \mathbf{f}^m \end{pmatrix} \quad (3.75)$$

Note that there is no vector  $\boldsymbol{\tau}_b$  because it would involve boundary integrals of bubble functions, which vanish. Also note that there are just as many ‘bubble-unknowns’ as there are triangles in the triangulation. We denote this number by  $n_T$ . Let  $n$  denote the number of  $\mathbb{P}^1$ -unknowns, which corresponds to the number of nodes in the triangulation. Then  $S_b$  is a  $2n_T \times 2n_T$  matrix that can be subdivided into four blocks like in (3.66). Moreover  $D_b$  is a  $n \times 2n_T$  matrix. The per-element contributions of the newly introduced bubble matrices and vectors can be found in Section A.3.

### Elimination of Bubbles

Instead of solving (3.75), one can eliminate the bubble equations to obtain a smaller system.

**Theorem 5.** *If  $(\mathbf{u}_b^m, \mathbf{u}^m, \mathbf{p}^m)^\top$  solves (3.75), then  $(\mathbf{u}^m, \mathbf{p}^m)^\top$  solves:*

$$\begin{pmatrix} S_{el} & -D^\top \\ D & k\Delta tL + C_b \end{pmatrix} \begin{pmatrix} \mathbf{u}^m \\ \mathbf{p}^m \end{pmatrix} = \begin{pmatrix} 0 & 0 \\ D & C_b \end{pmatrix} \begin{pmatrix} \mathbf{u}^{m-1} \\ \mathbf{p}^{m-1} \end{pmatrix} + \begin{pmatrix} \mathbf{g}^m + \boldsymbol{\tau}^m \\ \Delta t\mathbf{f}^m - B_b(\mathbf{g}_b^m - \mathbf{g}_b^{m-1}) \end{pmatrix}, \quad (3.76)$$

where

$$C_b = D_b S_b^{-1} D_b^\top, \quad B_b = D_b S_b^{-1}. \quad (3.77)$$

*Proof.* Note that we can rewrite the first block-equation of (3.75) to:

$$S_b \mathbf{u}_b^m = D_b^\top \mathbf{p}^m + \mathbf{g}_b^m.$$

And thus

$$\begin{aligned} D_b \mathbf{u}_b^m &= D_b S_b^{-1} S_b \mathbf{u}_b^m = D_b S_b^{-1} D_b^\top \mathbf{p}^m + D_b S_b^{-1} \mathbf{g}_b^m \\ &= C_b \mathbf{p}^m + B_b \mathbf{g}_b^m. \end{aligned}$$

We will later give an explicit expression for  $S_b^{-1}$ . Substituting the above expression into the third block-equation of (3.75) yields:

$$D\mathbf{u}^m + (k\Delta tL + C_b)\mathbf{p}^m = D\mathbf{u}^{m-1} + C_b\mathbf{p}^{m-1} + \Delta t\mathbf{f}^m - B_b(\mathbf{g}_b^m - \mathbf{g}_b^{m-1}).$$

Thus, we have eliminated the first block-equation and obtain the system in (3.76).  $\square$

We give an explicit expression for  $S_b^{-1}$ . We have already seen that  $S_b$  consists of four diagonal blocks. From (A.37) and (A.38) it is evident that  $S_b^{12} = S_b^{21}$ . In other words  $S_b$  is symmetric. We have the following Lemma:

**Lemma 4.** *Let  $S_b$  be the matrix as defined by (A.36)-(A.39), and let*

$$\Lambda = S_b^{11} S_b^{22} - (S_b^{12})^2. \quad (3.78)$$

*If  $\Lambda$  is invertible, then*

$$S_b^{-1} = \begin{pmatrix} \Lambda^{-1} & 0 \\ 0 & \Lambda^{-1} \end{pmatrix} \begin{pmatrix} S_b^{22} & -S_b^{12} \\ -S_b^{12} & S_b^{11} \end{pmatrix} \quad (3.79)$$

*Proof.* Note that since all blocks are diagonal and have the same dimensions, they all commute.

$$\begin{aligned} S^{-1}S &= \begin{pmatrix} \Lambda^{-1} & 0 \\ 0 & \Lambda^{-1} \end{pmatrix} \begin{pmatrix} S_b^{22} & -S_b^{12} \\ -S_b^{12} & S_b^{11} \end{pmatrix} \begin{pmatrix} S_b^{11} & S_b^{12} \\ S_b^{12} & S_b^{22} \end{pmatrix} \\ &= \begin{pmatrix} \Lambda^{-1} & 0 \\ 0 & \Lambda^{-1} \end{pmatrix} \begin{pmatrix} S_b^{11} S_b^{22} - (S_b^{12})^2 & 0 \\ 0 & S_b^{11} S_b^{22} - (S_b^{12})^2 \end{pmatrix} \\ &= \begin{pmatrix} I & 0 \\ 0 & I \end{pmatrix} \end{aligned}$$

$\square$

Note that we need  $\Lambda$  to be invertible. In practice we find that this is always the case, but we have yet to prove this fact.

### 3.3.3 Diffusive Stabilization

A different method used to stabilize the system is to introduce a small diffusive perturbation to the flow equation, similar to the one-dimensional case. In [1] it is shown that on a uniform grid the perturbed flow equation is given by:

$$\nabla \cdot \frac{\partial \mathbf{u}}{\partial t} - \nabla \cdot (k \nabla p) = f + \beta h^2 \frac{\partial}{\partial t} (\nabla^2 p) \quad (3.80)$$

Note that  $\beta$  is a tuning parameter. In [1] it is shown that in the one-dimensional case the optimal value of  $\beta$  is given by

$$\beta = \frac{1}{4(\lambda_\varepsilon + \mu_\varepsilon)}. \quad (3.81)$$

Numerical experiments show that this value is also optimal in 2D on a uniform grid [1, Figure 5]. Note that an error of order  $h^2$  is introduced as a result of the added perturbation in (3.80). Since the finite element error is already of order  $h^2$  when using linear-linear elements, the effect of the perturbation will be minimal. Discretization on a uniform grid using linear-linear elements and applying implicit Euler yields a system of the form

$$\begin{pmatrix} S & -D^\top \\ D & (k\Delta t + \beta h^2)L \end{pmatrix} \begin{pmatrix} \mathbf{u}^m \\ \mathbf{p}^m \end{pmatrix} = \begin{pmatrix} 0 & 0 \\ D & \beta h^2 L \end{pmatrix} \begin{pmatrix} \mathbf{u}^{m-1} \\ \mathbf{p}^{m-1} \end{pmatrix} + \begin{pmatrix} \mathbf{g}^m + \boldsymbol{\tau}^m \\ \Delta t \mathbf{f}^m \end{pmatrix}. \quad (3.82)$$

Here  $L$  is a Laplace matrix:

$$L_{ij} = \int_{\Omega_h} \nabla \psi_i \cdot \nabla \psi_j \, d\Omega. \quad (3.83)$$

The stabilization term  $\beta h^2 L$  can be generalized to non-uniform grids by replacing it with the matrix  $C_s$ , given by:

$$(C_s)_{ij} = \beta \sum_{T \in \mathcal{T}} h_T^2 \int_T \nabla \psi_i \cdot \nabla \psi_j \, d\Omega. \quad (3.84)$$

Here  $\mathcal{T}$  is a (non-uniform) triangulation of  $\Omega$ , and  $h_T$  denotes the diameter of a triangle  $T$ . System (3.82) is then given by

$$\begin{pmatrix} S & -D^\top \\ D & k\Delta t L + C_s \end{pmatrix} \begin{pmatrix} \mathbf{u}^m \\ \mathbf{p}^m \end{pmatrix} = \begin{pmatrix} 0 & 0 \\ D & C_s \end{pmatrix} \begin{pmatrix} \mathbf{u}^{m-1} \\ \mathbf{p}^{m-1} \end{pmatrix} + \begin{pmatrix} \mathbf{g}^m + \boldsymbol{\tau}^m \\ \Delta t \mathbf{f}^m \end{pmatrix}. \quad (3.85)$$

Note that if  $\mathcal{T}$  is uniform, we again obtain  $C_s = \beta h^2 L$ .

### 3.3.4 Comparison Between Bubble- and Diffusive Stabilization

We test both bubble stabilization (3.76) and diffusive stabilization (3.85) and compare the corresponding finite element solutions for  $p$  to the case where no stabilization is applied. To this end we solve equations (3.47)-(3.48) on  $\Omega = (0, 1)^2$ . We set  $\mathbf{g} = (0, -1)^\top$  and  $\mathbf{f} = 0$ . Concerning the boundary conditions (3.50)-(3.51), we let  $\Gamma_2 = \{(x, y) : x = 0, y \in [0, 1]\}$  and  $\Gamma_1 = \partial\Omega \setminus \Gamma_2$ . The stress vector  $\boldsymbol{\tau}$  on  $\Gamma_2$  is set to  $\mathbf{0}$ . The parameter  $k$  is made small to cause oscillations in the non-stabilized solution: we set  $k = 10^{-8}$ . To ensure only small deformations occur, the Lamé constants are set to  $\mu_\varepsilon = \lambda_\varepsilon = 1000$ . A uniform triangulation is used with 21 grid nodes in each coordinate direction. Thus the meshwidth is equal to  $h = 0.05\sqrt{2}$ . A timestep of  $\Delta t = 10^{-2}$  is used.

Figures 3.11, 3.12 and 3.13 show the finite element solutions for  $p$  using respectively no stabilization, bubble stabilization and diffusive stabilization. That is, in Figure 3.11 the system in (3.62) is solved, in Figure 3.12 the system in (3.76) is solved, and in Figure 3.13 the system in (3.85) is solved.

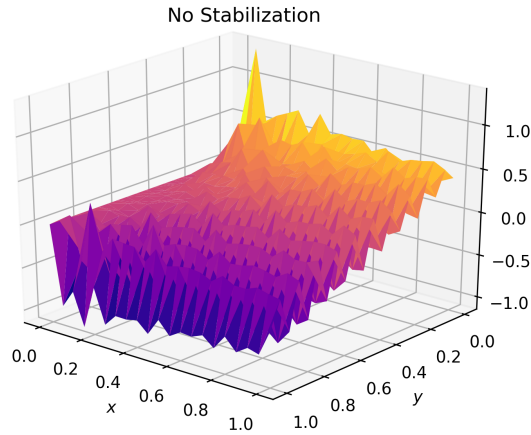


Figure 3.11: Finite element solution for  $p$  after one timestep of size  $\Delta t = 10^{-2}$ , using linear-linear elements.

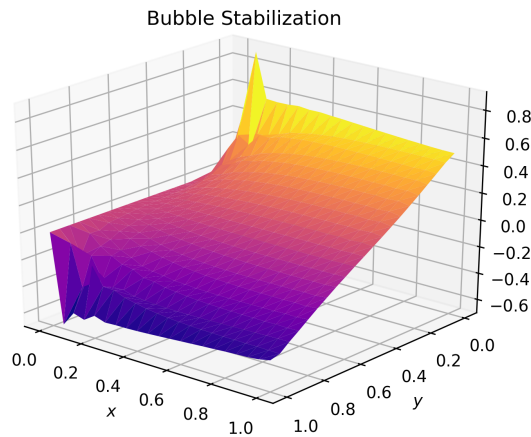


Figure 3.12: Finite element solution for  $p$  after one timestep of size  $\Delta t = 10^{-2}$ , using linear-linear elements and bubble stabilization.

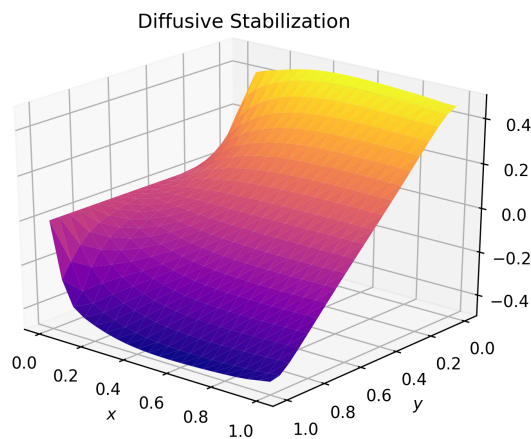


Figure 3.13: Finite element solution for  $p$  after one timestep of size  $\Delta t = 10^{-2}$ , using linear-linear elements and diffusive stabilization.

It is clear that diffusive stabilization does its job well; the solution is completely smooth. More interesting is the behaviour of the bubble stabilization. In the interior of the domain the solution is smooth, but some oscillations occur near the boundary.



### 3.4 Visco-Poroelasticity on a Deforming Domain

The visco-poroelastic system is given by the following two equations:

$$\rho \left( \frac{D\mathbf{v}}{Dt} + (\nabla \cdot \mathbf{v})\mathbf{v} \right) - \nabla \cdot \boldsymbol{\sigma}(\mathbf{u}, \mathbf{v}, p) = \mathbf{g} \quad \text{on } \Omega(t), \quad (3.86)$$

$$\nabla \cdot \mathbf{v} - \nabla \cdot (k\nabla p) = f \quad \text{on } \Omega(t). \quad (3.87)$$

The domain  $\Omega(t)$  again depends on the displacement vector  $\mathbf{u}$  at time  $t$ . Furthermore,  $\mathbf{v}$  is the material derivative of  $\mathbf{u}$ . We impose the following boundary conditions:

$$\mathbf{u} = \mathbf{0}, \quad k\nabla p \cdot \mathbf{n} = 0, \quad \text{on } \Gamma_1(t), \quad (3.88)$$

$$\boldsymbol{\sigma} \cdot \mathbf{n} = \boldsymbol{\tau}, \quad p = 0, \quad \text{on } \Gamma_2(t), \quad (3.89)$$

The stress tensor is given by

$$\boldsymbol{\sigma}(\mathbf{u}, \mathbf{v}, p) = \boldsymbol{\sigma}_{\text{el}}(\mathbf{u}) + \boldsymbol{\sigma}_{\text{vis}}(\mathbf{v}) - p\mathbf{I}, \quad (3.90)$$

where the elastic and viscous parts are given respectively by (2.38) and (2.39).

#### 3.4.1 Derivation of System of Equations

We derive the weak formulations of (3.86) and (3.87). Let  $\boldsymbol{\varphi}(\cdot, t)$  be a Lagrangian test field on  $\Omega(t)$  that vanishes on  $\Gamma_1(t)$ . The weak form of equation (3.86) is given by:

$$\begin{aligned} \rho \frac{d}{dt} \int_{\Omega(t)} \mathbf{v} \cdot \boldsymbol{\varphi} \, d\Omega + \int_{\Omega(t)} \boldsymbol{\sigma}(\mathbf{u}, \mathbf{v}) : \boldsymbol{\varepsilon}(\boldsymbol{\varphi}) \, d\Omega - \int_{\Omega(t)} p(\nabla \cdot \boldsymbol{\varphi}) \, d\Omega \\ = \int_{\Omega(t)} \mathbf{g} \cdot \boldsymbol{\varphi} \, d\Omega + \int_{\Gamma_2(t)} \boldsymbol{\tau} \cdot \boldsymbol{\varphi} \, d\Gamma. \end{aligned} \quad (3.91)$$

We have used Theorem 3 to simplify the first term of (3.86). Let  $\psi(\cdot, t)$  be a scalar-valued test field on  $\Omega(t)$  that vanished on  $\Gamma_2(t)$ . Then the weak form of equation (3.87) is given by:

$$\int_{\Omega(t)} (\nabla \cdot \mathbf{v})\psi \, d\Omega + \int_{\Omega(t)} k\nabla p \cdot \nabla \psi \, d\Omega = \int_{\Omega(t)} f\psi \, d\Omega. \quad (3.92)$$

Let  $\mathcal{T}_h(t)$  be a triangulation of the approximate domain  $\Omega_h(t)$ . We use linear basis functions with respect to the triangles in  $\mathcal{T}_h(t)$  for  $\mathbf{u}$ ,  $\mathbf{v}$  and  $\mathbf{p}$ . We then obtain the following Galerkin equations corresponding to the weak forms (3.91) and (3.92):

$$\rho \frac{d}{dt} (M_v \mathbf{v}) + S_{\text{el}} \mathbf{u} + S_{\text{vis}} \mathbf{v} - D^\top \mathbf{p} = \mathbf{g}, \quad (3.93)$$

$$D\mathbf{v} + kL\mathbf{p} = \mathbf{f}. \quad (3.94)$$

Note that  $\mathbf{u}$ ,  $\mathbf{v}$  and  $\mathbf{p}$  now denote numerical vectors containing the values of their corresponding unknowns in the grid points of  $\mathcal{T}_h(t)$ . Moreover, the vector  $\mathbf{g}$  contains both the domain and boundary integral in the right-hand side of (3.91). We write  $\mathbf{u} = (\mathbf{u}^1, \mathbf{u}^2)^\top$  and  $\mathbf{v} = (\mathbf{v}^1, \mathbf{v}^2)^\top$ , then the matrices  $M_v$ ,  $S_{\text{el}}$ ,  $S_{\text{vis}}$  and  $D$  can be subdivided into blocks:

$$M_v = \begin{pmatrix} M & 0 \\ 0 & M \end{pmatrix}, \quad S_{\text{el}} = \begin{pmatrix} S_{\text{el}}^{11} & S_{\text{el}}^{12} \\ S_{\text{el}}^{21} & S_{\text{el}}^{22} \end{pmatrix}, \quad S_{\text{vis}} = \begin{pmatrix} S_{\text{vis}}^{11} & S_{\text{vis}}^{12} \\ S_{\text{vis}}^{21} & S_{\text{vis}}^{22} \end{pmatrix}, \quad D = (D^1 \quad D^2). \quad (3.95)$$

For component-wise definitions of the  $S_{\text{el}}$  and  $S_{\text{vis}}$  blocks, we refer to 2.3.4. The divergence blocks are given by (3.67). The matrix  $L$  is a Laplacian matrix, it is given by (3.71). Finally the mass matrix  $M$  is given element-wise by:

$$M_{ij} = \int_{\Omega(t)} \varphi_i \varphi_j \, d\Omega, \quad (3.96)$$

where  $\varphi_i$  is a basis function corresponding to the  $i$ 'th grid point in  $\mathcal{T}_h(t)$ . We discretize equations (3.93) and (3.94) in time using implicit Euler:

$$\rho M_v^m \mathbf{v}^m + \Delta t S_{\text{el}}^m \mathbf{u}^m + \Delta t S_{\text{vis}}^m \mathbf{v}^m - \Delta t (D^m)^\top \mathbf{p}^m = \rho M_v^{m-1} \mathbf{v}^{m-1} + \Delta t \mathbf{g}^m \quad (3.97)$$

$$D^m \mathbf{v}^m + kL^m \mathbf{p}^m = \mathbf{f}^m. \quad (3.98)$$

The superscript  $(\cdot)^m$  denotes a time-level. Note that the coefficient matrices also depend on time because their entries are integrations over the current domain. To eliminate  $\mathbf{u}^m$  from the system we substitute

$$\mathbf{u}^m = \mathbf{u}^{m-1} + \Delta t \mathbf{v}^m, \quad (3.99)$$

into (3.97) to obtain:

$$(\rho M_v^m + \Delta t S_{\text{vis}}^m + \Delta t^2 S_{\text{el}}^m) \mathbf{v}^m - \Delta t (D^m)^\top \mathbf{p}^m = \rho M_v^{m-1} \mathbf{v}^{m-1} - \Delta t S_{\text{el}}^m \mathbf{u}^{m-1} + \Delta t \mathbf{g}^m \quad (3.100)$$

$$D^m \mathbf{v}^m + kL^m \mathbf{p}^m = \mathbf{f}^m. \quad (3.101)$$

In block matrix form we get:

$$\begin{pmatrix} \rho M_v^m + \Delta t S_{\text{vis}}^m + \Delta t^2 S_{\text{el}}^m & -\Delta t (D^m)^\top \\ D^m & kL^m \end{pmatrix} \begin{pmatrix} \mathbf{v}^m \\ \mathbf{p}^m \end{pmatrix} = \begin{pmatrix} \rho M_v^{m-1} & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} \mathbf{v}^{m-1} \\ \mathbf{p}^{m-1} \end{pmatrix} + \begin{pmatrix} \Delta t (\mathbf{g}^m - S_{\text{el}}^m \mathbf{u}^{m-1}) \\ \mathbf{f}^m \end{pmatrix}. \quad (3.102)$$

The above system must be solved using a Picard-type iterative scheme. We use the following convergence criterion:

$$\frac{\|\mathbf{v}^{\ell+1} - \mathbf{v}^\ell\|}{\|\mathbf{v}^\ell\|} + \frac{\|\mathbf{p}^{\ell+1} - \mathbf{p}^\ell\|}{\|\mathbf{p}^\ell\|} \leq \delta. \quad (3.103)$$

Here  $\ell$  does *not* denote a time-level, but a sub-iteration within one timestep. We find that choosing  $\delta$  between  $10^{-4}$  and  $10^{-8}$  produces good results, depending on the meshwidth.

### 3.4.2 Stabilization

We observe oscillations in the pressure solution for small values of  $k$ . We investigate what is causing this behaviour. Write system (3.102) in the following abbreviated form:

$$\begin{pmatrix} Q & -\Delta t D^\top \\ D & kL \end{pmatrix} \begin{pmatrix} \mathbf{v} \\ \mathbf{p} \end{pmatrix} = \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix} \quad (3.104)$$

Here we have dropped the time-level superscripts and called the right-hand side vectors  $\mathbf{a}$  and  $\mathbf{b}$  respectively. Moreover,  $Q = \rho M_v + \Delta t S_{\text{vis}} + \Delta t^2 S_{\text{el}}$ . Rearrange the first equation to

$$\mathbf{v} = \Delta t Q^{-1} D^\top \mathbf{p} + Q^{-1} \mathbf{a}. \quad (3.105)$$

Substitute this expression into the second equation to find

$$(kL + \Delta t D Q^{-1} D^\top) \mathbf{p} = \mathbf{b} - D Q^{-1} \mathbf{a}. \quad (3.106)$$

As the oscillations occur for small values of  $k$ , it appears that the matrix  $\Delta t D Q^{-1} D^\top$  has unfavorable properties. Note that the oscillations can also be fixed by reducing  $\Delta t$ . This solution quickly becomes computationally unfeasible due the high number of time steps necessary. In Chapter 4 we construct a diffusive stabilization scheme that is optimal in a certain sense.

### 3.4.3 Numerical Experiments

We test the model in a similar situation to the one in section 2.3.6. Let  $\Omega(0) = (0, 1)^2$  and set  $\Gamma_1(0) = \{(x, y) : x = 0\}$ . Then from boundary condition (3.89) it is clear that  $\Gamma_1(t) = \Gamma_1(0)$  for all  $t > 0$ . Furthermore, we let  $\Gamma_2(t) = \partial\Omega(t) \setminus \Gamma_1(t)$ . The body force is set to  $\mathbf{g} = (0, -0.1)^\top$ , the pressure source is set to zero:  $f = 0$ , and the force on  $\Gamma_2(t)$  is also set to zero:  $\boldsymbol{\tau} = \mathbf{0}$ . The following parameter values are used:

$$\rho = 1, \quad \lambda_\varepsilon = \mu_\varepsilon = \lambda_v = \mu_v = 1, \quad k = 10^{-2}. \quad (3.107)$$

We use a meshwidth of  $h = 0.1$  and a timestep of  $\Delta t = 0.05$ . Figures 3.14a - 3.14d show the displacement and pressure solution at various timesteps. Note that the pressure is shown as a contour plot superimposed on the displaced mesh. The initial mesh is shown for reference in red.

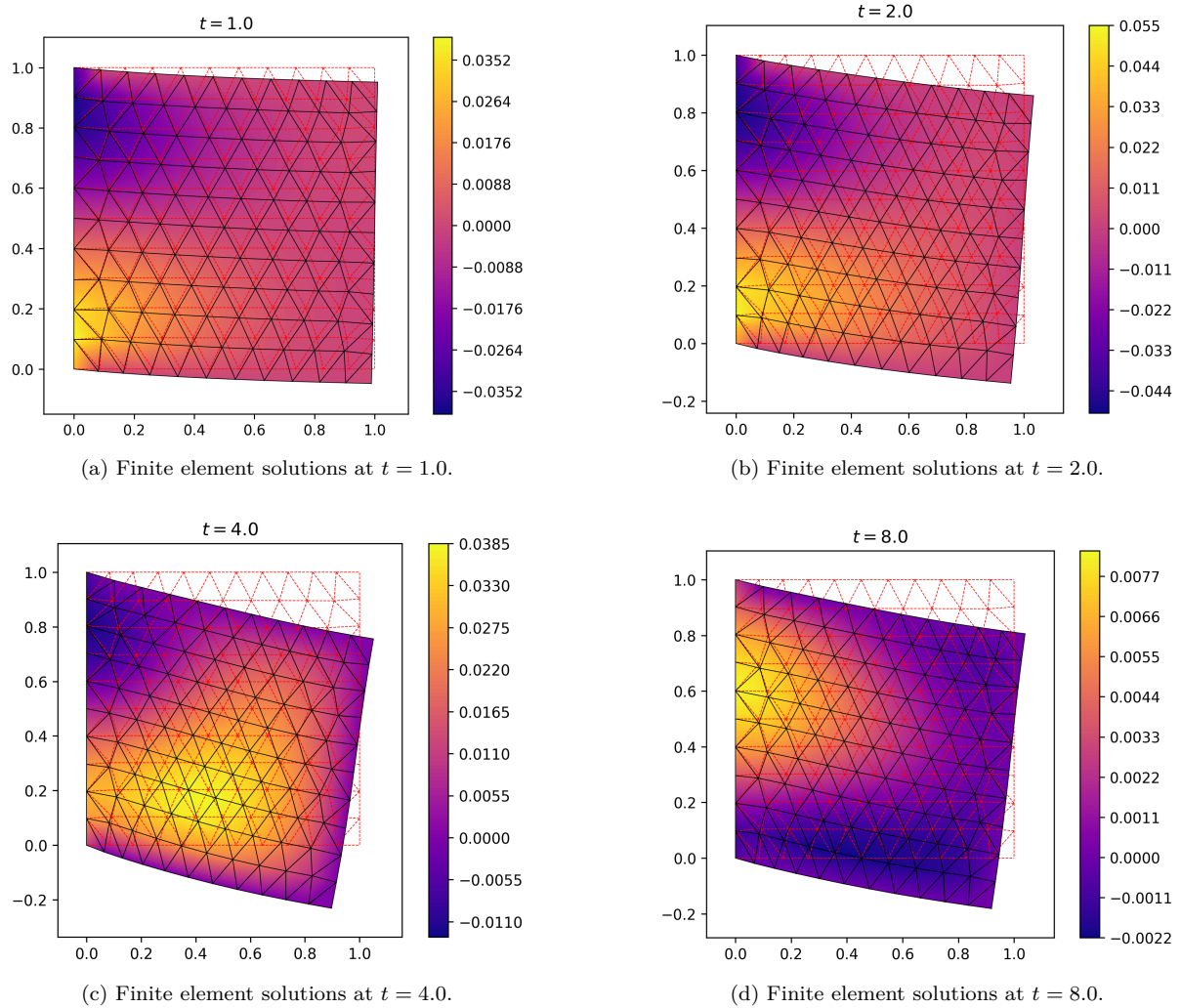
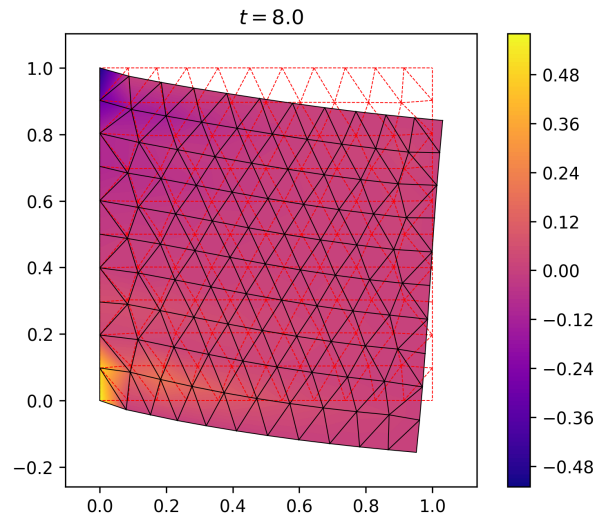


Figure 3.14: Plots showing the finite element solution at various points in time.

Note that the mesh displacement looks very similar to those in Figures 2.6a - 2.6d. This is due to a relatively high values of  $k$ : the pressure does not have a big effect on the displacement (velocity). Figure 3.15 shows the finite element solutions at  $t = 8$  for a lower value  $k = 10^{-5}$ , whereas the other parameters are kept constant. Compared to Figures 3.14d and 2.6d, we can see a small difference in displacement. Note that lowering the value of  $k$  even further will result in more severe oscillations. They can already be seen forming near the top and bottom left corners.

Figure 3.15: Finite element solutions at  $t = 8$  for  $k = 10^{-5}$ 

### 3.5 Comparison Between Viscoelastic and Visco-Poroelastic Models

We compare the behaviour of the visco-poroelastic model from Section 3.4 to the viscoelastic model from Section 2.3. We again apply the body force

$$\mathbf{g}(\mathbf{x}, t) = \begin{pmatrix} 0.2 \sin(\pi t) \mathbb{1}_{[0,1]}(t) \\ 0 \end{pmatrix} \quad (3.108)$$

to the unite square like in Section 2.5. The mass source function  $f$  is set to zero, and the following mechanical parameters are used:

$$\mu_\varepsilon = \lambda_\varepsilon = 1, \quad \mu_v = \lambda_v = 0.5, \quad \rho = 0.5. \quad (3.109)$$

We use the (initially) uniform finite element mesh from Figure 2.10 and a timestep of size  $\Delta t = 0.05$ . Figures 3.16a and 3.16b show the maximum width and minimum height of the domain respectively as functions of time. Two different values of  $k$  are used. For  $k = 10^{-5}$ , the stabilized model is used. One can clearly see that the inclusion of poroelasticity has an effect on the evolution of the domain. This effect is more pronounced for smaller values of  $k$ .

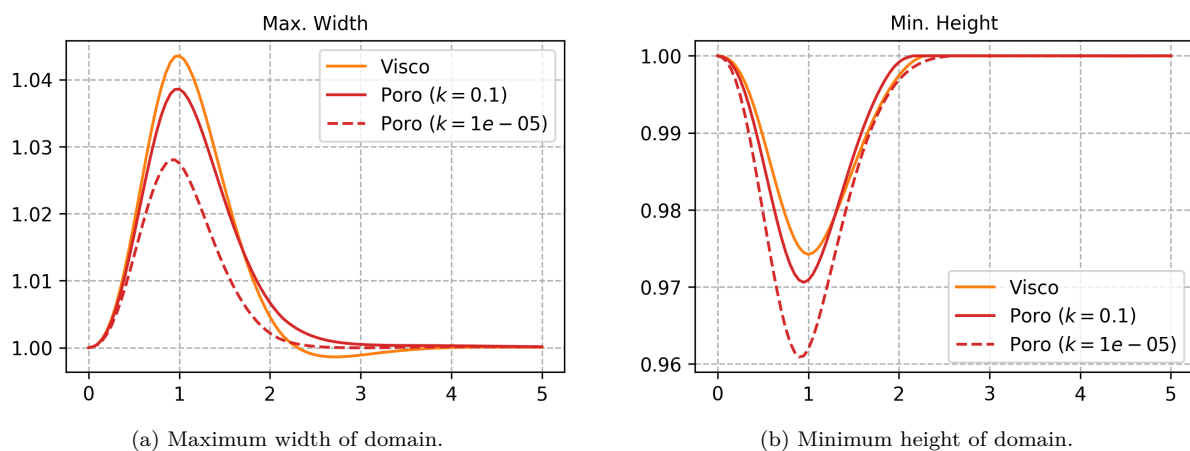


Figure 3.16: Maximum domain width and minimum domain height as a function of time for two visco-poroelastic models compared to the viscoelastic model.

We give an intuitive explanation from a physical point of view for the effect of  $k$  on the solution. Recall that  $k$  is the quotient of the permeability of the medium and the dynamic viscosity of the poro-fluid. If

$k$  is large, then either the medium is very permeable or the fluid has a low viscosity. In both cases, the fluid is able to move relatively freely. Thus, whenever the medium is deformed, the fluid will quickly adapt to the new situation. On the other hand, if  $k$  is small, either the permeability is low or the fluid has a high viscosity. In this case the fluid has more trouble to move through the medium, hence it will react slower to deformations. This inability of the poro-fluid to adapt to the new situation causes it have a more pronounced effect on the displacement.

## Part II

# 4

## Morpho-Poroelasticity and Analysis of Stabilized Finite Element Methods

In this chapter we present new work. We first analyze the properties of diffusive stabilization applied to poroelastic systems. The aim is to give a rigorous derivation of the optimal tuning parameter  $\beta$ . This derivation is performed in Sections 4.2, 4.3 and 4.4 for a one-dimensional visco-poroelastic system. For a two- or higher dimensional systems, the computations become too hard to do exactly. For this reason, the result obtained in the one-dimensional case is ‘extrapolated’ to higher dimensional poroelastic systems. The second part of this chapter contains the construction of a finite element model for a morpho-poroelastic system. Diffusive stabilization from earlier sections is also applied to this model, since the problematic equation has a very similar structure.

### 4.1 Stabilizing Poroelastic Systems

For small  $k$ , the visco-poroelastic system as presented in (3.102), as well as the morpho-poroelastic system (4.76) will start to show non-physical oscillations in the pressure solution. To eliminate them a diffusive stabilization term can be added to the pressure equation, such as we have done for the poroelastic system in Section 3.3.3. Note that in both the visco-poroelastic and morpho-poroelastic systems the last equation in the fully discretized system is given by

$$D^m \mathbf{v}^m + kL^m \mathbf{p}^m = \mathbf{f}^m. \quad (4.1)$$

As we have seen in Section 3.3.3, the idea is to introduce a small diffusive perturbation. This causes the solution to remain smooth. It is possible to only add this perturbation to the left-hand side of (4.1), effectively increasing  $k$ . Equation (4.1) is then replaced by

$$D^m \mathbf{v}^m + (k + \beta)L^m \mathbf{p}^m = \mathbf{f}^m. \quad (4.2)$$

Another option would be to add a perturbation to both sides, thereby approximating some kind of time derivative of the Laplacian of  $p$ . In this case, equation (4.1) becomes

$$D^m \mathbf{v}^m + (k + \beta)L^m \mathbf{p}^m = \mathbf{f}^m + \beta L^{m-1} \mathbf{p}^{m-1}. \quad (4.3)$$

Both (4.2) and (4.3) will have a similar effect on the stability of the solution, since that is determined by the matrix  $(k + \beta)L$ . The two approaches differ from each other in the right-hand side, and therefore will lead to different errors compared to (4.1). To gain insight in the error, we determine the differential equations corresponding to (4.2) and (4.3). Recall that (4.1) is the discretized version of the differential equation

$$\nabla \cdot \mathbf{v} - k \nabla^2 p = f. \quad (4.4)$$

Thus, we can immediately tell that (4.2) discretizes the following perturbed differential equation:

$$\nabla \cdot \mathbf{v} - k \nabla^2 p = f + \beta \nabla^2 p. \quad (4.5)$$

The corresponding differential equation of (4.3) is less obvious. See the following theorem:

**Theorem 6.** Assume that  $\mathbf{v}$  and  $\mathbf{p}$  satisfy

$$\nabla \cdot \mathbf{v} - k\nabla^2 p = f + \Delta t \beta \left( \frac{D(\nabla^2 p)}{Dt} + (\nabla \cdot \mathbf{v})(\nabla^2 p) \right). \quad (4.6)$$

Then the corresponding finite element solution vectors  $\mathbf{v}^m$  and  $\mathbf{p}^m$  satisfy (4.3), provided Lagrangian basis functions are used and the time discretization is done using implicit Euler.

*Proof.* Assume we have homogeneous boundary conditions on  $p$ , such as in (3.88)-(3.89). Let  $q$  be a Lagrangian basis function. Then the weak form of (4.6) is given by

$$\int_{\Omega(t)} (\nabla \cdot \mathbf{v})q \, d\Omega + \int_{\Omega(t)} k\nabla p \cdot \nabla q \, d\Omega = \int_{\Omega(t)} f q \, d\Omega + \Delta t \beta \int_{\Omega(t)} \left( \frac{D(\nabla^2 p)}{Dt} + (\nabla \cdot \mathbf{v})(\nabla^2 p) \right) q \, d\Omega.$$

The last term can be simplified using Theorem 3:

$$\int_{\Omega(t)} (\nabla \cdot \mathbf{v})q \, d\Omega + \int_{\Omega(t)} k\nabla p \cdot \nabla q \, d\Omega = \int_{\Omega(t)} f q \, d\Omega + \Delta t \beta \frac{d}{dt} \int_{\Omega(t)} (\nabla^2 p)q \, d\Omega.$$

By applied partial integration once more, we obtain

$$\int_{\Omega(t)} (\nabla \cdot \mathbf{v})q \, d\Omega + \int_{\Omega(t)} k\nabla p \cdot \nabla q \, d\Omega = \int_{\Omega(t)} f q \, d\Omega - \Delta t \beta \frac{d}{dt} \int_{\Omega(t)} \nabla p \cdot \nabla q \, d\Omega.$$

By computing the Galerkin equations, the semi-discretized system is obtained:

$$D\mathbf{v} + kL\mathbf{p} = \mathbf{f} - \Delta t \beta \frac{d}{dt}(L\mathbf{p}),$$

where  $D$  is the divergence matrix and  $L$  is the Laplacian matrix. Applying implicit Euler yields

$$D^m \mathbf{v}^m + kL^m \mathbf{p}^m = \mathbf{f}^m - \beta (L^m \mathbf{p}^m - L^{m-1} \mathbf{p}^{m-1}).$$

Equation (4.3) is obtained by rearranging terms. □

To summarize, we have found that adding a one-sided stabilization to the fully discretized system such as in (4.2) is equivalent to solving (4.5) using finite elements. Similarly, adding the two-sided stabilization such as in (4.3) amounts to solving (4.6). Assume the system converges to a steady state, thus  $\mathbf{v} \rightarrow \mathbf{0}$  as  $t \rightarrow \infty$ . Then the stationary pressure equations corresponding to both (4.4) and (4.6) become

$$-k\nabla^2 p = f. \quad (4.7)$$

Whereas for (4.5) it becomes

$$-(k + \beta)\nabla^2 p = f. \quad (4.8)$$

Based on this, we argue that a two-sided stabilization is the way to go.

## 4.2 Analysis of One-Dimensional Visco-Poroelastic System

In both (4.2) and (4.3) the stability of the solution is affected by the tuning parameter  $\beta$ . If  $\beta$  is too small the solution will still contain oscillations. On the other hand, if  $\beta$  is too big the solution will be excessively smoothed, leading to a large error. In the following sections we derive an optimal tuning parameter  $\beta$  for a one-dimensional visco-poroelastic system on a fixed, uniform grid. This optimal choice of  $\beta$  is valid in the limit of  $h \rightarrow 0$ , where  $h$  is the meshwidth. In this framework system (3.102) becomes

$$\begin{pmatrix} \rho M_v + (\Delta t(\mu_v + \lambda_v) + \Delta t^2(\mu_\varepsilon + \lambda_\varepsilon))L_v & -\Delta t D^\top \\ D & kL_p \end{pmatrix} \begin{pmatrix} \mathbf{v}^m \\ \mathbf{p}^m \end{pmatrix} = \begin{pmatrix} \rho M_v & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} \mathbf{v}^{m-1} \\ \mathbf{p}^{m-1} \end{pmatrix} + \begin{pmatrix} \Delta t(\mathbf{g}^m - L_v \mathbf{u}^{m-1}) \\ \mathbf{f}^m \end{pmatrix}. \quad (4.9)$$

where  $S_{\text{vis}}$  is replaced with  $(\mu_v + \lambda_v)L_v$  and  $S_{\text{el}}$  with  $(\mu_\varepsilon + \lambda_\varepsilon)L_v$ , where  $L_v$  is a Laplace matrix incorporating the boundary conditions imposed on  $v$ . We define

$$\theta = \Delta t(\mu_v + \lambda_v) + \Delta t^2(\mu_\varepsilon + \lambda_\varepsilon), \quad \text{and} \quad A = \rho M_v + \theta L_v. \quad (4.10)$$



By isolating  $\mathbf{v}^m$  in the first equation of (4.9) and substituting the result into the second, we obtain the following evolution equation for  $\mathbf{p}^m$ :

$$(kL_p + \Delta tDA^{-1}D^\top)\mathbf{p}^m = \mathbf{f}^m - \rho DA^{-1}M_v\mathbf{v}^{m-1} + \Delta tDA^{-1}(L_v\mathbf{u}^{m-1} - \mathbf{g}^m). \quad (4.11)$$

Let  $C = DA^{-1}D^\top$ , then the coefficient matrix  $kL_p + \Delta tC$  is important: it causes the notorious oscillations in  $p$ . From numerical tests we know that oscillations only occur for small values of  $k$ . Since  $L_p$  is an M-matrix, the matrix  $\Delta tDA^{-1}D^\top$  must be the cause of the oscillations. Equation (4.11) shows that the effect of the problematic matrix can be reduced by decreasing  $\Delta t$ . Numerical tests show that the system is again stable if we have  $\Delta t \leq k$  whenever  $k$  is small. In this section the matrix  $C$  is approximated. With this approximation an approximate optimal tuning parameter  $\beta$  can be determined.

### 4.3 Approximating the Inverse of the Discrete Reaction-Laplacian Operator

The inverse of  $A$  is approximated for a certain set of boundary conditions. To this end, we consider a related boundary value problem with a delta function as source. This problem can be solved exactly but also using finite elements. By using the latter method we encounter simplified versions of the matrices  $M_v$  and  $L_v$ . Comparing the exact solution to the finite element solution allows us to approximate  $A^{-1}$ . Consider the following reaction-Laplacian equation and corresponding boundary conditions:

$$\rho u - \theta u'' = \delta(x - \alpha), \quad x \in [0, 1] \quad (4.12)$$

$$u'(0) = 0, \quad (4.13)$$

$$u(1) = 0, \quad (4.14)$$

where  $\delta(\cdot)$  denotes the Dirac delta-function, and  $\alpha \in [0, 1)$ . Moreover, the parameters  $\rho$  and  $\theta$  are strictly positive. Using the Laplace transform, the following exact (weak) solution is obtained:

$$\begin{aligned} u(x) &= \frac{1}{\sqrt{\rho\theta}} \left[ \frac{\sinh(\lambda(1-\alpha))}{\cosh(\lambda)} \cosh(\lambda x) - H(x-\alpha) \sinh(\lambda(x-\alpha)) \right] \\ &= \frac{\lambda}{\rho \cosh(\lambda)} \begin{cases} \sinh(\lambda(1-\alpha)) \cosh(\lambda x), & x \leq \alpha \\ \sinh(\lambda(1-x)) \cosh(\lambda\alpha), & x > \alpha \end{cases} \end{aligned} \quad (4.15)$$

where  $\lambda = \sqrt{\rho/\theta}$ . Note that (4.15) is not a classical solution, since it is not twice continuously differentiable on  $(0, 1)$ . Figure 4.1 shows a plot of this solution for certain values of  $\alpha$ ,  $\rho$  and  $\theta$ .

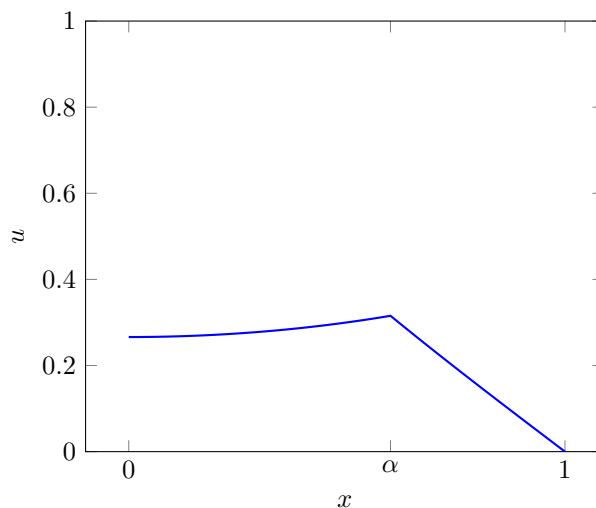


Figure 4.1: Plot of exact solution as given in (4.15) for  $\alpha = 0.6$  and  $\rho = \theta = 1$ .

#### 4.3.1 Existence and Uniqueness of Weak Solution to Reaction-Laplacian Equation

We prove that a solution to the reaction-Laplacian problem (4.12)-(4.14) exists in the weak sense. Define the function space

$$H_0^1(0, 1) = \{f \in H^1(0, 1) : f(1) = 0\}, \quad (4.16)$$

which inherits its norm from  $H^1(0, 1)$ :

$$\|f\|_{H^1}^2 = \|f\|_{L^2}^2 + \|f'\|_{L^2}^2 = \int_0^1 f^2 + (f')^2 dx.$$

The weak form corresponding to problem (4.12)-(4.14) is given by:

Find  $u \in H_0^1(0, 1)$  such that for all  $v \in H_0^1(0, 1)$  we have:

$$a(u, v) = \delta_\alpha(v), \quad (4.17)$$

where

$$a(u, v) = \int_0^1 \rho uv + \theta u'v' dx, \quad \text{and} \quad \delta_\alpha(v) = v(\alpha). \quad (4.18)$$

To show that a solution can be obtained using the finite element method, it must first be shown that the weak form (4.17) has a unique solution. This is summarized in the following theorem:

**Theorem 7.** *There exists a unique  $u \in H_0^1(0, 1)$  that satisfies*

$$a(u, v) = \delta_\alpha(v)$$

for all  $v \in H_0^1(0, 1)$ , where  $a(\cdot, \cdot)$  and  $\delta_\alpha(\cdot)$  are given in (4.18).

*Proof.* The statement can be proved using the Lax-Milgram theorem. First, the necessary conditions on  $a$  are proved. Note that we view  $a(\cdot, \cdot)$  as the bilinear form

$$a : H_0^1(0, 1) \times H_0^1(0, 1) \rightarrow \mathbb{R}.$$

It is clear from its definition that  $a$  is symmetric. Furthermore,

- $a$  is bounded (and hence continuous). Indeed:

$$\begin{aligned} |a(u, v)| &= \left| \int_0^1 \rho uv + \theta u'v' dx \right| \\ &\leq \rho \int_0^1 |uv| dx + \theta \int_0^1 |u'v'| dx \\ &\leq \rho \sqrt{\int_0^1 u^2 dx} \sqrt{\int_0^1 v^2 dx} + \theta \sqrt{\int_0^1 (u')^2 dx} \sqrt{\int_0^1 (v')^2 dx} \\ &= \rho \|u\|_{L^2} \|v\|_{L^2} + \theta \|u'\|_{L^2} \|v'\|_{L^2} \\ &\leq (\rho + \theta) \|u\|_{H^1} \|v\|_{H^1}. \end{aligned}$$

In the third step the Cauchy-Schwarz inequality is applied.

- $a$  is coercive:

$$\begin{aligned} a(u, u) &= \int_0^1 \rho u^2 + \theta (u')^2 dx \\ &= \rho \|u\|_{L^2}^2 + \theta \|u'\|_{L^2}^2 \\ &\geq \min\{\rho, \theta\} (\|u\|_{L^2}^2 + \|u'\|_{L^2}^2) \\ &= \min\{\rho, \theta\} \|u\|_{H^1}^2. \end{aligned}$$

In order for Lax-Milgram to be applicable, the right-hand side of (4.17) must be continuous. Note that we view  $\delta_\alpha(\cdot)$  as a functional acting on  $H_0^1(0, 1)$ . Before continuity is proved, the issue of  $\delta_\alpha$  being well-defined must be addressed. After all, since  $v \in H_0^1(0, 1) \subset L^2(0, 1)$ , pointwise evaluation of  $v$  is not defined. Let  $(a, b)$  be a real interval, then we have the inclusion<sup>1</sup>  $H^1(a, b) \subset AC(a, b)$ , where  $AC(a, b)$  is the set of absolutely continuous functions on  $(a, b)$ . A proof of this fact is provided in [2]. A well known

<sup>1</sup>The inclusion holds almost everywhere: for each  $f \in H^1(a, b)$  there is an  $\tilde{f} \in AC(a, b)$  such that  $f = \tilde{f}$  a.e.

property of absolutely continuous functions is that the fundamental theorem of integration holds [9]. Let  $f \in AC(a, b)$ , then for  $x \in (a, b)$ :

$$f(x) = f(a) + \int_a^x f'(t) dt,$$

where  $f'$  is a weak derivative of  $f$ . We can now prove the continuity of  $\delta_\alpha$ . It suffices to show  $\delta_\alpha$  is bounded due to its linearity. Let  $u \in H_0^1(0, 1)$ , then we have

$$|\delta_\alpha(u)| = |u(\alpha)| = \left| u(0) + \int_0^\alpha u' dx \right|$$

To eliminate  $u(0)$ , note that  $u(1) = 0$ , and therefore:

$$u(0) + \int_0^1 u' dx = 0.$$

It follows that:

$$|\delta_\alpha(u)| = \left| \int_0^\alpha u' dx - \int_0^1 u' dx \right| = \left| \int_\alpha^1 u' dx \right| \leq \int_\alpha^1 |u'| dx$$

Apply the Cauchy-Schwarz inequality:

$$\begin{aligned} \int_\alpha^1 |u'| dx &\leq \sqrt{1-\alpha} \sqrt{\int_\alpha^1 (u')^2 dx} \\ &\leq \sqrt{1-\alpha} \sqrt{\int_0^1 (u')^2 dx} \\ &= \sqrt{1-\alpha} \|u'\|_{L^2} \\ &\leq \sqrt{1-\alpha} \|u\|_{H^1} \end{aligned}$$

Which completes the proof of the boundedness of  $\delta_\alpha$ . All conditions of the Lax-Milgram theorem are met, thus problem (4.17) has a unique solution.  $\square$

### 4.3.2 Error Analysis

In finite elements, the solution- and test space (which are equal in our case) are replaced by finite dimensional spaces. Let  $V = H_0^1(0, 1)$  be the original test space, and let  $V_h$  be a finite dimensional subspace of  $V$ . Note that  $V_h$  will be specified later on. In the remainder of this section, we derive the  $L^2$ -error of the finite element solution. Assume that  $u \in V$  satisfies

$$a(u, v) = \delta_\alpha(v), \quad \forall v \in V, \quad (4.19)$$

and  $u_h \in V_h$  satisfies

$$a(u_h, v_h) = \delta_\alpha(v_h), \quad \forall v_h \in V_h. \quad (4.20)$$

It follows from the proof of Theorem 7 that  $a(\cdot, \cdot)$  defines an inner product on  $V$ . We say that  $v, w \in V$  are  $a$ -orthogonal if  $a(v, w) = 0$ . Note that the finite element error  $u - u_h$  is  $a$ -orthogonal to the test space  $V_h$ :

$$a(u - u_h, v_h) = 0, \quad \forall v_h \in V_h. \quad (4.21)$$

Indeed:

$$\begin{aligned} a(u - u_h, v_h) &= a(u, v_h) - a(u_h, v_h) \\ &= \delta_\alpha(v_h) - \delta_\alpha(v_h) = 0 \end{aligned}$$

Before the  $L^2$  norm of the error can be computed, its energy norm is analyzed.

**Definition 4.** The  $a$ -induced *energy norm* is defined by:

$$\|u\|_a^2 = a(u, u) \quad (4.22)$$

We have the following well-known result:

**Lemma 5** (Céa). *Let  $u \in V$  and  $u_h \in V_h$  satisfy (4.19) and (4.20) respectively. Then:*

$$\|u - u_h\|_a \leq \|u - v_h\|_a, \quad \forall v_h \in V_h. \quad (4.23)$$

*Proof.* Let  $v_h \in V_h$  be arbitrary. The proof follows from the  $a$ -orthogonality property of  $u - u_h$ :

$$\begin{aligned} \|u - u_h\|_a^2 &= a(u - u_h, u - u_h) \\ &= a(u - u_h, u - v_h + v_h - u_h) \\ &= a(u - u_h, u - v_h) + a(u - u_h, u_h - v_h) \\ &= a(u - u_h, u - v_h) \\ &\leq \|u - u_h\|_a \|u - v_h\|_a \end{aligned}$$

We have used the fact that  $v_h - u_h \in V_h$ . In the last step, the Cauchy-Schwarz inequality is applied. Dividing both sides by  $\|u - u_h\|_a$  yields the desired inequality.  $\square$

Céa's Lemma tells us that  $u_h$  is the best approximation to  $u$  in the energy norm. It also allows for an error bound in the  $H^1$ -norm:

**Corollary 1.** *Let  $u \in V$  and  $v_h \in V_h$  satisfy (4.19) and (4.20) respectively. Then:*

$$\|u - u_h\|_{H^1} \leq \sqrt{\frac{\rho + \theta}{\min\{\rho, \theta\}}} \|u - v_h\|_{H^1}, \quad \forall v_h \in V_h. \quad (4.24)$$

*Proof.* Let  $v_h \in V_h$  be arbitrary. Using the coercivity and boundedness property of  $a$  (see the proof of Theorem 7), and Céa's Lemma, we have:

$$\begin{aligned} \min\{\rho, \theta\} \|u - u_h\|_{H^1}^2 &\leq a(u - u_h, u - u_h) \\ &= \|u - u_h\|_a^2 \\ &\leq \|u - v_h\|_a^2 \\ &= a(u - v_h, u - v_h) \\ &\leq (\rho + \theta) \|u - v_h\|_{H^1}^2. \end{aligned}$$

Rearrange terms and take the square root on both sides to obtain the desired inequality.  $\square$

Céa's Lemma and its corollary tell us that the finite element solution  $u_h$  is the 'best' in some sense. We can exploit this property by constructing a function in  $V_h$  that approximates  $u$  really well. By Céa's Lemma,  $u_h$  will be an even better approximation (in the energy norm). We will use Lagrangian linear interpolation polynomials.

**Lemma 6** (Linear Interpolation). *Let  $\mathcal{T}_h$  be a mesh on  $(0, 1)$ , that is:*

$$\mathcal{T}_h = \{T_i = (x_i, x_{i+1}) \mid i \in \{1, \dots, n\}\}, \quad \text{such that } \bigcup_{i=1}^n \overline{T_i} = (0, 1).$$

*Let  $w \in H^1(0, 1)$  such that for each  $T \in \mathcal{T}_h$  we have  $w \in C^2(T)$  and  $w''$  is bounded on  $T$ . Here we use the following shorthand:  $w_i := w|_{T_i}$ . Let  $\pi_h w$  be the piecewise linear interpolant of  $w$  such that  $\pi_h w(x_i) = w(x_i)$  for each grid point  $x_i$ . Then there exists some  $c > 0$  such that:*

$$\|w - \pi_h w\|_{H^1} \leq cKh,$$

where

$$h = \max_i |T_i|,$$

and

$$K = \max_i \{K_i\}, \quad \text{for } K_i = \sup_{x \in T_i} \{w_i''(x)^2\}. \quad (4.25)$$

*Proof.* Define the interpolation error  $e(x) = w(x) - \pi_h w(x)$ . Consider the element  $T_i = (x_i, x_{i+1})$ . Since  $e(x_i) = e(x_{i+1}) = 0$ , by Rolle's theorem there is a  $\xi \in T_i$  such that  $e'(\xi) = 0$ . Taking into account that  $(\pi_h w)''(x) = 0$  on  $T_i$ , we get:

$$e'(x) = e'(x) - e'(\xi) = \int_{\xi}^x e''(t) dt = \int_{\xi}^x w_i''(t) dt.$$

By the Cauchy-Schwarz inequality we find for  $x \in T_i$ :

$$|e'(x)| \leq \int_{\xi}^x |w_i''(t)| dt \leq \int_{T_i} |w_i''(t)| dt \leq \sqrt{|T_i|} \sqrt{\int_{T_i} w_i''(t)^2 dt}.$$

It follows that

$$\begin{aligned} \|e'\|_{L^2}^2 &= \int_0^1 e'(x)^2 dx \\ &= \sum_i \int_{T_i} e'(x)^2 dx \\ &\leq \sum_i \int_{T_i} |T_i| \int_{T_i} w_i''(t)^2 dt dx \\ &= \sum_i |T_i|^2 \int_{T_i} (w_i''(t))^2 dt \\ &\leq \max_i |T_i|^2 \sum_i \int_{T_i} w_i''(t)^2 dt \end{aligned}$$

To eliminate the sum in the right-hand side, note that

$$\sum_i \int_{T_i} w_i''(t)^2 dt \leq \sum_i \int_{T_i} K_i dt \leq \sum_i \int_{T_i} K dt = K.$$

And thus we get:

$$\|e'\|_{L^2}^2 \leq Kh^2.$$

The  $L^2$ -norm of  $e$  can be computed in a similar fashion. Note that for  $x \in T_i$ :

$$e(x) = e(x) - e(x_i) = \int_{x_i}^x e'(t) dt.$$

And thus by Cauchy-Schwarz:

$$|e(x)| \leq \int_{x_i}^x |e'(t)| dt \leq \int_{T_i} |e'(t)| dt \leq \sqrt{|T_i|} \sqrt{\int_{T_i} e'(t)^2 dt}.$$

It follows that:

$$\begin{aligned} \|e\|_{L^2}^2 &= \int_0^1 e(x)^2 dx \\ &= \sum_i \int_{T_i} e(x)^2 dx \\ &\leq \sum_i \int_{T_i} |T_i| \int_{T_i} e'(t)^2 dt dx \\ &\leq h^2 \int_0^1 e'(t)^2 dt \\ &= h^2 \|e'\|_{L^2}^2 \leq Kh^4. \end{aligned}$$

Having computed both the  $L^2$ -norm of  $e$  and  $e'$ , the  $H^1$  can now be determined:

$$\|e\|_{H^1}^2 = \|e\|_{L^2}^2 + \|e'\|_{L^2}^2 \leq (h^4 + h^2)K.$$

Note that in our case  $h < 1$ , and thus  $h^4 + h^2 \leq 2h^2$ . Let  $c = \sqrt{2}$ , then we conclude that

$$\|e\|_{H^1} \leq cKh.$$

□

Note that our bilinear form  $a(\cdot, \cdot)$  is uniformly elliptic. Assume  $w$  satisfies

$$a(w, v) = (f, v), \quad \forall v \in V, \quad (4.26)$$

for some  $f \in L^2(0, 1)$ . Then by the uniform ellipticity of  $a$ , we have  $w \in H^2(0, 1)$  and

$$\|w\|_{H^2} \leq c\|f\|_{L^2}. \quad (4.27)$$

Since  $w \in H^2(0, 1)$  the proof of Lemma 6 can be slightly adapted to obtain

$$\|w - \pi_h w\|_{H^1} \leq c_1 h \|w''\|_{L^2}.$$

From which it follows that

$$\|w - \pi_h w\|_{H^1} \leq c_1 h \|w''\|_{L^2} \leq c_2 h \|w\|_{H^2} \leq c_3 h \|f\|_{L^2}. \quad (4.28)$$

Note that the constants  $c_1$ ,  $c_2$  and  $c_3$  do not depend on  $h$ . This leads us to the following theorem.

**Theorem 8** ( $L^2$ -error of FE-solution). *Let  $\mathcal{T}_h$  be a mesh on  $(0, 1)$  such that  $\alpha$  coincides with a grid node. Let  $V_h \subset V = H_0^1(0, 1)$  be the space of piecewise linear functions w.r.t.  $\mathcal{T}_h$ . Assume  $u \in V$  is the solution to (4.19), then  $u$  is also given by (4.15). Let  $u_h \in V_h$  satisfy the finite element problem (4.20), then:*

$$\|u - u_h\|_{L^2} \leq cKh^2 \quad (4.29)$$

for some  $c > 0$ , where  $h = \max_{T \in \mathcal{T}_h} |T|$  and  $K$  is given by (4.25) for  $w = u$ .

*Proof.* We use Nitsche's trick to prove the error bound. To this end, the dual problem is introduced. Assume  $\varphi \in V$  satisfies

$$a(\varphi, v) = (u - u_h, v), \quad \forall v \in V.$$

Where  $a$  is the same bilinear form as in (4.19) and  $(\cdot, \cdot)$  denotes the  $L^2$  inner product. Let  $\pi_h \varphi \in V_h$  be the linear interpolant of  $\varphi$  such that  $\pi_h \varphi(x_i) = \varphi(x_i)$  for each grid node  $x_i$ . Then, by the definition of the dual problem and the  $a$ -orthogonality of  $u - u_h$  to  $V_h$ , we have

$$\begin{aligned} \|u - u_h\|_{L^2}^2 &= (u - u_h, u - u_h) \\ &= a(\varphi, u - u_h) \\ &= a(\varphi - \pi_h \varphi, u - u_h) \\ &\leq \|\varphi - \pi_h \varphi\|_a \|u - u_h\|_a. \end{aligned}$$

In the last step, we have applied Cauchy-Schwarz. In the proof of Corollary 1 it is shown that the energy norm is bounded by a constant times the  $H^1$ -norm. Then by (4.28) we get:

$$\|\varphi - \pi_h \varphi\|_a \leq c_1 \|\varphi - \pi_h \varphi\|_{H^1} \leq c_2 h \|\varphi''\|_{L^2} \leq c_3 h \|u - u_h\|_{L^2}.$$

The reason that  $\|u - u_h\|_{L^2}$  appears is because it is the right-hand side of the dual problem solved by  $\varphi$ . By Céa's Lemma 5 we know that  $\|u - u_h\|_a \leq \|u - \pi_h u\|_a$ , where  $\pi_h u$  is the linear interpolant of  $u$  w.r.t.  $\mathcal{T}_h$ . Because we let  $\alpha$  coincide with a grid node of  $\mathcal{T}_h$ ,  $u$  meets the conditions of Lemma 6. Thus we get the following inequality:

$$\|u - u_h\|_a \leq \|u - \pi_h u\|_a \leq c_4 \|u - \pi_h u\|_{H^1} \leq c_5 Kh.$$

To summarize, we have found

$$\|u - u_h\|_{L^2}^2 \leq cKh^2 \|u - u_h\|_{L^2}.$$

Where  $c = c_3 c_5$ . Divide both sides by  $\|u - u_h\|_{L^2}$  to obtain the desired error bound. □



*Proof.* Let  $T_i = (x_i, x_{i+1})$  and define the error function  $e_i(x) = |u(x) - u_h(x)|$  on  $T_i$ . We claim that  $e_i \in C^2(T_i)$  with second derivative  $e_i''(x) = |u_i''(x)|$ , where  $u_i = u|_{T_i}$ . Indeed, note that  $u$  satisfies the differential equation

$$\rho u - \theta u'' = 0$$

on  $T_i$ . Thus we find

$$u'' = \frac{\rho}{\theta} u.$$

From expression (4.15) it can be deduced that  $u$  is non-negative and  $C^2$  on each open interval not containing  $\alpha$ . Therefore,  $u_i''$  is non-negative and  $C^2$  on  $T_i$ . Because  $u_h''(x) = 0$ , we get  $e_i''(x) = |u_i''(x)|$ . It follows that  $e_i$  satisfies the conditions of the trapezoid rule for approximating integrals. See for example [18, Theorem 5.3.3]. We get:

$$\begin{aligned} e_i(x_i) &\leq e_i(x_i) + e_i(x_{i+1}) \\ &= \frac{2}{h} \left| \frac{h}{2} e_i(x_i) + e_i(x_{i+1}) \right| \\ &\leq \frac{2}{h} \left( \int_{T_i} e_i \, dx + \gamma_i h^3 \right) \\ &\leq \frac{2}{h} \left( \sqrt{h} \sqrt{\int_{T_i} e_i^2 \, dx} + \gamma_i h^3 \right) \\ &\leq \frac{2}{h} \left( \sqrt{h} \|e\|_{L^2(0,1)} + \gamma_i h^3 \right) \\ &\leq \frac{2}{h} \left( c_1 h^2 \sqrt{h} + \gamma_i h^3 \right) \\ &= c_2 h^{3/2} + c_3 h^2. \end{aligned}$$

Here we have used the Cauchy-Schwarz inequality and Theorem 8.  $\square$

Although Lemma 7 shows that the pointwise error has order  $\mathcal{O}(h^{3/2})$ , we numerically verify that the pointwise error is actually of order  $\mathcal{O}(h^2)$ . See for example Figure 4.2, where we have plotted

$$\log_2 \left( \frac{E_{2n}}{E_n} \right), \quad (4.38)$$

where  $E_n$  is the maximum pointwise error between the exact solution and finite element solution on a uniform grid containing  $n$  elements.

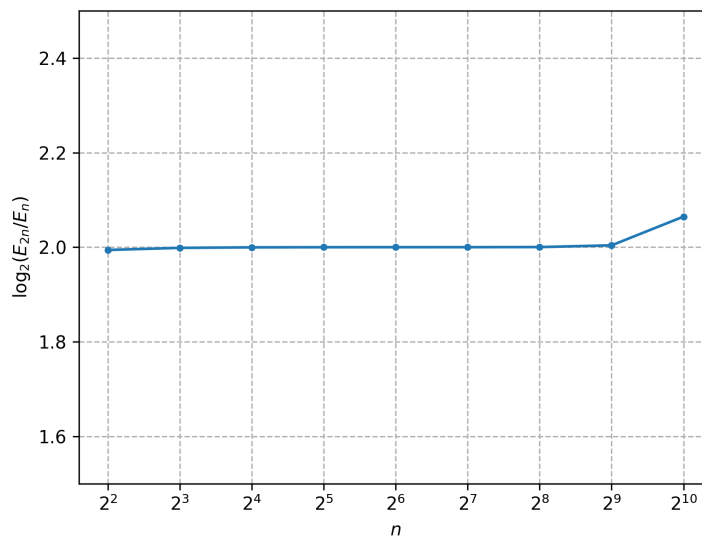


Figure 4.2: Numerical calculation of the order of the pointwise finite element error. We have set  $\rho = \theta = 1$  and  $i_\alpha = n/2$ , corresponding to  $\alpha = 0.5$ .

We have not been able to prove this pointwise error bound, and therefore continue this section with the result of Lemma 7 in mind.



**Theorem 9.** For each  $i, j \in \{0, \dots, n-1\}$ :

$$\tilde{A}_{ij}^{-1} - A_{ij}^{-1} = \mathcal{O}(h^{3/2}).$$

*Proof.* Assume  $A\mathbf{u}_i = \mathbf{e}_i$  and  $\tilde{A}\tilde{\mathbf{u}}_i = \mathbf{e}_i$ , where  $\tilde{A} = (\tilde{A}^{-1})^{-1}$ . By the definition and symmetry of  $\tilde{A}^{-1}$ , we have

$$(\tilde{\mathbf{u}}_i)_j = \left(\tilde{A}^{-1}\mathbf{e}_i\right)_j = \tilde{A}_{ji}^{-1} = \tilde{A}_{ij}^{-1} = u_i(x_j).$$

It then follows that

$$\tilde{A}_{ij}^{-1} - A_{ij}^{-1} = \tilde{A}_{ji}^{-1} - A_{ji}^{-1} = \left(\tilde{A}^{-1}\mathbf{e}_i\right)_j - \left(A^{-1}\mathbf{e}_i\right)_j = u_i(x_j) - (\mathbf{u}_i)_j = \mathcal{O}(h^{3/2})$$

In the last step the result from Lemma 7 is applied.  $\square$

We will see later that it is sufficient to have an element-wise approximation of  $A^{-1}$ .

## 4.4 Determining the Optimal Tuning Parameter

Recall that our aim is to find an approximation of the matrix  $C = DA^{-1}D^\top$ . Using the results from the previous section, we can now approximate  $C$ . We first express its entries in terms of those of  $A^{-1}$ .

**Lemma 8.** Let  $D$  be the finite element divergence matrix in the setting of Section 4.2. Then it is given by:

$$D = \frac{1}{2} \begin{pmatrix} -1 & 1 & & & \\ -1 & 0 & 1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 0 & 1 \\ & & & -1 & 0 \end{pmatrix}, \quad (4.39)$$

where  $D$  is an  $n \times n$  matrix, and we again use the index set  $\{0, \dots, n-1\}$ . Define the approximation  $\tilde{C} = D\tilde{A}^{-1}D^\top$  to  $C$ . Then for  $i, j \in \{1, \dots, n-2\}$ :

$$\tilde{C}_{ij} = \frac{1}{4} \left[ \tilde{A}_{i+1,j+1}^{-1} - \tilde{A}_{i-1,j+1}^{-1} - \tilde{A}_{i+1,j-1}^{-1} + \tilde{A}_{i-1,j-1}^{-1} \right]. \quad (4.40)$$

*Proof.* The proof follows from combining the following two identities that hold for any  $n \times n$  matrix  $B$ :

$$(DB)_{ij} = \frac{1}{2} \begin{cases} B_{1,j} - B_{0,j}, & i = 0 \\ B_{i+1,j} - B_{i-1,j}, & 1 \leq i \leq n-2, \\ -B_{n-2,j}, & i = n-1 \end{cases} \quad (4.41)$$

and

$$(BD^\top)_{ij} = \frac{1}{2} \begin{cases} B_{i,1} - B_{i,0}, & j = 0 \\ B_{i,j+1} - B_{i,j-1}, & 1 \leq j \leq n-2, \\ -B_{i,n-2}, & j = n-1 \end{cases} \quad (4.42)$$

Thus, for  $i, j \in \{1, \dots, n-2\}$  we find

$$\begin{aligned} \left(D\tilde{A}^{-1}D^\top\right)_{ij} &= \frac{1}{2} \left[ \left(D\tilde{A}^{-1}\right)_{i,j+1} - \left(D\tilde{A}^{-1}\right)_{i,j-1} \right] \\ &= \frac{1}{2} \left[ \frac{1}{2} \left[ \tilde{A}_{i+1,j+1}^{-1} - \tilde{A}_{i-1,j+1}^{-1} \right] - \frac{1}{2} \left[ \tilde{A}_{i+1,j-1}^{-1} - \tilde{A}_{i-1,j-1}^{-1} \right] \right] \\ &= \frac{1}{4} \left[ \tilde{A}_{i+1,j+1}^{-1} - \tilde{A}_{i-1,j+1}^{-1} - \tilde{A}_{i+1,j-1}^{-1} + \tilde{A}_{i-1,j-1}^{-1} \right] \end{aligned}$$

$\square$

Using Lemma 8 and equation (4.37), the entries of  $\tilde{C}$  can be computed exactly.

**Theorem 10.** Let  $\tilde{C}$  again be defined by  $\tilde{C} = D\tilde{A}^{-1}D^\top$ , and let  $i \in \{1, \dots, n-2\}$ , then:

$$\tilde{C}_{ii} = \frac{\lambda}{4\rho \cosh(\lambda)} \left[ 2h\lambda \cosh(\lambda(1-h)) + 2h^2\lambda^2 \sinh(\lambda(1-2x_i)) \right] + \mathcal{O}(h^3). \quad (4.43)$$

$$\tilde{C}_{i,i-1} = \frac{\lambda}{4\rho \cosh(\lambda)} \left[ h\lambda \cosh(\lambda(1-2h)) + 2h^2\lambda^2 \sinh(\lambda(1-2x_i+h)) \right] + \mathcal{O}(h^3). \quad (4.44)$$

$$\tilde{C}_{i,i+1} = \frac{\lambda}{4\rho \cosh(\lambda)} \left[ h\lambda \cosh(\lambda(1-2h)) + 2h^2\lambda^2 \sinh(\lambda(1-2x_i-h)) \right] + \mathcal{O}(h^3) \quad (4.45)$$

Futhermore, for  $1 \leq j \leq i-2$  we have:

$$\tilde{C}_{ij} = \frac{h^2\lambda^3}{2\rho \cosh(\lambda)} \left[ \sinh(\lambda(1-2x_i+\delta h)) - \sinh(\lambda(1-\delta h)) \right] + \mathcal{O}(h^4), \quad \delta = i-j. \quad (4.46)$$

Likewise, for  $i+2 \leq j \leq n-2$  we have:

$$\tilde{C}_{ij} = \frac{h^2\lambda^3}{2\rho \cosh(\lambda)} \left[ \sinh(\lambda(1-2x_i-\delta h)) - \sinh(\lambda(1-\delta h)) \right] + \mathcal{O}(h^4), \quad \delta = j-i. \quad (4.47)$$

*Proof.* To avoid unnecessary repetition, we will only derive the expression for  $\tilde{C}_{ii}$ . The other entries can be proved using a similar trick. Using (4.37) and Lemma 8, we find

$$\begin{aligned} \tilde{C}_{ii} = \frac{\lambda}{8\rho \cosh(\lambda)} \left[ 2(\sinh(\lambda) - \sinh(\lambda(1-2h))) \right. \\ \left. + \sinh(\lambda(1-2x_{i-1})) - 2\sinh(\lambda(1-2x_i)) + \sinh(\lambda(1-2x_{i+1})) \right]. \end{aligned}$$

The first two terms between brackets look like a finite difference approximation to a first derivative. Similarly, the last three terms look like a finite difference approximation of a second derivative. Define

$$\psi(x) = \sinh(\lambda(1-x)),$$

then

$$\psi'(x) = -\lambda \cosh(\lambda(1-x)), \quad \text{and} \quad \psi''(x) = \lambda^2 \sinh(\lambda(1-x)).$$

We use the following finite difference approximations:

$$\psi'(x) = \frac{\psi(x+h) - \psi(x-h)}{2h} + \mathcal{O}(h^2), \quad \text{and} \quad \psi''(x) = \frac{\psi(x+h) - 2\psi(x) + \psi(x-h)}{h^2} + \mathcal{O}(h^2).$$

From which it follows that

$$\begin{aligned} \sinh(\lambda) - \sinh(\lambda(1-2h)) &= \psi(0) - \psi(2h) \\ &= -2h\psi'(h) + \mathcal{O}(h^3) \\ &= 2h\lambda \cosh(\lambda(1-h)) + \mathcal{O}(h^3). \end{aligned}$$

Moreover:

$$\begin{aligned} \sinh(\lambda(1-2x_{i-1})) - 2\sinh(\lambda(1-2x_i)) + \sinh(\lambda(1-2x_{i+1})) &= \psi(2x_{i-1}) - 2\psi(2x_i) + \psi(2x_{i+1}) \\ &= \psi(2x_i - 2h) - 2\psi(2x_i) + \psi(2x_i + 2h) \\ &= 4h^2\psi''(x) + \mathcal{O}(h^4) \\ &= 4h^2\lambda^2 \sinh(\lambda(1-2x_i)) + \mathcal{O}(h^4). \end{aligned}$$

Thus we find

$$\tilde{C}_{ii} = \frac{\lambda}{4\rho \cosh(\lambda)} \left[ 2h\lambda \cosh(\lambda(1-h)) + 2h^2\lambda^2 \sinh(\lambda(1-2x_i)) \right] + \mathcal{O}(h^3)$$

Which completes the proof for  $\tilde{C}_{ii}$ . To prove the identities for  $\tilde{C}_{i,i-1}$  and  $\tilde{C}_{i,i+1}$  one must also recognize finite difference quotients for a first and second derivative of  $\psi$ . To compute the remaining identities (4.46) and (4.47), two second derivative quotients should be used.  $\square$

Theorem 10 only considers the interior of  $\tilde{C}$  and not the first and last rows or columns. Note that we can expand Lemma 8 to also include  $i, j \in \{0, n-1\}$  using the expressions (4.41) and (4.42). Because this leads to a large number of distinct cases, we have omitted these derivations from this report. It has been verified that expression such as the ones in Theorem 10 also hold for the entries in the first and last rows or columns of  $\tilde{C}$ . Numerical experiments show that the entries become more accurate as they are further away from the diagonal.

From Theorem 9 we know that the entries of  $\tilde{A}^{-1}$  are accurate up to order  $\mathcal{O}(h^{3/2})$ . Thus, the order  $\mathcal{O}(h^2)$  and higher terms in the entries of  $\tilde{C}$  will not contribute to a more accurate result. Let  $Q$  be the order  $h$  part of  $\tilde{C}$ . By dividing the expressions in Theorem 10 by  $h$  and taking the limit  $h \rightarrow 0$ , we obtain for  $i \in \{1, \dots, n-2\}$ :

$$(Q\mathbf{p})_i = \frac{h\lambda^2}{4\rho} (p_{i-1} + 2p_i + p_{i+1}) = \frac{h}{4\theta} (p_{i-1} + 2p_i + p_{i+1}). \quad (4.48)$$

And, as a direct result of Theorem 9:

$$(Q\mathbf{p})_i = (DA^{-1}D^\top \mathbf{p})_i + \mathcal{O}(h^{3/2}). \quad (4.49)$$

Consider again the coefficient matrix in the left-hand side of (4.11). Let  $i \in \{1, \dots, n-2\}$ , then:

$$\begin{aligned} [(kL_p + \Delta t C)\mathbf{p}]_i &= [(kL_p + \Delta t Q)\mathbf{p}]_i + \mathcal{O}(h^2) \\ &= \left(\frac{h\Delta t}{4\theta} - \frac{k}{h}\right) p_{i-1} + \left(\frac{h\Delta t}{2\theta} + \frac{2k}{h}\right) p_i + \left(\frac{h\Delta t}{4\theta} - \frac{k}{h}\right) p_{i+1} + \mathcal{O}(h^{3/2}) \end{aligned} \quad (4.50)$$

Stabilization is achieved by replacing  $kL_p$  with  $(k + \beta)L_p$ , we then get:

$$\begin{aligned} [((k + \beta)L_p + \Delta t C)\mathbf{p}]_i &= \\ &= \left(\frac{h\Delta t}{4\theta} - \frac{k + \beta}{h}\right) p_{i-1} + \left(\frac{h\Delta t}{2\theta} + \frac{2(k + \beta)}{h}\right) p_i + \left(\frac{h\Delta t}{4\theta} - \frac{k + \beta}{h}\right) p_{i+1} + \mathcal{O}(h^{3/2}) \end{aligned} \quad (4.51)$$

Let

$$\beta^* = \frac{h^2 \Delta t}{4\theta} = \frac{h^2}{4(\mu_v + \lambda_v + \Delta t(\mu_\varepsilon + \lambda_\varepsilon))}. \quad (4.52)$$

Then  $\beta^*$  is the smallest value of  $\beta$  such that  $(k + \beta)L_p + \Delta t Q$  is an M-matrix:

$$[((k + \beta^*)L_p + \Delta t C)\mathbf{p}]_i = -\frac{k}{h} p_{i-1} + \left(\frac{h\Delta t}{\theta} + \frac{2k}{h}\right) p_i - \frac{k}{h} p_{i+1} + \mathcal{O}(h^{3/2}) \quad (4.53)$$

Note that from (4.48) it follows that  $\rho$  does not affect the order  $h$  part of  $\tilde{C}$ . Numerical experiments have confirmed that increasing or decreasing  $\rho$  has little to no effect on the occurrence of non-physical oscillations.

If we assume that the two-sided stabilization of (4.3) is used, and  $\beta$  is set to  $\beta^*$ , then the perturbation in the corresponding differential equation (4.6) is - at least - of order  $\mathcal{O}(\Delta t^2 h^{3/2})$ . If the phenomenon seen in Figure 4.2 is taken into account, the order of the perturbation can be further sharpened to  $\mathcal{O}(\Delta t^2 h^2)$ .

## 4.5 Morpho-Poroelasticity

Combining the morpho- and poroelastic systems from Section 2.4 and 3.4 respectively yields the following system:

$$\rho \left( \frac{D\mathbf{v}}{Dt} + (\nabla \cdot \mathbf{v})\mathbf{v} \right) - \nabla \cdot \boldsymbol{\sigma}(\mathbf{v}, \boldsymbol{\varepsilon}, p) = \mathbf{g}, \quad (4.54)$$

$$\frac{D\boldsymbol{\varepsilon}}{Dt} + \boldsymbol{\varepsilon} \text{skw}(\nabla \mathbf{v}) - \text{skw}(\nabla \mathbf{v})\boldsymbol{\varepsilon} + (\text{tr}(\boldsymbol{\varepsilon}) - 1) \text{sym}(\nabla \mathbf{v}) = -\mathbf{G}, \quad (4.55)$$

$$\nabla \cdot \mathbf{v} - \nabla \cdot (k\nabla p) = f. \quad (4.56)$$

Where  $\mathbf{G}$  is the growth tensor. The above equations must be satisfied on the moving domain  $\Omega(t)$ . The stress tensor  $\boldsymbol{\sigma}$  in this framework is given by:

$$\boldsymbol{\sigma}(\mathbf{v}, \boldsymbol{\varepsilon}, p) = \boldsymbol{\sigma}_{\text{el}}(\boldsymbol{\varepsilon}) + \boldsymbol{\sigma}_{\text{vis}}(\mathbf{v}) - p\mathbf{I}, \quad (4.57)$$

where

$$\begin{aligned}\boldsymbol{\sigma}_{\text{el}}(\boldsymbol{\varepsilon}) &= \mu_{\varepsilon}\boldsymbol{\varepsilon} + \lambda_{\varepsilon}\text{tr}(\boldsymbol{\varepsilon})\mathbf{I}, \\ \boldsymbol{\sigma}_{\text{vis}}(\mathbf{v}) &= \frac{1}{2}\mu_v(\nabla\mathbf{v} + \nabla\mathbf{v}^{\top}) + \lambda_v(\nabla\cdot\mathbf{v})\mathbf{I}.\end{aligned}$$

The same boundary conditions as for the visco- poroelastic system are imposed:

$$\mathbf{v} = \mathbf{0}, \quad k\nabla p \cdot \mathbf{n} = 0, \quad \text{on } \Gamma_1(t), \quad (4.58)$$

$$\boldsymbol{\sigma} \cdot \mathbf{n} = \boldsymbol{\tau}, \quad p = 0, \quad \text{on } \Gamma_2(t). \quad (4.59)$$

We assume that  $\Gamma(t)$  is the piecewise smooth boundary of  $\Omega(t)$ , such that  $\Gamma(t) = \Gamma_1(t) \cup \Gamma_2(t)$  and  $\Gamma_1(t) \cap \Gamma_2(t) = \emptyset$  for all  $t > 0$ .

### 4.5.1 Relation to Morphoelastic System

The morpho-poroelastic system is closely related to the morphoelastic system. We show that as  $k$  approaches infinity the solutions for  $\mathbf{v}$  and  $\boldsymbol{\varepsilon}$  of the morpho-poroelastic system approach the corresponding solutions to the morphoelastic system. Intuitively, increasing  $k$  means either the permeability of the porous medium increases, or the dynamic viscosity of the poro-fluid decreases. In both cases the fluid is able to move more freely throughout the medium as it deforms, having less impact on the domain deformation. This is formalised in the following theorem.

**Theorem 11.** *Let  $\mathbf{v}_k$ ,  $\boldsymbol{\varepsilon}_k$  and  $p_k$  be the solutions to the morpho-poroelastic system (4.54)-(4.56) with boundary conditions*

$$\mathbf{v} = \mathbf{0}, \quad k\nabla p \cdot \mathbf{n} = J, \quad \text{on } \Gamma_1(t), \quad (4.60)$$

$$\boldsymbol{\sigma} \cdot \mathbf{n} = \boldsymbol{\tau}, \quad p = p_0, \quad \text{on } \Gamma_2(t), \quad (4.61)$$

where  $J$  is bounded and  $p_0$  is constant. Let  $\mathbf{v}$  and  $\boldsymbol{\varepsilon}$  be the solutions to the morphoelastic system (4.58)-(4.59) with boundary conditions

$$\mathbf{v} = \mathbf{0}, \quad \text{on } \Gamma_1(t), \quad (4.62)$$

$$\boldsymbol{\sigma} \cdot \mathbf{n} = \boldsymbol{\tau}, \quad \text{on } \Gamma_2(t). \quad (4.63)$$

Then  $\mathbf{v}_k \rightarrow \mathbf{v}$  and  $\boldsymbol{\varepsilon}_k \rightarrow \boldsymbol{\varepsilon}$  as  $k \rightarrow \infty$ .

*Proof.* We give the outline of the proof. Some properties of the solutions are assumed without them being proved. In a future work, these technicalities should be addressed. Note that in the limit  $k \rightarrow \infty$ , equation (4.56) turns into

$$\nabla^2 p = 0,$$

with boundary conditions

$$\nabla p \cdot \mathbf{n} = 0, \quad \text{on } \Gamma_1(t), \quad p = p_0, \quad \text{on } \Gamma_2(t).$$

This holds provided  $\nabla \cdot \mathbf{v}$  remains bounded on  $\Omega(t)$  for all  $t$ . Note that since  $\mathbf{v}$  has been eliminated from (4.56), the equation has been decoupled from (4.54) and (4.55) but not vice versa, since  $p$  still occurs in (4.54). The resulting boundary value problem for  $p$  has the unique constant solution  $p(\mathbf{x}, t) = p_0$  on  $\Omega(t)$ . If we substitute this solution into the stress tensor (4.57) and take the divergence, we obtain:

$$\begin{aligned}\nabla \cdot \boldsymbol{\sigma}(\mathbf{v}, \boldsymbol{\varepsilon}, p_0) &= \nabla \cdot (\boldsymbol{\sigma}_{\text{el}}(\boldsymbol{\varepsilon}) + \boldsymbol{\sigma}_{\text{vis}}(\mathbf{v})) - \nabla p_0 \\ &= \nabla \cdot (\boldsymbol{\sigma}_{\text{el}}(\boldsymbol{\varepsilon}) + \boldsymbol{\sigma}_{\text{vis}}(\mathbf{v})),\end{aligned}$$

because  $\nabla p_0$  vanishes. Thus, we find that  $\mathbf{v}$  and  $\boldsymbol{\varepsilon}$  satisfy the morphoelastic system (4.58)-(4.59).  $\square$

### 4.5.2 Derivation of System of Equations

We present the weak forms corresponding to equations (4.54)-(4.56). Their complete derivations will not be discussed, as most of them have already been shown a number of times throughout this work. The weak form of (4.54) is similar to (3.91):

$$\begin{aligned}\rho \frac{d}{dt} \int_{\Omega(t)} \mathbf{v} \cdot \boldsymbol{\varphi} \, d\Omega + \int_{\Omega(t)} (\boldsymbol{\sigma}_{\text{el}}(\boldsymbol{\varepsilon}) + \boldsymbol{\sigma}_{\text{vis}}(\mathbf{v})) : \nabla \boldsymbol{\varphi} \, d\Omega - \int_{\Omega(t)} p(\nabla \cdot \boldsymbol{\varphi}) \, d\Omega \\ = \int_{\Omega(t)} \mathbf{g} \cdot \boldsymbol{\varphi} \, d\Omega + \int_{\Gamma_2(t)} \boldsymbol{\tau} \cdot \boldsymbol{\varphi} \, d\Gamma.\end{aligned} \quad (4.64)$$

Where  $\boldsymbol{\varphi}$  is a vector-valued Lagrangian test field that vanishes on  $\Gamma_1(t)$ . We show the derivation of the term involving the stress tensor. Using Green's formula, Gauss' divergence theorem and the boundary conditions, we get:

$$\begin{aligned} - \int_{\Omega(t)} (\nabla \cdot \boldsymbol{\sigma}) \cdot \boldsymbol{\varphi} \, d\Omega &= \int_{\Omega(t)} \boldsymbol{\sigma} : \nabla \boldsymbol{\varphi} - \nabla \cdot (\boldsymbol{\sigma} \cdot \boldsymbol{\varphi}) \, d\Omega \\ &= \int_{\Omega(t)} (\boldsymbol{\sigma}_{\text{el}}(\boldsymbol{\varepsilon}) + \boldsymbol{\sigma}_{\text{vis}}(\mathbf{v}) - p\mathbf{I}) : \nabla \boldsymbol{\varphi} \, d\Omega - \int_{\Gamma(t)} (\mathbf{n} \cdot \boldsymbol{\sigma}) \cdot \boldsymbol{\varphi} \, d\Gamma \\ &= \int_{\Omega(t)} (\boldsymbol{\sigma}_{\text{el}}(\boldsymbol{\varepsilon}) + \boldsymbol{\sigma}_{\text{vis}}(\mathbf{v})) : \nabla \boldsymbol{\varphi} \, d\Omega - \int_{\Omega(t)} p(\nabla \cdot \boldsymbol{\varphi}) \, d\Omega - \int_{\Gamma_2(t)} \boldsymbol{\tau} \cdot \boldsymbol{\varphi} \, d\Gamma. \end{aligned}$$

The weak form of (4.55) is derived in (2.106), it is given by:

$$\begin{aligned} \frac{d}{dt} \int_{\Omega(t)} \boldsymbol{\varepsilon} : \boldsymbol{\zeta} \, d\Omega \\ + \int_{\Omega(t)} \left( \boldsymbol{\varepsilon} \text{skw}(\nabla \mathbf{v}) - \text{skw}(\nabla \mathbf{v}) \boldsymbol{\varepsilon} + (\text{tr}(\boldsymbol{\varepsilon}) - 1) \text{sym}(\nabla \mathbf{v}) - (\nabla \cdot \mathbf{v}) \boldsymbol{\varepsilon} \right) : \boldsymbol{\zeta} \, d\Omega \\ = - \int_{\Omega(t)} \mathbf{G} : \boldsymbol{\zeta} \, d\Omega. \end{aligned} \quad (4.65)$$

Here  $\boldsymbol{\zeta}$  is a tensor-valued Lagrangian test field. Finally, the weak form of equation (4.56) is given by

$$\int_{\Omega(t)} (\nabla \cdot \mathbf{v}) \psi \, d\Omega + \int_{\Omega(t)} k \nabla p \cdot \nabla \psi \, d\Omega = \int_{\Omega(t)} f \psi \, d\Omega \quad (4.66)$$

Where  $\psi$  is a scalar-valued Lagrangian test function. Weak forms (4.64), (4.65) and (4.66) are converted to a semi-discretised algebraic system by computing the Galerkin equations. To this end, the approximate domain  $\Omega_h(t)$  is triangulated by  $\mathcal{T}_h(t)$ . The vectors  $\mathbf{v}$ ,  $\boldsymbol{\varepsilon}$  and  $\mathbf{p}$  contain the values of their corresponding unknown in the grid points. By following the usual steps we arrive at the following set of ordinary differential equations:

$$\rho \frac{d}{dt} (M_v \mathbf{v}) + S_{\text{vis}} \mathbf{v} + S_{\text{el}} \boldsymbol{\varepsilon} - D^\top \mathbf{p} = \mathbf{g} \quad (4.67)$$

$$\frac{d}{dt} (M_\varepsilon \boldsymbol{\varepsilon}) - B \mathbf{v} + \mathbf{N}(\mathbf{v}, \boldsymbol{\varepsilon}) = -\mathbf{G} \quad (4.68)$$

$$D \mathbf{v} + k L \mathbf{p} = \mathbf{f} \quad (4.69)$$

Here  $\mathbf{v} = (\mathbf{v}^1, \mathbf{v}^2)^\top$ ,  $\boldsymbol{\varepsilon} = (\boldsymbol{\varepsilon}^{11}, \boldsymbol{\varepsilon}^{12}, \boldsymbol{\varepsilon}^{22})^\top$  and  $\mathbf{p}$  are the vectors containing the unknowns, and  $\mathbf{g} = (\mathbf{g}^1, \mathbf{g}^2)^\top$  and  $\mathbf{G} = (\mathbf{G}^{11}, \mathbf{G}^{12}, \mathbf{G}^{22})^\top$  and  $\mathbf{f}$  are the right-hand side vectors. The symmetry of  $\boldsymbol{\varepsilon}$  has been used to eliminate  $\boldsymbol{\varepsilon}^{21}$  from the system. Moreover, the two mass matrices are given by:

$$M_v = \begin{pmatrix} M & 0 \\ 0 & M \end{pmatrix}, \quad M_\varepsilon = \begin{pmatrix} M & 0 & 0 \\ 0 & M & 0 \\ 0 & 0 & M \end{pmatrix}, \quad (4.70)$$

where  $M$  is the mass matrix given in (3.96). Note that here we assume  $\mathbf{v}$  and  $\boldsymbol{\varepsilon}$  both use linear basis functions. Since the divergence matrix  $D$  is the same as in (3.62), it again consists of the two blocks  $D = (D^1 \quad D^2)$ . From (2.109) we see that the matrices  $B$  and  $S_{\text{el}}$  can be expressed in terms of  $D^1$  and  $D^2$ :

$$B = \begin{pmatrix} D^1 & 0 \\ \frac{1}{2} D^2 & \frac{1}{2} D^1 \\ 0 & D^2 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & \frac{1}{2} \\ 0 & 0 \end{pmatrix} \otimes D^1 + \begin{pmatrix} 0 & 0 \\ \frac{1}{2} & 0 \\ 0 & 1 \end{pmatrix} \otimes D^2, \quad (4.71)$$

and

$$S_{\text{el}}^\top = \begin{pmatrix} (\mu_\varepsilon + \lambda_\varepsilon) D^1 & \lambda_\varepsilon D^2 \\ \mu_\varepsilon D^2 & \mu_\varepsilon D^1 \\ \lambda_\varepsilon D^1 & (\mu_\varepsilon + \lambda_\varepsilon) D^2 \end{pmatrix} = \begin{pmatrix} \mu_\varepsilon + \lambda_\varepsilon & 0 \\ 0 & \mu_\varepsilon \\ \lambda_\varepsilon & 0 \end{pmatrix} \otimes D^1 + \begin{pmatrix} 0 & \lambda_\varepsilon \\ \mu_\varepsilon & 0 \\ 0 & \mu_\varepsilon + \lambda_\varepsilon \end{pmatrix} \otimes D^2. \quad (4.72)$$

Here  $\otimes$  denotes the Kronecker product. Note that the elasticity matrix  $S_{\text{el}}$  now acts on  $\boldsymbol{\varepsilon}$  instead of the displacement  $\mathbf{u}$ . The matrix  $S_{\text{vis}}$  is well-known viscous stiffness matrix of which its elements are given by

(2.71)-(2.74). The non-linear operator  $\mathbf{N}$  is defined in Section 2.4.2. We apply implicit Euler to system (4.67)-(4.69). Let a superscript  $(\cdot)^m$  denote a time-level. Note that all discrete operators depend on the current grid, and therefore also get a time-level superscript. We obtain the fully discrete system:

$$(\rho M_v^m + \Delta t S_{\text{vis}}^m) \mathbf{v}^m + \Delta t S_{\text{el}}^m \boldsymbol{\varepsilon}^m - \Delta t (D^m)^\top \mathbf{p}^m = \rho M_v^{m-1} \mathbf{v}^{m-1} + \Delta t \mathbf{g}^m, \quad (4.73)$$

$$M_\varepsilon^m \boldsymbol{\varepsilon}^m - \Delta t B^m \mathbf{v}^m + \Delta t \mathbf{N}^m(\mathbf{v}^m, \boldsymbol{\varepsilon}^m) = M_\varepsilon^{m-1} \boldsymbol{\varepsilon}^{m-1} - \Delta t \mathbf{G}^m, \quad (4.74)$$

$$D^m \mathbf{v}^m + k L^m \mathbf{p}^m = \mathbf{f}^m. \quad (4.75)$$

Equivalently, in block-matrix notation:

$$\begin{pmatrix} \rho M_v^m + \Delta t S_{\text{vis}}^m & \Delta t S_{\text{el}}^m & -\Delta t (D^m)^\top \\ -\Delta t B^m & M_\varepsilon^m & 0 \\ D^m & 0 & k L^m \end{pmatrix} \begin{pmatrix} \mathbf{v}^m \\ \boldsymbol{\varepsilon}^m \\ \mathbf{p}^m \end{pmatrix} + \Delta t \begin{pmatrix} 0 \\ \mathbf{N}^m(\mathbf{v}^m, \boldsymbol{\varepsilon}^m) \\ 0 \end{pmatrix} = \begin{pmatrix} \rho M_v^{m-1} \mathbf{v}^{m-1} \\ M_\varepsilon^{m-1} \boldsymbol{\varepsilon}^{m-1} \\ 0 \end{pmatrix} + \begin{pmatrix} \Delta t \mathbf{g}^m \\ -\Delta t \mathbf{G}^m \\ \mathbf{f}^m \end{pmatrix}. \quad (4.76)$$

Note that to solve system (4.76), the operators and right-hand side vectors at time-level  $m$  must be known. However, these depend on the grid at time-level  $m$  and thus on the solution of the system. As such, an iterative method is used to solve (4.76). Details regarding this method can be found in section 2.4.3.

### 4.5.3 Stabilization

As was the case in the previous poroelastic systems, equation (4.56) will cause non-physical oscillations in the finite element solution for small values of  $k$ . To find out which matrix is causing this behaviour, write system (4.76) in the following abbreviated form:

$$\begin{pmatrix} A & \Delta t S_{\text{el}} & -\Delta t D^\top \\ -\Delta t B & M_\varepsilon & 0 \\ D & 0 & k L \end{pmatrix} \begin{pmatrix} \mathbf{v} \\ \boldsymbol{\varepsilon} \\ \mathbf{p} \end{pmatrix} + \Delta t \begin{pmatrix} 0 \\ \mathbf{N}(\mathbf{v}, \boldsymbol{\varepsilon}) \\ 0 \end{pmatrix} = \begin{pmatrix} \mathbf{r}^1 \\ \mathbf{r}^2 \\ \mathbf{r}^3 \end{pmatrix}. \quad (4.77)$$

Here we have dropped the time-level superscripts and set  $A = \rho M_v + \Delta t S_{\text{vis}}$ . Solve the first equation for  $\mathbf{v}$ :

$$\mathbf{v} = \Delta t A^{-1} D^\top \mathbf{p} - \Delta t A^{-1} S_{\text{el}} \boldsymbol{\varepsilon} + A^{-1} \mathbf{r}^1. \quad (4.78)$$

Substitute the above result in the third equation to obtain:

$$(k L + \Delta t D A^{-1} D^\top) \mathbf{p} = \Delta t D A^{-1} S_{\text{el}} \boldsymbol{\varepsilon} - D A^{-1} \mathbf{r}^1 + \mathbf{r}^3. \quad (4.79)$$

The left-hand side of the above equation has the same structure as (4.11), but with a different matrix  $A$ . By repeating the derivation of Sections 4.2, 4.3 and 4.4 for the morpho-poroelastic system, we obtain the following optimal tuning parameter:

$$\beta^* = \frac{h^2}{4(\mu_v + \lambda_v)}. \quad (4.80)$$

This result is then extrapolated to the two-dimensional morpho-poroelastic system on a non-uniform grid. The complete stabilized version of system (4.76) is given by:

$$\begin{pmatrix} \rho M_v^m + \Delta t S_{\text{vis}}^m & \Delta t S_{\text{el}}^m & -\Delta t (D^m)^\top \\ -\Delta t B^m & M_\varepsilon^m & 0 \\ D^m & 0 & k L^m + C_s^m \end{pmatrix} \begin{pmatrix} \mathbf{v}^m \\ \boldsymbol{\varepsilon}^m \\ \mathbf{p}^m \end{pmatrix} + \Delta t \begin{pmatrix} 0 \\ \mathbf{N}^m(\mathbf{v}^m, \boldsymbol{\varepsilon}^m) \\ 0 \end{pmatrix} = \begin{pmatrix} \rho M_v^{m-1} \mathbf{v}^{m-1} \\ M_\varepsilon^{m-1} \boldsymbol{\varepsilon}^{m-1} \\ C_s^{m-1} \mathbf{p}^{m-1} \end{pmatrix} + \begin{pmatrix} \Delta t \mathbf{g}^m \\ -\Delta t \mathbf{G}^m \\ \mathbf{f}^m \end{pmatrix}, \quad (4.81)$$

where the stabilization matrix  $C_s$  is given by:

$$(C_s)_{ij}^m = \frac{1}{4(\mu_v + \lambda_v)} \sum_{T \in \mathcal{T}^m} h_T^2 \int_T \nabla \varphi_i \cdot \nabla \varphi_j \, d\Omega. \quad (4.82)$$

Here  $\mathcal{T}^m$  denotes the triangulation of the domain at time-level  $m$ .

## 4.6 Comparison Between Morphoelastic and Morpho-Poroelastic Models

We compare the behaviour of the morpho-poroelastic model compared to the morphoelastic model from Section 2.4. To this end, a similar configuration as in Section 3.5 is used, where we compare the viscoporoelastic model to the viscoelastic model. Thus, the following body force is applied to the domain:

$$\mathbf{g}(\mathbf{x}, t) = \begin{pmatrix} 0.2 \sin(\pi t) \mathbb{1}_{[0,1]}(t) \\ 0 \end{pmatrix}. \quad (4.83)$$

The initial domain is again the unit square. The mass source function  $f$  is set to zero, and the following mechanical parameters are used:

$$\mu_\varepsilon = \lambda_\varepsilon = 1, \quad \mu_v = \lambda_v = 0.5, \quad \rho = 0.5. \quad (4.84)$$

Furthermore, we set the growth parameter  $\alpha$  equal to 1 for both models. We use the (initially) uniform finite element mesh from Figure 2.10 and a timestep of size  $\Delta t = 0.05$ . Figures 4.3a and 4.3b show the maximum domain width and minimum domain height as functions of time, for various morphoelastic models. The red lines correspond to the morpho-poroelastic simulations for various values of  $k$ , whereas the green line corresponds to the morphoelastic model from Section 2.4.

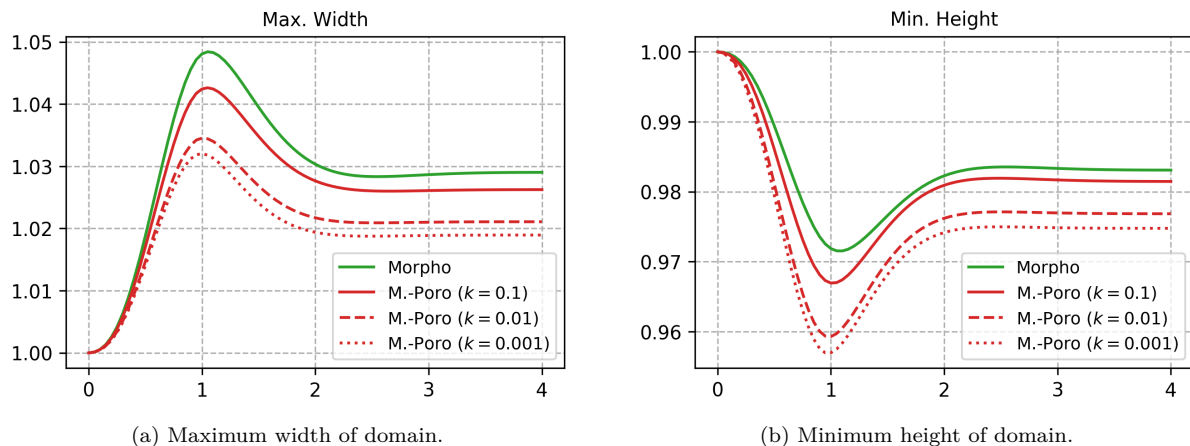


Figure 4.3: Maximum domain width and minimum domain height as functions of time for morphoelastic model and three morpho-poroelastic models.

Note that for smaller values of  $k$ , there is less extension in the  $x$ -direction (width), but more compression in the  $y$ -direction (height). It also seems that the viscoelastic ‘overshoot’ property is less pronounced for lower values of  $k$ . For higher values of  $k$ , we see that the solution to the morpho-poroelastic system approaches the solution to the morphoelastic system. This is in line with the theory from Section 4.5.1

# 5

## Tumor Growth

In this chapter we apply the developed morpho-poroelastic model to tumor growth. To this end, the framework of [14] is used. We will first give a brief summary of this framework.

### 5.1 Mathematical Model

The model proposed in [14] is based on the idea that a tumor grows better in areas where there is a high concentration of oxygen. The presence of oxygen will therefore encourage tumor growth. The new tissue that is created will exert a force on its surroundings, which has to be incorporated in the force balance equation. Let  $\eta$  be the volume of new tissue created per unit volume of tissue. To model the force exerted by the new tissue on the medium, the following term is added to the stress tensor:

$$-\left(\frac{\mu_\varepsilon}{3} + \lambda_\varepsilon\right)\eta\mathbf{I} \quad (5.1)$$

The parameters  $\lambda_\varepsilon$  and  $\mu_\varepsilon$  are the familiar Lamé constants, such that the constant  $\lambda_\varepsilon + \mu_\varepsilon/3$  equals the bulk modulus of the drained medium. Moreover, depending on whether the underlying model is purely elastic, viscoelastic or morphoelastic, the elastic stress tensor  $\boldsymbol{\sigma}_{\text{el}}$  can be adjusted. In the latter two cases, an inertial term must be included. The change in time of  $\eta$  depends on the oxygen concentration  $c$  and the fraction of the medium occupied by cells, which is equal to the solid part of the medium  $\Phi$ . Therefore,  $1 - \Phi$  is the fraction of space that is occupied by fluid. It is important to note that  $\Phi$  is not a constant; it can change in time and space according to the following equation [14, Eq. 1]:

$$\frac{\partial\Phi}{\partial t} + \nabla \cdot (\Phi\mathbf{v}) = S. \quad (5.2)$$

The function  $S$  is the cell creation source. It also acts as a source function for  $\eta$  in the sense that  $\partial\eta/\partial t = S$ . In [14] the following cell creation source function is proposed:

$$S = F(\bar{\sigma})\Phi(1 - \Phi)N_G(c) - r_d\Phi, \quad (5.3)$$

where  $N_G(c)$  is the Michaelis-Menten cell growth rate as a function of the nutrient concentration. It is given by:

$$N_G(c) = \frac{G_{\text{max}}c}{c_G + c}. \quad (5.4)$$

The parameter  $G_{\text{max}}$  is the maximum cell growth rate when there is no shortage of oxygen, and  $c_G$  is the Michaelis constant, which is equal to the oxygen concentration for which the cell creation rate is half of  $G_{\text{max}}$ . Looking back at (5.3), the parameter  $r_d$  is the death rate of tumor cells. The term  $r_d$  models the fact that tumor cells die in the absence of oxygen. The factor  $F(\bar{\sigma})$  is a stress-dependent growth factor. In [14],  $F$  is chosen equal to  $F(\bar{\sigma}) = 1 - \beta\bar{\sigma}$ . The parameter  $\beta$  is chosen to equal

$$\beta = \frac{3}{\sqrt[3]{2}} \left( \frac{1}{\mu_\varepsilon + 3\lambda_\varepsilon} + \frac{1}{2\mu_\varepsilon} \right). \quad (5.5)$$

This seemingly arbitrary choice is derived in [14] by calculating the stress needed for a spherical cell to increase its radius. If we repeat this derivation for a circular cell in two dimensions, we find:

$$\beta = \sqrt{2} \left( \frac{1}{\mu_\varepsilon + 3\lambda_\varepsilon} + \frac{1}{2\mu_\varepsilon} \right), \quad (5.6)$$



under the assumption that [14, Eq (12)] holds for circular cells. Note the factor  $\sqrt{2} = 2/\sqrt{2}$  instead of  $3/\sqrt[3]{2}$ , which is the result of using area (squares) instead of volume (cubes). Moreover,  $\bar{\sigma}$  is the average of the radial and circumferential stresses. In three dimensions, we would have:

$$\bar{\sigma} = \frac{\sigma_{rr} + \sigma_{\theta\theta} + \sigma_{\varphi\varphi}}{3}, \quad (5.7)$$

where  $\theta$  and  $\varphi$  are respectively the azimuthal and polar angle. Using the fact that the trace of a matrix is invariant under a change of basis, we also have

$$\bar{\sigma} = \frac{\sigma_{xx} + \sigma_{yy} + \sigma_{zz}}{3}. \quad (5.8)$$

Similarly, in two dimensions:

$$\bar{\sigma} = \frac{\sigma_{xx} + \sigma_{yy}}{2}. \quad (5.9)$$

Note that  $\bar{\sigma}$  depends on the elastic model that is used. For example:

- For the visco-poroelastic model from Section 3.4, we have

$$\boldsymbol{\sigma} = \boldsymbol{\sigma}_{\text{el}}(\mathbf{u}) + \boldsymbol{\sigma}_{\text{vis}}(\mathbf{v}) - p\mathbf{I} - \left(\frac{\mu_\varepsilon}{3} + \lambda_\varepsilon\right)\eta\mathbf{I}, \quad (5.10)$$

where the viscous stress tensor is  $\boldsymbol{\sigma}_{\text{vis}}(\mathbf{v}) = \mu_v \text{sym}(\nabla \mathbf{v}) + \lambda_v (\nabla \cdot \mathbf{v})\mathbf{I}$ . In this case

$$\bar{\sigma} = \left(\frac{\mu_\varepsilon}{2} + \lambda_\varepsilon\right) (\nabla \cdot \mathbf{u}) + \left(\frac{\mu_v}{2} + \lambda_v\right) (\nabla \cdot \mathbf{v}) - p - \left(\frac{\mu_\varepsilon}{3} + \lambda_\varepsilon\right)\eta. \quad (5.11)$$

- For the morpho-poroelastic model from Section 4.5, we would have

$$\boldsymbol{\sigma} = \boldsymbol{\sigma}_{\text{el}}(\boldsymbol{\varepsilon}) + \boldsymbol{\sigma}_{\text{vis}}(\mathbf{v}) - p\mathbf{I} - \left(\frac{\mu_\varepsilon}{3} + \lambda_\varepsilon\right)\eta\mathbf{I}. \quad (5.12)$$

The elasticity tensor is now given by  $\boldsymbol{\sigma}_{\text{el}}(\boldsymbol{\varepsilon}) = \mu_\varepsilon \boldsymbol{\varepsilon} + \lambda_\varepsilon \text{tr}(\boldsymbol{\varepsilon})\mathbf{I}$ . We then get

$$\bar{\sigma} = \left(\frac{\mu_\varepsilon}{2} + \lambda_\varepsilon\right) \text{tr}(\boldsymbol{\varepsilon}) + \left(\frac{\mu_v}{2} + \lambda_v\right) (\nabla \cdot \mathbf{v}) - p - \left(\frac{\mu_\varepsilon}{3} + \lambda_\varepsilon\right)\eta. \quad (5.13)$$

Finally, consider the effect of  $\Phi$  on  $S$ . If  $\Phi$  is equal to 0, then there are no tumor cells present that can multiply, hence no growth is possible. On the other hand, if  $\Phi$  is equal to 1, the fluid can no longer reach the tumor to deliver oxygen. The transport of oxygen throughout the medium is modeled using a convection-diffusion equation. Convection is driven by the fluid velocity  $\mathbf{v}_F$ , which is related to the Darcy flux as follows:

$$(1 - \Phi)\mathbf{v}_F = -k\nabla p. \quad (5.14)$$

To model the absorption of oxygen into the tissue, a nutrient absorption function  $N_A(c)$  that acts as a sink is added to the equation. We assume that oxygen is the only nutrient that allows tumor cells to grow. Oxygen absorption also follows a Michaelis-Menten type equation:

$$N_A(c) = \frac{A_{\text{max}}c}{c_A + c}, \quad (5.15)$$

where  $A_{\text{max}}$  is the maximum absorption rate, and  $c_A$  is the concentration for which the absorption rate is half of its maximum. By also incorporating the deformation of the medium, and a diffusive term, we get the following transport equation for  $c$ :

$$\frac{\partial c}{\partial t} + \nabla \cdot (\mathbf{v}c) - \nabla \cdot (k\nabla pc) = \lambda_N \nabla^2 c - \Phi N_A(c) \quad (5.16)$$

The parameter  $\lambda_N$  is the oxygen (nutrient) diffusion constant. In the following sections, we derive a finite element model to solve convection-diffusion equations such as (5.16). First, a simpler equation is solved on a fixed grid and constant fluid velocity. To avoid oscillations in convection dominated systems, the Streamline Upwind Petrov-Galerkin (SUPG) method can be used. We have found this to be unnecessary in our numerical experiments, because the convection remained small compared to the diffusion, especially near the domain boundaries. In Sections 5.2.1 and 5.2.2 we consider two ‘zero-dimensional’ systems, to investigate the behaviour of nutrient absorption and its relation to tissue growth. In subsequent sections, we build up to the complete model that combines a (visco/morpho-) elastic model with oxygen transport and cell creation.

### 5.1.1 Complete Model Equations

In this section we combine the mechanical model (that is, morpho-/visco-/poroelasticity), with the biochemical tumor growth model. On the one hand, we use the model from Section 4.5 with an additional term that models the stress exerted by tumor cells on the tissue. The mechanical model equations are given by:

$$\rho \left( \frac{D\mathbf{v}}{Dt} + (\nabla \cdot \mathbf{v})\mathbf{v} \right) - \nabla \cdot \boldsymbol{\sigma}(\mathbf{v}, \boldsymbol{\varepsilon}, p, \eta) = \mathbf{g}, \quad (5.17)$$

$$\frac{D\boldsymbol{\varepsilon}}{Dt} + \boldsymbol{\varepsilon} \text{skw}(\nabla \mathbf{v}) - \text{skw}(\nabla \mathbf{v})\boldsymbol{\varepsilon} + (\text{tr}(\boldsymbol{\varepsilon}) - 1) \text{sym}(\nabla \mathbf{v}) = -\mathbf{G}, \quad (5.18)$$

$$\nabla \cdot \mathbf{v} - \nabla \cdot (k \nabla p) = f. \quad (5.19)$$

The stress tensor is given by:

$$\boldsymbol{\sigma}(\mathbf{v}, \boldsymbol{\varepsilon}, p, \eta) = \boldsymbol{\sigma}_{\text{el}}(\boldsymbol{\varepsilon}) + \boldsymbol{\sigma}_{\text{vis}}(\mathbf{v}) - p\mathbf{I} - \left( \frac{\mu_\varepsilon}{3} + \lambda_\varepsilon \right) \eta \mathbf{I} \quad (5.20)$$

The elastic and viscous stress tensor are given by (2.89) and (2.39) respectively. We set the growth tensor equal to  $\mathbf{G} = \alpha \boldsymbol{\varepsilon}$  for  $\alpha \geq 0$ , such that we can model permanent deformations. On the other hand, the biological model from [14] is used to model oxygen transport and tumor growth.

$$\frac{Dc}{Dt} + (\nabla \cdot \mathbf{v})c - \nabla \cdot (kc \nabla p + \lambda_N \nabla c) = -\Phi N_A(c), \quad (5.21)$$

$$\frac{D\Phi}{Dt} + (\nabla \cdot \mathbf{v})\Phi = F(\bar{\sigma})\Phi(1 - \Phi)N_G(c) - r_d \Phi. \quad (5.22)$$

$$\frac{\partial \eta}{\partial t} = F(\bar{\sigma})\Phi(1 - \Phi)N_G(c) - r_d \Phi. \quad (5.23)$$

Thus, we have six equations and the six corresponding unknowns  $\mathbf{v}$ ,  $\boldsymbol{\varepsilon}$ ,  $p$ ,  $c$ ,  $\Phi$  and  $\eta$ . Note however that there are nine scalar unknowns;  $\mathbf{v}$  has components  $v^1, v^2$ , and  $\boldsymbol{\varepsilon}$  has the (unique) components  $\varepsilon^{11}, \varepsilon^{12}, \varepsilon^{22}$ . We use linear basis functions with respect to a triangular mesh for all unknowns. The finite element discretization of equations (5.17)-(5.19) can be found in Section 4.5. The discretization of the term involving  $\nabla \eta$  in (5.17) is done similarly to the  $\nabla p$  term. Discretization of (5.21) and (5.22) is discussed in Section 5.3. Equation (5.23) is discussed in Section 5.3.3.

It is important to note that the mechanical and biological models are two-way coupled. The mechanical model produces the solid velocity  $\mathbf{v}$  and the pressure  $p$ , which are necessary to compute the oxygen concentration. Moreover, the average stress  $\bar{\sigma}$  is needed to calculate the volume fraction  $\Phi$  and new tissue volume  $\eta$ . On the other hand, the biological model also affects the mechanics through the  $\nabla \eta$  term in the force balance equation.

Due to the interconnectedness of the equations, we choose a monolithic approach over some kind of staggered approach. A possible staggered approach would be to first solve the mechanical model, and use its output to solve the biological model. In the next step the output of the biological model is used to solve the mechanical model again. The advantage of such a method is that two smaller systems are solved instead of one large system. However, the coefficient matrix of the mechanical model is already  $6n \times 6n$  where  $n$  is the number of meshpoints. The biological model is much smaller in comparison; only  $3n \times 3n$ . Thus, a monolithic approach will amount to solving a  $9n \times 9n$  system. Additionally, a staggered approach is often very useful for converting a non-linear system into two linear systems. If we were to use a staggered method then the variable  $p$  in equation (5.21) would be known, effectively making this term linear. However, both the mechanical and biological systems are inherently non-linear, a staggered approach will not help in this sense.

### 5.1.2 Testing Configurations

We present the different testing configurations used to test the model. Note that due to the nature of equation (5.18) and (5.22), the unknowns  $\boldsymbol{\varepsilon}$ ,  $\Phi$  and  $\eta$  do not need boundary conditions.

- (i) The first configuration is a square with one rigid edge. We set

$$\Omega(0) = (0, 1)^2, \quad \Gamma_1(0) = \{(x, y) : x = 0\}, \quad \Gamma_2(0) = \partial\Omega(0) \setminus \Gamma_1(0).$$

On  $\Gamma_1(t)$ , we impose either the boundary conditions:

$$\mathbf{v} = \mathbf{0}, \quad \mathbf{n} \cdot k \nabla p = 0, \quad c = c_b. \quad (5.24)$$

or

$$\mathbf{v} = \mathbf{0}, \quad \mathbf{n} \cdot k \nabla p = 0, \quad \mathbf{n} \cdot (ck \nabla p + \lambda_N \nabla c) = 0 \quad (5.25)$$

Thus, this part of the boundary is fixed. In the case of conditions (5.24), the nutrient concentration  $c$  is prescribed on  $\Gamma_1(t)$ . This means that nutrients can flow into and out of the domain. On  $\Gamma_2(t)$ , we impose boundary conditions

$$\boldsymbol{\sigma} \cdot \mathbf{n} = \boldsymbol{\tau}, \quad p = 0, \quad \mathbf{n} \cdot (ck \nabla p + \lambda_N \nabla c) + \kappa(c - c_0) = 0. \quad (5.26)$$

Note that  $\boldsymbol{\sigma}$  is the complete stress tensor given in (5.20). The parameter  $\kappa$  should be non-negative. Note that if  $\kappa = 0$ , we have a zero-flux boundary condition for  $c$  on  $\Gamma_2(t)$ . If in addition we use condition (5.25) on  $\Gamma_1(t)$ , nutrients can not flow into or out of the domain at any point.

- (ii) This configuration is very similar to configuration (i). The same initial domain  $\Omega(0) = (0, 1)^2$  is used. On  $\Gamma_1(t)$ , we impose

$$\mathbf{v} = \mathbf{0}, \quad \mathbf{n} \cdot k \nabla p = 0, \quad \mathbf{n} \cdot (ck \nabla p + \lambda_N \nabla c) = 0, \quad (5.27)$$

and on  $\Gamma_2(t)$  we impose

$$\boldsymbol{\sigma} \cdot \mathbf{n} + \gamma \mathbf{u} = \mathbf{0}, \quad p = 0, \quad \mathbf{n} \cdot (ck \nabla p + \lambda_N \nabla c) = 0. \quad (5.28)$$

The parameter  $\gamma$  should be non-negative. If it is zero, then  $\Gamma_2(t)$  is a free boundary and we are back in configuration (i) with  $\boldsymbol{\tau} = \mathbf{0}$ . The boundary condition on  $\boldsymbol{\sigma}$  essentially models a mass-spring system: the force on the boundary is opposite and proportional to the displacement. This models the force of the surrounding tissue as it is pressing on the domain. The weak form of equation (5.17) must also be adapted, it now becomes

$$\frac{d}{dt} \int_{\Omega(t)} \rho \mathbf{v} \cdot \boldsymbol{\varphi} \, d\Omega + \int_{\Omega(t)} \boldsymbol{\sigma} : \nabla \boldsymbol{\varphi} \, d\Omega + \int_{\Gamma_2(t)} \gamma \mathbf{u} \cdot \boldsymbol{\varphi} \, d\Gamma = \int_{\Omega(t)} \mathbf{g} \cdot \boldsymbol{\varphi} \, d\Omega. \quad (5.29)$$

This leads to a semi-discrete system of the form

$$\frac{d}{dt} (\rho M_v \mathbf{v}) + S_{\text{vis}} \mathbf{v} + S_{\text{el}} \boldsymbol{\varepsilon} - D^\top \mathbf{p} + \tilde{C} \mathbf{u} = \mathbf{g}, \quad (5.30)$$

where  $\mathbf{u} = (\mathbf{u}^1, \mathbf{u}^2)^\top$  is the vector containing the domain displacement in the grid nodes, and  $\tilde{C}$  is the block matrix

$$\tilde{C} = \begin{pmatrix} C & 0 \\ 0 & C \end{pmatrix},$$

where  $C$  is given element-wise by

$$C_{ij} = \gamma \int_{\Gamma_2(t)} \varphi_i \varphi_j \, d\Gamma. \quad (5.31)$$

Note that the displacement does not appear explicitly in any of the other equations. We use the displacement velocity to update it:

$$\mathbf{u}^m = \mathbf{u}^{m-1} + \Delta t \mathbf{v}^m. \quad (5.32)$$

Thus, integrating (5.30) in time using implicit Euler yields

$$\begin{aligned} & \left( \rho M_v^m + \Delta t S_{\text{vis}}^m + \Delta t^2 \tilde{C}^m \right) \mathbf{v}^m + \Delta t S_{\text{el}}^m \boldsymbol{\varepsilon}^m - \Delta t (D^m)^\top \mathbf{p}^m \\ & = \rho M_v^{m-1} \mathbf{v}^{m-1} - \Delta t \tilde{C}^m \mathbf{u}^{m-1} + \Delta t \mathbf{g}^m \end{aligned} \quad (5.33)$$

- (iii) Finally, we consider a quarter circle with radius  $R$ :

$$\Omega(0) = \{(x, y) : x, y > 0, x^2 + y^2 < R^2\},$$

with boundaries

$$\begin{aligned}\Gamma_{11}(0) &= \{(x, y) : 0 < x < R, y = 0\}, \\ \Gamma_{12}(0) &= \{(x, y) : x = 0, 0 < y < R\}, \\ \Gamma_2(0) &= \{(x, y) : x, y > 0, x^2 + y^2 = R^2\}.\end{aligned}$$

Figure 5.1 shows a sketch of this initial domain. The following boundary conditions are imposed:

$$\text{on } \Gamma_{11}(t) : v^2 = 0, \quad \mathbf{n} \cdot k \nabla p = 0, \quad \mathbf{n} \cdot (ck \nabla p + \lambda_N \nabla c) = 0. \quad (5.34)$$

$$\text{on } \Gamma_{12}(t) : v^1 = 0, \quad \mathbf{n} \cdot k \nabla p = 0, \quad \mathbf{n} \cdot (ck \nabla p + \lambda_N \nabla c) = 0. \quad (5.35)$$

$$\text{on } \Gamma_2(t) : \boldsymbol{\sigma} \cdot \mathbf{n} = \mathbf{0}, \quad p = 0, \quad c = c_b. \quad (5.36)$$

Thus, on  $\Gamma_{11}(t)$ , the domain can only move in the  $x$ -direction. Similarly, on  $\Gamma_{12}(t)$ , the domain can only move in the  $y$ -direction. The circle arc  $\Gamma_2(t)$  is a free boundary. In this configuration, we are essentially modelling a full circle in which there are symmetries across its vertical and horizontal radius.

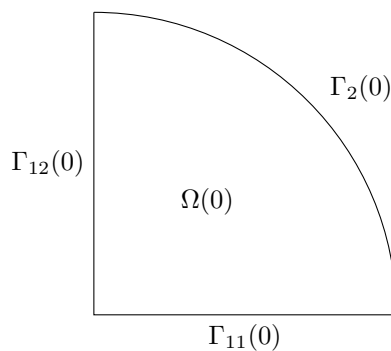


Figure 5.1: Sketch of initial domain with corresponding boundaries.

## 5.2 Analysis of Zero-Dimensional Nutrient Decay & Tissue Growth

In the following sections, equations (5.21) and (5.22) are analyzed in a zero-dimensional setting. We consider the oxygen concentration  $c$ , and later also the tumor volume fraction  $\Phi$ , at a single point and assume there is no transport. Hence the evolution of these quantities only depends on their respective source functions and their coupling.

### 5.2.1 Analysis of Nutrient Absorption

We first investigate the nutrient absorption behaviour as modeled by the Michaelis-Menten equation. To this end, we analyze the following differential equation:

$$\frac{dc}{dt} = -\Phi N_A(c) = -\frac{\Phi A_{\max} c}{c_A + c}, \quad (5.37)$$

where the concentration  $c$  is only a function of time:  $c = c(t)$ . This effectively mimics equation (5.21) in the case where  $\mathbf{v} = \mathbf{0}$ ,  $\nabla p$  is constant and  $\lambda_N = 0$ ; there is no domain deformation, no Darcy flux and no nutrient diffusion respectively. In other words, there is no transport of oxygen and the evolution of the concentration is only determined by absorption. In this section, we assume  $\Phi$  is constant in time and treat  $\Phi A_{\max}$  as being one parameter. Later it will make sense to separate them, as  $\Phi$  will no longer be constant in time. Equation (5.37) is separable:

$$\int \frac{c_A}{c} + 1 \, dc = - \int \Phi A_{\max} \, dt. \quad (5.38)$$

Which is solved by

$$c_A \ln |c| + c = -\Phi A_{\max} t + K, \quad (5.39)$$

where  $K$  is an integration constant. From the original differential equation it is clear that  $dc/dt < 0$  for  $c > 0$  and  $dc/dt = 0$  if and only if  $c = 0$ . Therefore if  $c(0) > 0$  we have  $c(t) > 0$  for all  $t \geq 0$ . With this in mind we can drop the absolute value and eliminate  $K$  in favour of the (positive) initial condition  $c(0) = c_0$ :

$$c_A \ln \left( \frac{c(t)}{c_0} \right) + c(t) - c_0 = -\Phi A_{\max} t. \quad (5.40)$$

We can not solve the above equation for  $c(t)$  in terms of elementary functions. However, we can either analyze the curve in  $t$ - $c$  space defined by the implicit relation, or we can solve it using the Lambert  $W$  function [4]. This function is defined as the solution  $W(z)$  to

$$W(z)e^{W(z)} = z. \quad (5.41)$$

The solution is uniquely determined for real  $z > 0$ . To solve equation (5.40) in terms of the Lambert  $W$  function, we will need the following result.

**Lemma 9.** *Let the function  $f(z)$  be implicitly defined by*

$$f(z)^n e^{f(z)} = z, \quad (5.42)$$

for  $z > 0$  and  $n > 0$ . Then

$$f(z) = nW \left( \frac{z^{1/n}}{n} \right). \quad (5.43)$$

*Proof.* We can prove the statement by simply substituting (5.43) into (5.42):

$$\begin{aligned} f(z)^n e^{f(z)} &= \left( nW \left( \frac{z^{1/n}}{n} \right) \right)^n \exp \left[ nW \left( \frac{z^{1/n}}{n} \right) \right] \\ &= \left( nW \left( \frac{z^{1/n}}{n} \right) \exp \left[ W \left( \frac{z^{1/n}}{n} \right) \right] \right)^n \\ &= \left( n \frac{z^{1/n}}{n} \right)^n \\ &= z \end{aligned}$$

Here we have used definition (5.41). □

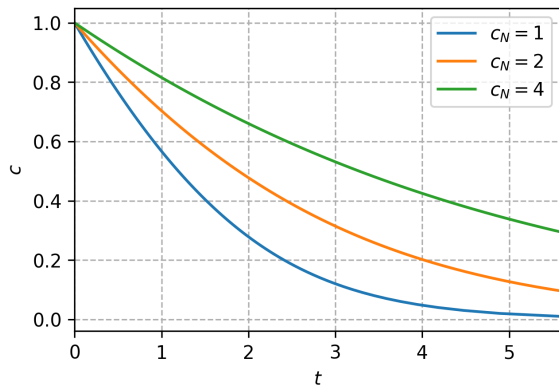
To apply the result of Lemma 9, we first write (5.40) in the form

$$c(t)^{c_A} e^{c(t)} = c_0^{c_A} \exp [c_0 - \Phi A_{\max} t], \quad (5.44)$$

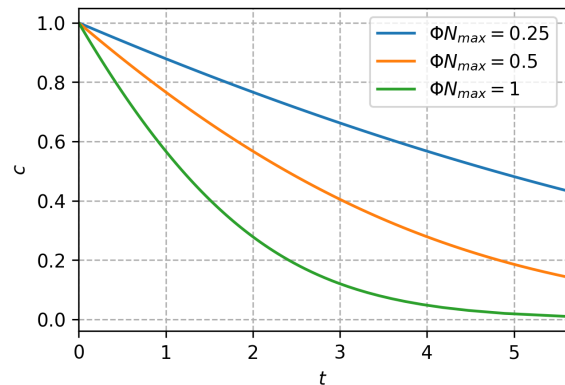
by taking the exponential on both sides. It then follows that

$$c(t) = c_A W \left( \frac{c_0}{c_A} \exp \left[ \frac{c_0 - \Phi A_{\max} t}{c_A} \right] \right) \quad (5.45)$$

Figure 5.2a and 5.2b show the effect of  $c_A$  and  $\Phi A_{\max}$  on the evolution of the concentration. The other parameters are kept at 1 for the sake of demonstration.



(a) Concentration for various values of  $c_A$ .



(b) Concentration for various values of  $\Phi A_{\max}$ .

## 5.2.2 Coupling Nutrient Absorption and Volume Fraction Growth

To extend the model of the previous section, the solid volume fraction  $\Phi$  is added to the system. We use equation (5.3) to model its growth or decay. Since  $c$  and  $\Phi$  are assumed to only be functions of time, the stress dependent factor  $F(\bar{\sigma})$  is assumed to be constant. Note that because there is no convective transport,  $\Phi$  and  $\eta$  are equal. Hence,  $\Phi$  is a measure of tissue growth. We get the following coupled system of differential equations:

$$\frac{dc}{dt} = -\Phi N_A(c) \quad (5.46)$$

$$\frac{d\Phi}{dt} = F\Phi(1 - \Phi)N_G(c) - r_d\Phi \quad (5.47)$$

The system is numerically integrated using the explicit Runge-Kutta 4 method. We use the initial conditions  $c(0) = 1$  and  $\Phi(0) = 0.01$  to simulate a situation in which there is initially almost no cell tissue present, such that the available nutrients will cause tissue growth. We set the cell creation- and nutrient absorption parameters to 1:

$$G_{\max} = c_G = A_{\max} = c_A = 1, \quad (5.48)$$

and we also set  $F = 1$ . The death rate  $r_d$  is varied to investigate the resulting behaviour. The system is integrated from  $t = 0$  to  $t = 100$  using a timestep of size  $\Delta t = 0.1$ . Figure 5.2 shows the solutions corresponding to  $r_d = 0$ ,  $r_d = 0.01$  and  $r_d = 0.1$ .

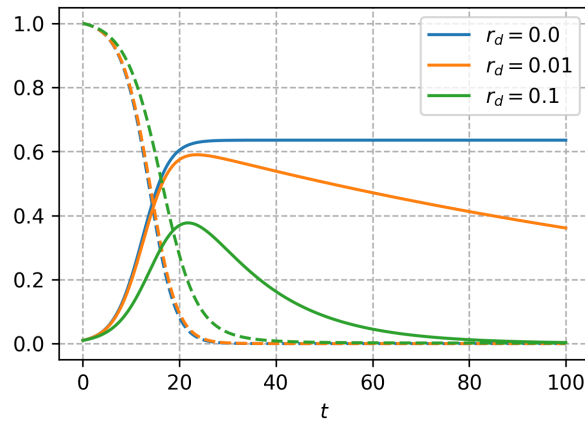


Figure 5.2: Plot of concentration (dashed line) and volume fraction (solid line) as a function of time for various values of  $r_d$ .

Initially  $\Phi$  starts growing due to the high amount of nutrients. After the nutrient concentration gradually declines, growth slows down and the death rate takes over. In the case of  $r_d = 0$ ,  $\Phi$  converges to a positive asymptotic value, whereas for  $r_d > 0$  it converges to 0.

## 5.3 Finite Element Approach

In this section, the weak form and corresponding Galerkin equations of equations (5.21), (5.22) and (5.23) are derived.

### 5.3.1 Nutrient Transport Equation

We start with the nutrient transport equation (5.21). Note that it is a convection-diffusion equation with a non-linear sink term. We can distinguish three types of terms in equation (5.21); a term accounting for the domain deformation, a concentration flux, and a concentration sink:

$$\underbrace{\frac{Dc}{Dt}}_{\text{deformation}} + (\nabla \cdot \mathbf{v})c - \underbrace{\nabla \cdot (ck\nabla p + \lambda_N \nabla c)}_{\text{flux}} = \underbrace{-\Phi N_A(c)}_{\text{sink}}, \quad (5.49)$$

where the nutrient absorption rate  $N_A(c)$  is given in (5.15). For the derivation of the weak form we use a zero-flux boundary condition on the entire boundary:

$$\mathbf{n} \cdot (ck\nabla p + \lambda_N \nabla c) = 0. \quad (5.50)$$

This boundary condition is also imposed in testing configuration (ii) in Section 5.1.2. Note that if a Robin boundary condition is imposed, such as in configuration (i), then some additional terms will appear in the weak form. Let  $\xi$  be a Lagrangian basis function. Multiply the above equation by  $\xi$  and integrate over  $\Omega(t)$  to obtain:

$$\frac{d}{dt} \int_{\Omega(t)} c \xi \, d\Omega - \int_{\Omega(t)} \nabla \cdot (ck \nabla p + \lambda_N \nabla c) \xi \, d\Omega = - \int_{\Omega(t)} \Phi N_A(c) \xi \, d\Omega. \quad (5.51)$$

Here we have applied Theorem 3 to obtain the first term. Use partial integration and Gauss' divergence theorem to find

$$\begin{aligned} \frac{d}{dt} \int_{\Omega(t)} c \xi \, d\Omega + \int_{\Omega(t)} ck \nabla p \cdot \nabla \xi + \lambda_N \nabla c \cdot \nabla \xi \, d\Omega - \int_{\Gamma(t)} \mathbf{n} \cdot (ck \nabla p + \lambda_N \nabla c) \xi \, d\Gamma \\ = - \int_{\Omega(t)} \Phi N_A(c) \xi \, d\Omega. \end{aligned} \quad (5.52)$$

Applying the boundary conditions yields

$$\frac{d}{dt} \int_{\Omega(t)} c \xi \, d\Omega + \int_{\Omega(t)} ck \nabla p \cdot \nabla \xi + \lambda_N \nabla c \cdot \nabla \xi \, d\Omega = - \int_{\Omega(t)} \Phi N_A(c) \xi \, d\Omega. \quad (5.53)$$

This weak form introduces two difficulties: we obtain a non-linear term involving both  $c$  and  $p$ , and on the right-hand side we get a non-linear term involving  $\Phi$  and a (non-linear) function of  $c$ . Using linear triangular elements, we can write

$$c(\mathbf{x}, t) = \sum_{j=1}^n c_j(t) \varphi_j(\mathbf{x}; t), \quad \Phi(\mathbf{x}, t) = \sum_{k=1}^n \Phi_k(t) \varphi_k(\mathbf{x}; t), \quad p(\mathbf{x}, t) = \sum_{\ell=1}^n p_\ell(t) \varphi_\ell(\mathbf{x}; t), \quad (5.54)$$

where the  $\varphi_j$  are linear basis functions with respect to a mesh  $\mathcal{T}_h$ . Set  $\xi = \varphi_i$  for some  $i \in \{1, \dots, n\}$ , then we get

$$\int_{\Omega(t)} ck \nabla p \cdot \nabla \xi + \lambda_N \nabla c \cdot \nabla \xi \, d\Omega = \sum_{j=1}^n c_j \int_{\Omega(t)} k \varphi_j \sum_{\ell=1}^n p_\ell \nabla \varphi_\ell \cdot \nabla \varphi_i + \lambda_N \nabla \varphi_j \cdot \nabla \varphi_i \, d\Omega. \quad (5.55)$$

Regarding the right-hand side, we approximate it in the following way:

$$\begin{aligned} \int_{\Omega(t)} \Phi N_A(c) \varphi_i \, d\Omega &= \int_{\Omega(t)} \frac{A_{\max} c \Phi \varphi_i}{c_A + c} \, d\Omega \\ &= \sum_{j=1}^n c_j \int_{\Omega(t)} \frac{A_{\max} \Phi \varphi_i \varphi_j}{c_A + c} \, d\Omega \\ &\approx \sum_{j=1}^n c_j \sum_{T \in \mathcal{T}_h} \tilde{A}_T(\mathbf{c}) \int_T \Phi \varphi_i \varphi_j \, d\Omega, \\ &= \sum_{j=1}^n c_j \sum_{T \in \mathcal{T}_h} \tilde{A}_T(\mathbf{c}) \sum_{\ell \in \{1, 2, 3\}} \Phi_{\ell, T} \int_T \varphi_{\ell, T} \varphi_i \varphi_j \, d\Omega. \end{aligned} \quad (5.56)$$

Here  $\tilde{A}_T(\mathbf{c})$  is an approximation of  $A_{\max}/(c_A + c)$  on the element  $T$ :

$$\tilde{A}_T(\mathbf{c}) = \frac{1}{3} \left( \frac{A_{\max}}{c_A + c_{1, T}} + \frac{A_{\max}}{c_A + c_{2, T}} + \frac{A_{\max}}{c_A + c_{3, T}} \right), \quad (5.57)$$

Here  $\Phi_{\ell, T}$  and  $c_{\ell, T}$  for  $\ell \in \{1, 2, 3\}$  denote the values of  $\Phi$  and  $c$  at the vertices of  $T$ . In conclusion, we get the following semi-discretized system of Galerkin equations:

$$\frac{d}{dt} (M \mathbf{c}) + [kA(\mathbf{p}) + \lambda_N L + N_A(\Phi, \mathbf{c})] \mathbf{c} = \mathbf{0}, \quad (5.58)$$

where the coefficient matrices are given element-wise by:

$$M_{ij} = \int_{\Omega(t)} \varphi_i \varphi_j \, d\Omega, \quad (5.59)$$

$$L_{ij} = \int_{\Omega(t)} \nabla \varphi_i \cdot \nabla \varphi_j \, d\Omega, \quad (5.60)$$

$$A(\mathbf{p})_{ij} = \sum_{\ell=1}^n p_\ell \int_{\Omega(t)} \varphi_j \nabla \varphi_\ell \cdot \nabla \varphi_i \, d\Omega. \quad (5.61)$$

$$N_A(\Phi, \mathbf{c})_{ij} = \sum_{T \in \mathcal{T}_h} \tilde{A}_T(\mathbf{c}) \sum_{\ell \in \{1,2,3\}} \Phi_{\ell,T} \int_T \varphi_{\ell,T} \varphi_i \varphi_j \, d\Omega. \quad (5.62)$$

for  $i, j \in \{1, \dots, n\}$ . The time-discretization is done in the usual way using implicit Euler:

$$M^m \mathbf{c}^m + \Delta t [kA^m(\mathbf{p}^m) + \lambda_N L^m + N_A^m(\Phi^m, \mathbf{c}^m)] \mathbf{c}^m = M^{m-1} \mathbf{c}^{m-1}. \quad (5.63)$$

### 5.3.2 Volume Fraction Evolution

In this section, we discretize equation (5.22). For readability purposes, we repeat the equation:

$$\frac{D\Phi}{Dt} + (\nabla \cdot \mathbf{v})\Phi = F(\bar{\sigma})\Phi(1 - \Phi)N_G(c) - r_d \Phi. \quad (5.64)$$

By multiplying the equation with  $\varphi_i$  for  $i = 1, \dots, n$  and integrating over  $\Omega(t)$ , we obtain a system of the form

$$\frac{d}{dt}(M\Phi) = \mathcal{S} - r_d M\Phi. \quad (5.65)$$

The first term is again the result of applying Theorem 3. Moreover, the vector  $\mathcal{S}$  is given element-wise by:

$$\begin{aligned} \mathcal{S}_i = \sum_{T \in \mathcal{T}_h} F(\bar{\sigma}_T) \tilde{G}_T(\mathbf{c}) \sum_{j=1}^n c_j \left[ \sum_{k \in \{1,2,3\}} \Phi_{k,T} \int_T \varphi_i \varphi_j \varphi_{k,T} \, d\Omega \right. \\ \left. - \sum_{k \in \{1,2,3\}} \sum_{\ell \in \{1,2,3\}} \Phi_{k,T} \Phi_{\ell,T} \int_T \varphi_i \varphi_j \varphi_{k,T} \varphi_{\ell,T} \, d\Omega \right]. \end{aligned} \quad (5.66)$$

The element matrices corresponding to the integrals in the above expression, as well as in (5.56), are computed in Section A.4. Similar to (5.57), we use

$$\tilde{G}_T(\mathbf{c}) = \frac{1}{3} \left( \frac{G_{\max}}{c_G + c_{1,T}} + \frac{G_{\max}}{c_G + c_{2,T}} + \frac{G_{\max}}{c_G + c_{3,T}} \right) \quad (5.67)$$

as an approximation of  $\tilde{G}(c) = G_{\max}/(c + c_G)$  on  $T$ . The factor  $\bar{\sigma}_T$  is the average over a triangle  $T$  of the total stress (5.13). Thus we have

$$\begin{aligned} \bar{\sigma}_T = \sum_{\ell \in \{1,2,3\}} \left[ \frac{1}{3} \left( \frac{\mu_\varepsilon}{2} + \lambda_\varepsilon \right) (\varepsilon_{\ell,T}^{11} + \varepsilon_{\ell,T}^{22}) + \left( \frac{\mu_v}{2} + \lambda_v \right) (v_{\ell,T}^1 \alpha_{\ell,T} + v_{\ell,T}^2 \beta_{\ell,T}) \right. \\ \left. - \frac{1}{3} p_{\ell,T} - \frac{1}{3} \left( \frac{\mu_\varepsilon}{3} + \lambda_\varepsilon \right) \eta_{\ell,T} \right], \end{aligned} \quad (5.68)$$

where  $\alpha_{\ell,T}$  and  $\beta_{\ell,T}$  are respectively the  $x$ - and  $y$ -derivative of the linear basis function  $\varphi_{\ell,T}$ , and a subscript  $(\ell, T)$  denotes one of the three vertices of  $T$ . Discretizing the semi-discrete system (5.65) in time using implicit Euler yields:

$$(1 + \Delta t r_d) M^{m+1} \Phi^{m+1} - \Delta t F \mathbf{s}^{m+1}(\Phi^{m+1}, \mathbf{c}^{m+1}) = M^m \Phi^m \quad (5.69)$$



### 5.3.3 Updating New Tissue Volume

In this section, we will discuss the treatment of the new tissue volume  $\eta$ . Note that equations (5.22) and (5.23) are very similar. The evolution of  $\Phi$  takes into account the domain deformation, whereas the evolution of  $\eta$  does not. This corresponds to their physical interpretations;  $\Phi$  represents the fraction of space occupied by tumor cells, whereas  $\eta$  represents the volume of new tumor cells *per unit of volume*. Thus, where  $\Phi$  is affected by the flow of the domain,  $\eta$  is not. We derive the weak form of equation (5.23). Since we are using a Lagrangian framework - the value of the unknowns is tracked at the location of a moving grid point - the partial time-derivative  $\partial\eta/\partial t$  cannot be easily discretized, because we must compensate for the displacement velocity  $\mathbf{v}$ . To this end, let  $\xi$  be a scalar-valued linear Lagrangian basis function, then:

$$\begin{aligned} \int_{\Omega(t)} \frac{\partial\eta}{\partial t} \xi \, d\Omega &= \int_{\Omega(t)} \left( \frac{\partial\eta}{\partial t} + \nabla \cdot (\mathbf{v}\eta) \right) \xi \, d\Omega - \int_{\Omega(t)} \nabla \cdot (\mathbf{v}\eta) \xi \, d\Omega \\ &= \frac{d}{dt} \int_{\Omega(t)} \eta \xi \, d\Omega - \int_{\Omega(t)} ((\nabla \cdot \mathbf{v})\eta + \mathbf{v} \cdot \nabla \eta) \xi \, d\Omega. \end{aligned} \quad (5.70)$$

Note that we added and subtracted the term  $\nabla \cdot (\mathbf{v}\eta)\xi$ . This then yields an integral to which Theorem 3 can be applied. We now write

$$\eta(\mathbf{x}, t) = \sum_{j=1}^n \eta_j(t) \varphi_j(\mathbf{x}; t),$$

where  $\varphi_j$  is a linear basis function centered in grid node  $j$ . Moreover, let  $\boldsymbol{\eta}$  be the vector containing the  $\eta_j$ . Set  $\xi = \varphi_i$  for  $i = 1, \dots, n$  to obtain the following semi-discrete system for (5.23):

$$\frac{d}{dt}(M\boldsymbol{\eta}) - Q(\mathbf{v})\boldsymbol{\eta} = \boldsymbol{\mathcal{S}} - r_d M \boldsymbol{\Phi}, \quad (5.71)$$

where the vector  $\boldsymbol{\mathcal{S}}$  is given element-wise by (5.66), and

$$Q(\mathbf{v})_{ij} = \int_{\Omega(t)} (\nabla \cdot \mathbf{v}) \varphi_i \varphi_j + (\mathbf{v} \cdot \nabla \varphi_j) \varphi_i \, d\Omega. \quad (5.72)$$

The per-element contributions of the matrix  $Q(\mathbf{v})$  are given by:

$$(Q(\mathbf{v})_T)_{ij} = \sum_{\ell \in \{1,2,3\}} \left[ (v_{\ell,T}^1 \alpha_{\ell,T} + v_{\ell,T}^2 \beta_{\ell,T}) \int_T \varphi_i \varphi_j \, d\Omega + (v_{\ell,T}^1 \alpha_j + v_{\ell,T}^2 \beta_j) \int_T \varphi_{\ell,T} \varphi_i \, d\Omega \right]. \quad (5.73)$$

Applying implicit Euler to (5.71) yields

$$(M^m - \Delta t Q^m(\mathbf{v}^m)) \boldsymbol{\eta}^m - \Delta t \boldsymbol{\mathcal{S}}^m + \Delta t r_d M^m \boldsymbol{\Phi}^m = M^{m-1} \boldsymbol{\eta}^{m-1} \quad (5.74)$$

Recall that  $\boldsymbol{\Phi}$  satisfies the semi-discrete system (5.65). By subtracting this equation from (5.71), we obtain

$$\frac{d}{dt} (M(\boldsymbol{\eta} - \boldsymbol{\Phi})) = Q(\mathbf{v})\boldsymbol{\eta}. \quad (5.75)$$

Thus, if  $\mathbf{v}$  remains small in both magnitude and gradient, the term on the right-hand side will also be small. In other words, if the domain deforms slowly, we expect  $\boldsymbol{\eta} - \boldsymbol{\Phi}$  to remain approximately constant in time. Moreover, if their initial values are equal, i.e.  $\boldsymbol{\eta}^0 = \boldsymbol{\Phi}^0$ , then  $\boldsymbol{\eta}$  and  $\boldsymbol{\Phi}$  will be roughly equal on a slowly deforming domain. This is illustrated in Figure 5.3. Here we have applied a body force  $\mathbf{g} = (1, 0)^\top$  to speed up the domain deformation. Both  $\boldsymbol{\eta}$  and  $\boldsymbol{\Phi}$  have been set to the constant initial value 0.5. It is clear that the difference between  $\boldsymbol{\eta}$  and  $\boldsymbol{\Phi}$  is very small, with the maximum error being less than 1% relative to both  $\boldsymbol{\eta}$  and  $\boldsymbol{\Phi}$ .

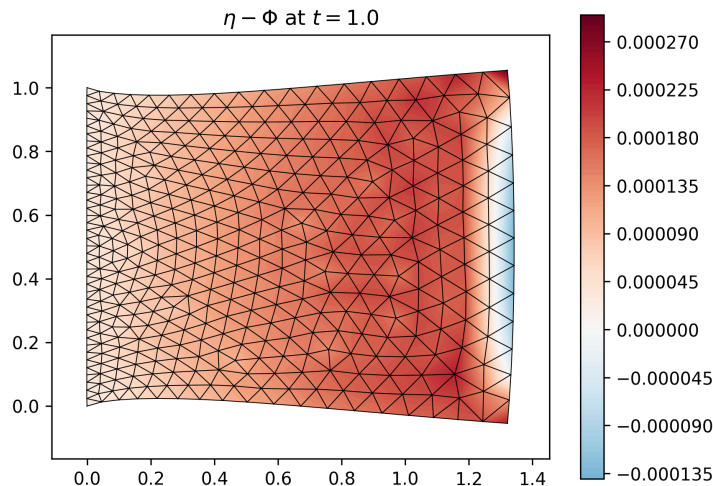


Figure 5.3: Contour plot of  $\eta - \Phi$  after 10 timesteps, superimposed on the deformed grid.

### 5.3.4 Summary of Finite Element Model

We present the complete finite element tumor growth model. First, consider the semi-discrete mechanical model corresponding to the partial differential equations (5.17) - (5.19):

$$\frac{d}{dt} (\rho M_v \mathbf{v}) + S_{\text{vis}} \mathbf{v} + S_{\text{el}} \boldsymbol{\varepsilon} - D^\top \mathbf{p} - \left( \frac{\mu_\varepsilon}{3} + \lambda_\varepsilon \right) D^\top \boldsymbol{\eta} = \mathbf{g}, \quad (5.76)$$

$$\frac{d}{dt} (M_\varepsilon \boldsymbol{\varepsilon}) - B \mathbf{v} + \alpha M_\varepsilon \boldsymbol{\varepsilon} + \mathcal{N}(\mathbf{v}, \boldsymbol{\varepsilon}) = \mathbf{0}, \quad (5.77)$$

$$D \mathbf{v} + k L \mathbf{p} = \mathbf{f}. \quad (5.78)$$

Here  $\mathbf{v}$  and  $\boldsymbol{\varepsilon}$  are the vectors  $\mathbf{v} = (\mathbf{v}^1, \mathbf{v}^2)^\top$  and  $\boldsymbol{\varepsilon} = (\boldsymbol{\varepsilon}^{11}, \boldsymbol{\varepsilon}^{12}, \boldsymbol{\varepsilon}^{22})^\top$ . The mass matrices  $M_v$  and  $M_\varepsilon$  denote block matrices with respectively 2 and 3 standard mass matrices  $M$  on their block-diagonals, and zeroes everywhere else. The coupled semi-discrete biochemical model, corresponding to equations (5.21) - (5.23), is given by:

$$\frac{d}{dt} (M \mathbf{c}) + [k A(\mathbf{p}) + \lambda_N L + N_A(\boldsymbol{\Phi}, \mathbf{c})] \mathbf{c} = \mathbf{0} \quad (5.79)$$

$$\frac{d}{dt} (M \boldsymbol{\Phi}) + r_d M \boldsymbol{\Phi} = \mathcal{S} \quad (5.80)$$

$$\frac{d}{dt} (M \boldsymbol{\eta}) - Q(\mathbf{v}) \boldsymbol{\eta} + r_d M \boldsymbol{\Phi} = \mathcal{S} \quad (5.81)$$

The vector  $\mathcal{S}$  depends on all unknowns;  $\mathcal{S} = \mathcal{S}(\mathbf{v}, \boldsymbol{\varepsilon}, \mathbf{p}, \mathbf{c}, \boldsymbol{\Phi}, \boldsymbol{\eta})$ . We have omitted this dependence in the above equations to increase readability. As was previously discussed, the two systems are solved simultaneously: we use a monolithic approach. The time-discretization is again done using implicit Euler. We will not give the fully discrete equations since they follow immediately from (5.76) - (5.81), and have already been shown in the previous sections. Note that due to the moving domain and the many non-linearities in the equations, an iterative method must be used to perform one timestep. We use a Picard-type iterative scheme, for more details see Section 2.4.3.

It should be noted that depending on the boundary conditions, additional terms must be added to one or more of the above equations. For example, if configuration (i) from Section 5.1.2 is used with  $\kappa > 0$ , equation (5.79) will gain additional terms involving boundary integrals over  $\Gamma_2(t)$ . Similarly, in configuration (ii), equation (5.76) will gain an additional term involving the boundary integral of the displacement, as is shown in (5.29).

## 5.4 Numerical Experiments

The tumor growth model is tested using the configurations from Section 5.1.2. The numerical experiments in this section should be considered as demonstrations of the finite element model. The simulations are not yet intended to model biologically correct tumor growth.

### 5.4.1 Configuration (i)

#### Domain Deformation Resulting from Tumor Cells

For our first experiment, we are interested in the effect of  $\Phi$  on the domain deformation. The experiment is performed in the framework of configuration (i). We choose the following initial profile for  $\Phi$ :

$$\Phi(x, y, 0) = 0.01 + 0.98 \cdot \mathbb{1}_{[0.5,1]}(x)(x - 0.5) \quad (5.82)$$

Here  $\mathbb{1}_A$  denotes the indicator function of the set  $A$ . For  $x < 0.5$ ,  $\Phi$  has the constant value of 0.01. It is deliberately not set to zero, because then the source function  $S$  would vanish on this part of the domain. For  $x \geq 0.5$ ,  $\Phi$  increases linearly to a value of 0.5 at  $x = 1$ . Note that this means that there are already tumor cells present. The initial concentration profile is set to the constant value 1. We expect the gradients to be higher on the right side of the domain. Therefore, the following variable initial meshwidth is used:

$$h(x, y) = 0.03 + 0.07(1 - x) \quad (5.83)$$

Figure 5.4 shows the initial volume fraction and initial finite element mesh.

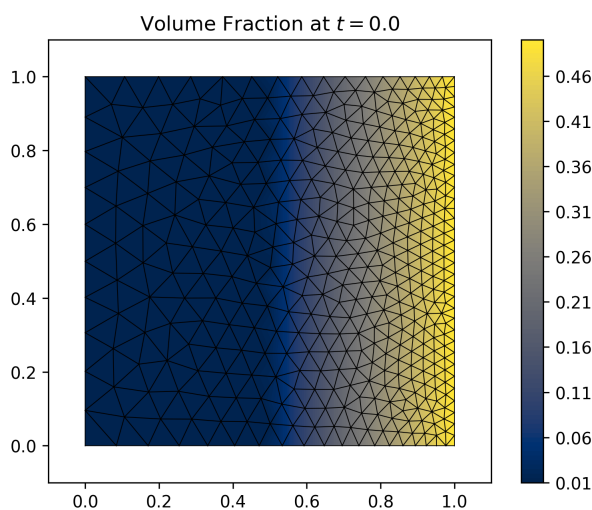


Figure 5.4: Contour plot of initial volume fraction profile (5.82) and initial grid using the variable meshwidth (5.83).

We run the morpho-visco-poroelastic tumor growth model using the simulation parameters in Tables 5.1, 5.2, 5.3.

Table 5.1: Mechanical parameters.

$\rho$	$\mu_\varepsilon$	$\lambda_\varepsilon$	$\mu_v$	$\lambda_v$	$k$	$\alpha$
1	1	1	1	1	1	1

Table 5.2: Biological parameters.

$\lambda_N$	$A_{\max}$	$c_A$	$G_{\max}$	$c_G$	$r_d$
0.01	1	1	1	1	0.1

Table 5.3: Source functions and boundary conditions.

$\mathbf{g}$	$f$	$\boldsymbol{\tau}$	$c_0$	$c_b$	$\kappa$
$\mathbf{0}$	0	$\mathbf{0}$	0	0	1

The model is run until  $t = 5$  using a timestep of size  $\Delta t = 0.1$ . Figures 5.5a and 5.5b show the finite element solutions for  $c$  and  $\Phi$ . Figure 5.5c shows the pressure solution. It is important to note that

there is no body force or shear stress applied to the medium; all deformations are a direct result of the presence of tumor cells.

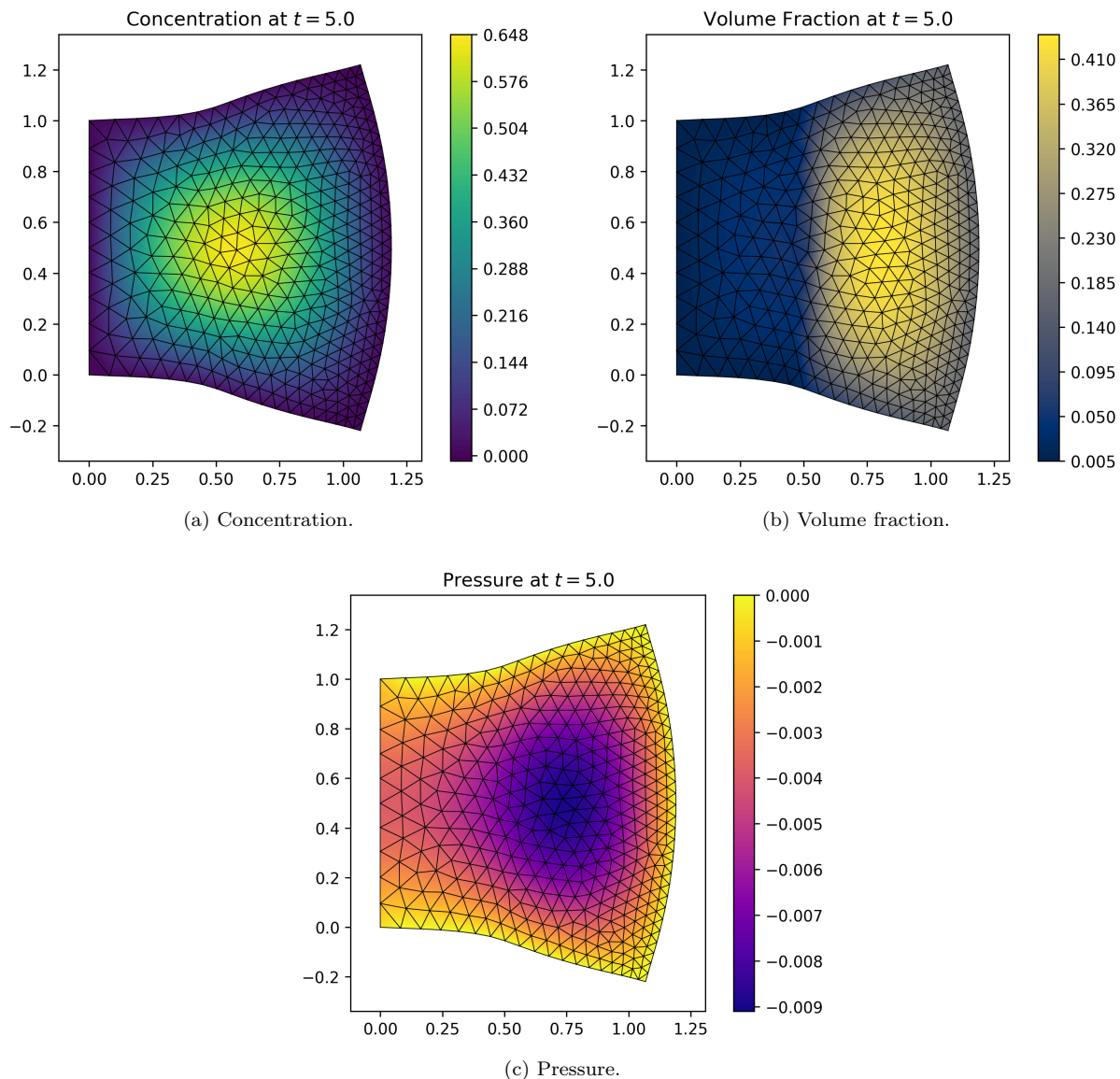


Figure 5.5: Finite element solutions for  $c$  and  $\Phi$ , superimposed on the deformed grid at  $t = 5$ . The corresponding simulation parameters are shown in Tables 5.1, 5.2 and 5.3.

### Measuring Total Concentration & Volume Fraction

For our second experiment we are interested in the evolution of the total concentration  $c_{\text{tot}}$  and total volume fraction  $\Phi_{\text{tot}}$  as functions of time. They are defined as:

$$c_{\text{tot}}(t) = \int_{\Omega(t)} c \, d\Omega, \quad \Phi_{\text{tot}}(t) = \int_{\Omega(t)} \Phi \, d\Omega. \quad (5.84)$$

The initial domain is again  $\Omega(0) = (0, 1)^2$ . We use a constant initial volume fraction/new tissue volume of  $\Phi(x, y, 0) = \eta(x, y, 0) = 0.01$ , and the initial concentration profile is given by

$$c(x, y, 0) = 16xy(1-x)(1-y). \quad (5.85)$$

We use the following simulation parameters:

Table 5.4: Mechanical parameters.

$\rho$	$\mu_\varepsilon$	$\lambda_\varepsilon$	$\mu_v$	$\lambda_v$	$k$	$\alpha$
1	1	1	1	1	1	0.1

Table 5.5: Biological parameters.

$\lambda_N$	$A_{\max}$	$c_A$	$G_{\max}$	$c_G$	$r_d$
0.01	2	1	2	1	0.1

Table 5.6: Source functions and boundary conditions.

$\mathbf{g}$	$f$	$\boldsymbol{\tau}$	$c_0$	$c_b$	$\kappa$
$\mathbf{0}$	0	$\mathbf{0}$	0	0	0

A constant (initial) meshwidth of  $h = 0.05\sqrt{2}$  is used. We run the model from  $t = 0$  to  $t = 50$  using a timestep of size  $\Delta t = 0.25$ . Figures 5.6a, 5.6b and 5.6c show the concentration, new tissue volume and pressure respectively at  $t = 10$ .

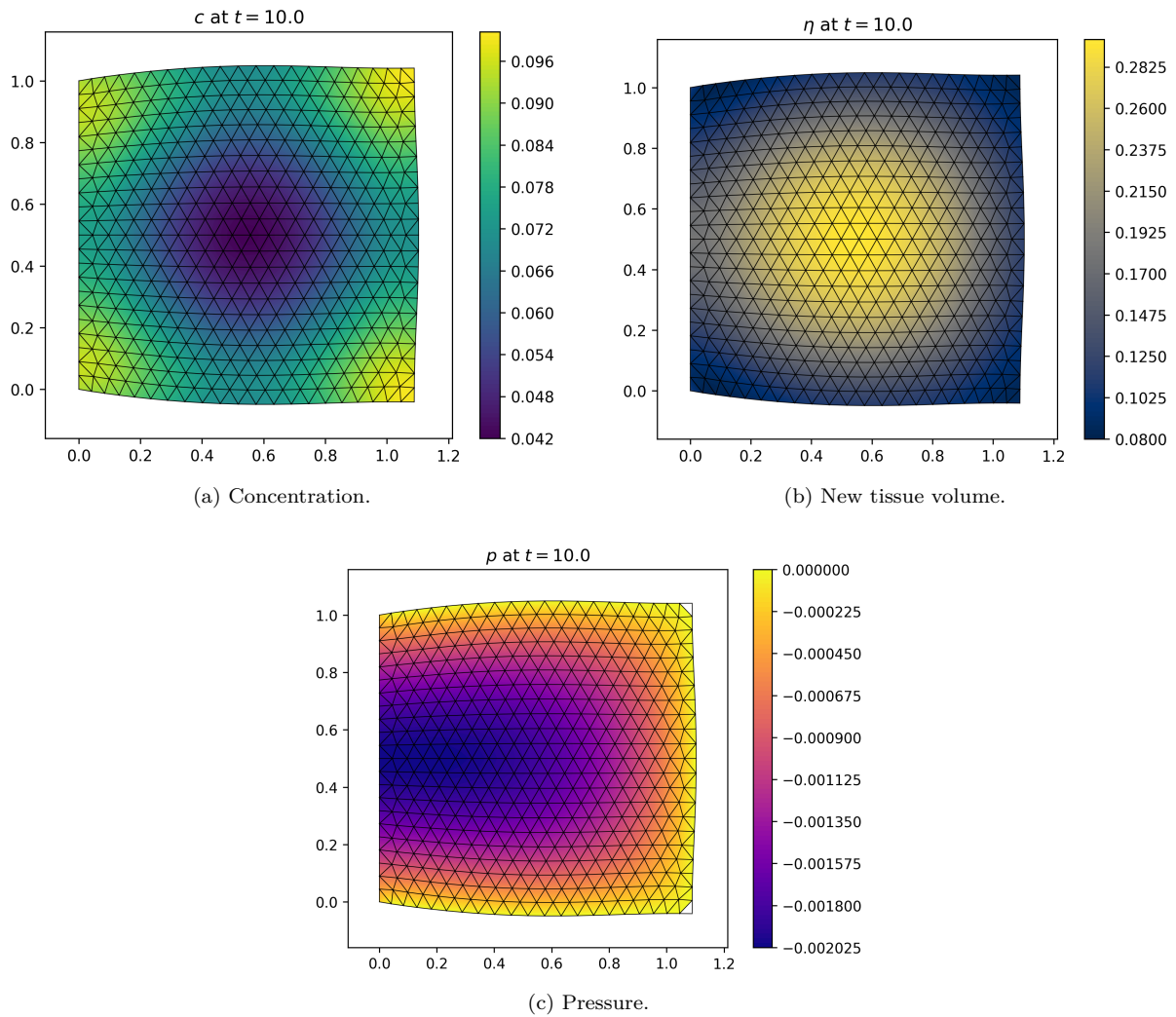


Figure 5.6: Contour plots of finite element solutions for  $c$ ,  $\eta$  and  $p$ , superimposed on the deformed grid at  $t = 10$ . The corresponding simulation parameters are shown in Tables 5.4, 5.5 and 5.6.

Figure 5.7 shows the evolution of the total concentration, total new tissue volume and domain area. Note that the evolution of  $c_{tot}$  and  $\eta_{tot}$  look qualitatively similar to that of  $c$  and  $\Phi$  in the zero-dimensional case, see Figure 5.2. Moreover, because we use  $\alpha = 0.1$ , the domain has become permanently deformed. Thus, even when the tumor cells have all died, the domain does not bounce back to its original state.

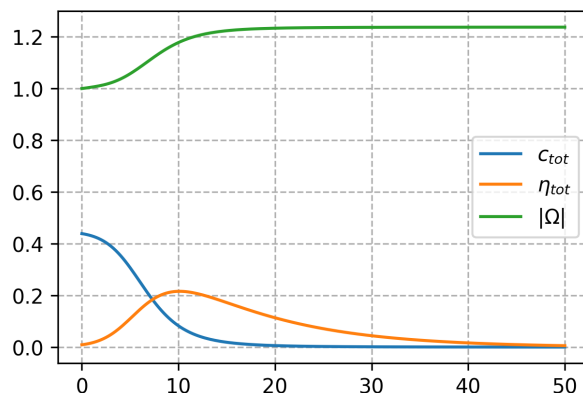


Figure 5.7: Total concentration, total new tissue volume and domain area as functions of time.

### 5.4.2 Configuration (ii)

Recall that compared to configuration (i), in configuration (ii) we replace the condition  $\boldsymbol{\sigma} \cdot \mathbf{n} = \boldsymbol{\tau}$  with  $\boldsymbol{\sigma} \cdot \mathbf{n} + \gamma \mathbf{u} = 0$  on  $\Gamma_2(t)$ . To investigate the impact of the ‘pushback’ force, the same initial domain, initial conditions and simulation parameters are used as in the previous section. They can be found in Tables 5.4, 5.5 and 5.6. Figure 5.8 shows the effect of  $\gamma$  on the evolution of the domain area. Note that for  $\gamma = 0$ , there is no shear force acting on the boundary as a result of the displacement. For higher values of  $\gamma$ , the shear force becomes larger, prohibiting growth and causing the domain to converge to its original state after  $\eta$  has vanished.

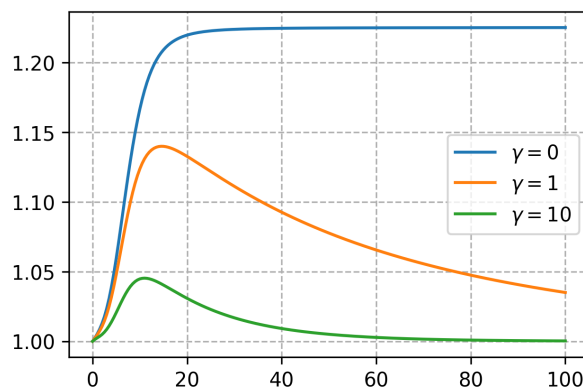


Figure 5.8: Domain area as a function of time for various values of  $\gamma$ .

### 5.4.3 Configuration (iii)

We test configuration (iii) using a quarter circle of radius  $R = 1$ . We use the following initial concentration profile:

$$c(0, x, y) = \frac{5}{\pi} (1 - (x^2 + y^2)^4) (1 + \sin(4 \arctan(y/x))). \quad (5.86)$$

The factor  $5/\pi$  ensures that  $c_{tot}(0) = 1$ . Moreover, note that  $\arctan(y/x)$  equals the angular coordinate  $\varphi$  in polar coordinates. Thus, we can also denote this concentration by:

$$c(0, r, \varphi) = \frac{5}{\pi} (1 - r^8) (1 + \sin(4\varphi)). \quad (5.87)$$

Figure 5.9 shows a contour plot of the above function. Note that for  $\varphi < \pi/4$  there is a higher concentration than for  $\varphi > \pi/4$ . We set  $c = 0$  on the free boundary  $\Gamma_2(t)$ . The term  $1 - r^8$  ensures that the

initial concentration goes smoothly to its boundary value, whilst at the same time keeping the oscillating behaviour of the sine function in the interior of the domain.

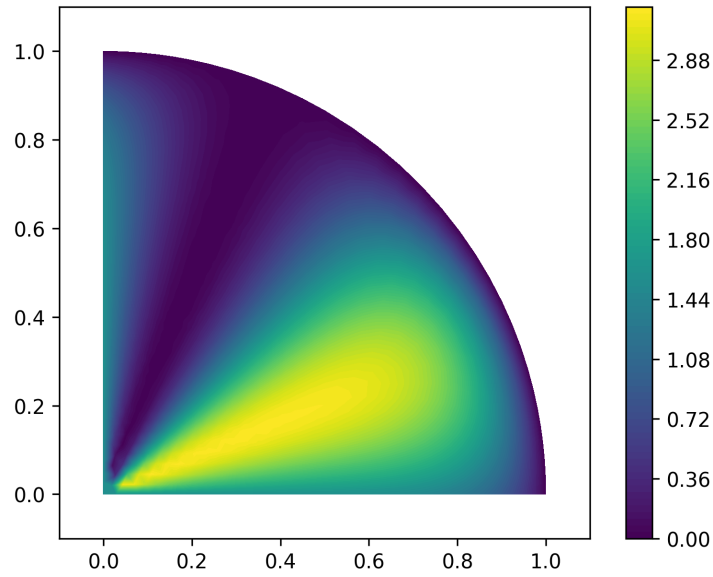


Figure 5.9: Contour plot of initial concentration (5.86)

Both  $\Phi$  and  $\eta$  are set to the constant value of 0.01 on the initial domain. The simulation is performed until  $t = 25$ , using the parameters from Tables 5.7, 5.8 and 5.9, and a timestep of size  $\Delta t = 0.25$ . A maximum initial meshwidth of  $h = 0.05\sqrt{2}$  is used. Figure 5.10 shows the domain area, total new tissue volume and total concentration as a function of time. Note that at  $t = 25$  both  $\eta_{tot}$  and  $c_{tot}$  are nearly zero. However, the domain area seems to converge to a larger value than its initial area. This is due to the inclusion of morphoelasticity by setting  $\alpha = 0.1$ . Although nearly all tumor cells have died, they have resulted in a permanent deformation in the tissue.

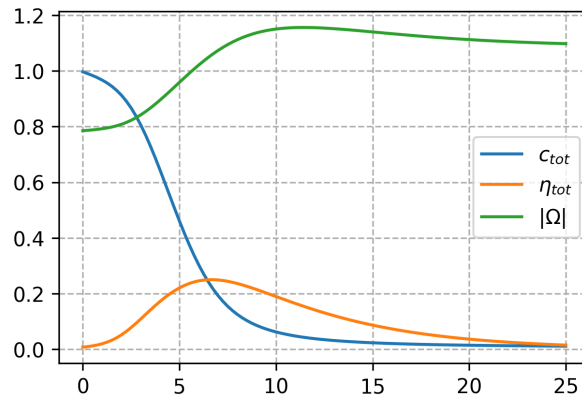


Figure 5.10: Domain area  $|\Omega|$ , total new tissue volume  $\eta_{tot}$  and total oxygen concentration  $c_{tot}$  as a function of time.

Figure 5.11 shows the strains at the halfway point of the simulation. Note that the domain is the most deformed near the  $x$ -axis. This is the result of tumor cells growing most rapidly at this position.

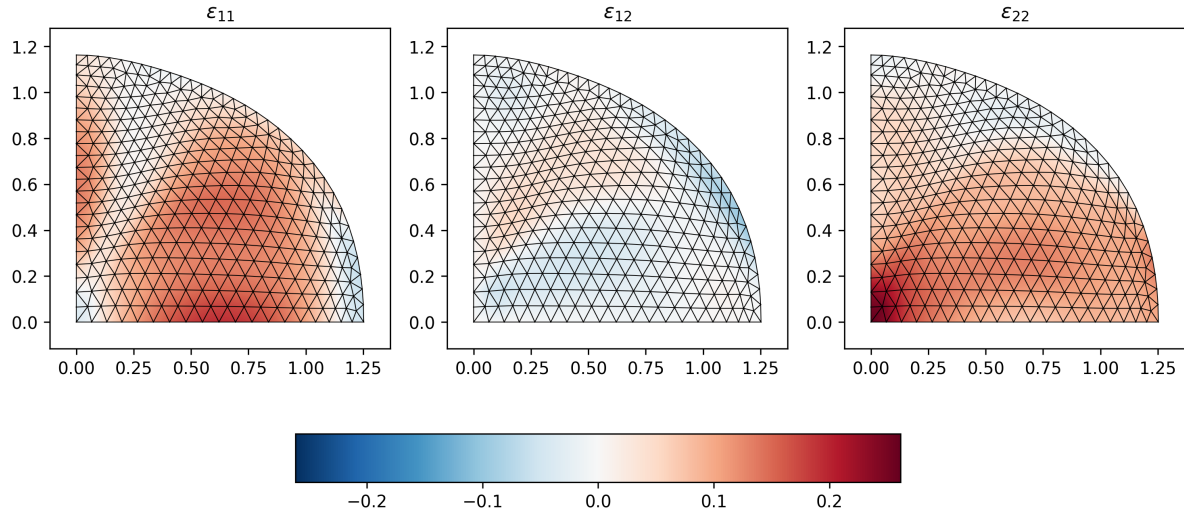


Figure 5.11: Components of the strain tensor at  $t = 12.5$  superimposed on the deformed grid.

An interesting way to visualize the growth of the domain is to consider the distribution of its radius. We view the outer boundary  $\Gamma_2(t)$  as the graph of an evolving function in polar coordinates:

$$\Gamma_2(t) = \left\{ (r(\varphi, t) \cos(\varphi), r(\varphi, t) \sin(\varphi)) \mid \varphi \in \left[0, \frac{\pi}{2}\right] \right\}. \quad (5.88)$$

Figure 5.12 shows the function  $r(\varphi, t)$  for various values of  $t$ .

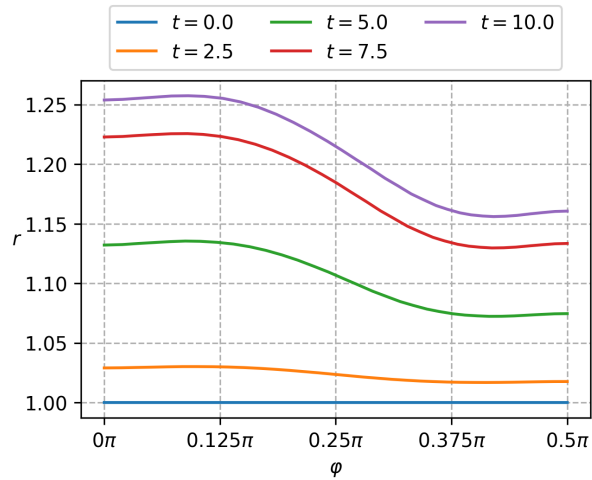


Figure 5.12: Distribution of the radius of the domain at various points in time.

Figures 5.13a, 5.13b and 5.13c show contour plots of the concentration  $c$ , new tissue volume  $\eta$  and pressure  $p$  at the halfway point of the simulation. Note that the concentration peak is now located at the top left side of the domain. At  $t = 0$ , there was a concentration dip in this area. Thus, no tumor cells have developed here and hence oxygen is absorbed at a slower rate.



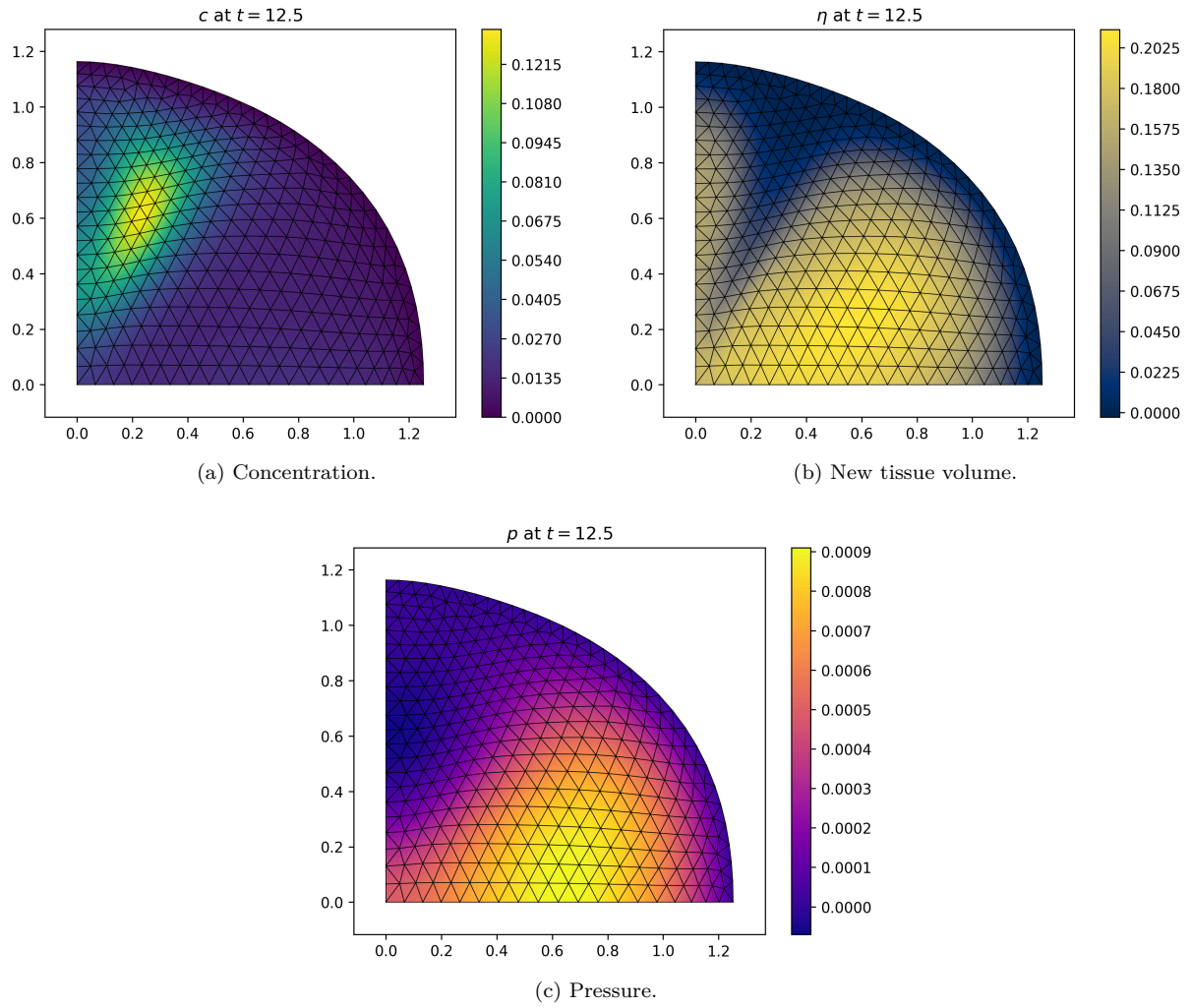


Figure 5.13: Contour plots of finite element solutions for  $c$ ,  $\eta$  and  $p$ , superimposed on the deformed grid at  $t = 12.5$ . The corresponding simulation parameters are found in Tables 5.7, 5.8 and 5.9.

The simulation parameters are presented in the following tables.

Table 5.7: Mechanical parameters.

$\rho$	$\mu_\varepsilon$	$\lambda_\varepsilon$	$\mu_v$	$\lambda_v$	$k$	$\alpha$
0.25	0.5	0.5	1	1	1	0.1

Table 5.8: Biological parameters.

$\lambda_N$	$A_{\max}$	$c_A$	$G_{\max}$	$c_G$	$r_d$
0.001	2	1	2	1	0.2

Table 5.9: Source functions and boundary conditions.

$\mathbf{g}$	$f$	$c_b$
$\mathbf{0}$	0	0

# 6

## Conclusion

In this work a finite element tumor growth model is developed, based on the combination of morpho-poroelasticity and a biochemical model. The aim of Part I of this thesis is to gain insight in the behaviour of various elastic and poroelastic models. Chapter 2 started with a two dimensional finite element model for pure elasticity. This model is then expanded to viscoelasticity, in which the elastic material is assumed to possess some viscous properties. The viscoelastic model is solved using a Lagrangian approach, in other words: a moving mesh is used. To compute and simplify the weak forms in a Lagrangian setting, Theorem 3 is used extensively throughout this work. The proof of this theorem relies on Reynold's transport theorem, and a result regarding Lagrangian basis functions. An important characteristic of the solution to a viscoelastic model is its 'overshoot' property: the domain oscillates around a stationary solution.

Next, morphoelasticity is introduced. In a morphoelastic framework the deformation tensor is decomposed in an elastic and plastic part. The plastic part causes permanent deformations in the material. The decomposition leads to an evolution equation for the strain tensor. Hence, an important difference between a visco- and morphoelastic model is that the displacement vector is no longer explicitly obtained in the latter model. To obtain the displacement vector, which is necessary to update the mesh, postprocessing is applied in every timestep. Solutions to the morphoelastic model also exhibit the 'overshoot' property, but converge to a different stationary solution due to permanent deformation in the material.

Chapter 3 introduces the concept of poroelasticity. In this framework the material is assumed to be porous and completely saturated with a poro-fluid. The presence of this fluid means that the fluid pressure must be included in the force balance. A well-known problem occurring in poroelastic modeling is the emergence of non-physical oscillations in the fluid pressure. To explore why this happens we first consider a one-dimensional simplified poroelastic model called the Terzaghi problem. We find that oscillations occur for timesteps smaller than a certain lower bound that is dependent on the order of the elements. We show that for the Terzaghi problem, these oscillations can be completely fixed by introducing a diffusive stabilization matrix.

In the two-dimensional poroelastic model these oscillations also occur for small timesteps, and for small values of the parameter  $k$ , which is the ratio between the dynamic viscosity of the fluid and the permeability of the medium. Two methods of stabilization are compared: bubble functions and diffusive stabilization. Some oscillations are still present near the domain boundaries when using bubble functions. Diffusive stabilization seems more promising, but requires the user to specify a tuning parameter. Lastly, a visco-poroelastic model is developed. Interestingly, in this framework the oscillations occur only for small values of  $k$ , and not for small timesteps.

In Part II of this thesis we show some novel results. First, the visco-poroelastic model is analyzed in order to obtain an approximately optimal tuning parameter for the diffusive stabilization. This tuning parameter is optimal in the sense that it is the smallest parameter for which the evolution matrix for the pressure is again an M-matrix. The key step of the derivation is finding an approximate inverse of the discrete version of a reaction-Laplacian operator. This is achieved by considering a closely related boundary value problem with a delta function in the right-hand side. Error bounds are proved using

functional analytic methods. The optimal tuning parameter turns out to be very similar to those found in the literature on poroelasticity.

Another novel result is the combination of morpho-and poroelasticity. We develop a finite element model that combines the permanent deformations of morphoelasticity with poroelastic behaviour. Moreover, it is shown that that when  $k$  approaches infinity, the solution to the morpho-poroelastic system tends to the solution of the morphoelastic system. However a more rigorous proof of this fact is desired. Using the same method as for the visco-poroelastic system, an approximately optimal tuning parameter can also be obtained for diffusive stabilization of the morpho-poroelastic system.

In the final chapter, the morpho-poroelastic model is combined with a biochemical tumor growth model developed by Roose et al. in [14]. This introduces three new unknowns into the system: the oxygen concentration, the tumor cell volume fraction and the new tissue volume. The oxygen concentration satisfies a convection-diffusion equation, where the convective transport is driven by the fluid pressure. The other two unknowns satisfy two very similar evolution equations with a non-linear source function. The tumor cell volume fraction takes into account the domain deformation, whereas the new tissue volume does not. We see that for slowly deforming domains, these two unknowns are almost equal. The qualitative behaviour of the oxygen concentration and tumor volume fraction is analyzed in a zero-dimensional setting. Finally, the tumor growth model is tested in a number of different testing configurations.

## 6.1 Suggestions for Further Work

In order to improve the model or further develop the mathematical theory, we put forward the following suggestions for further work:

- We suggest further research should be done to investigate the relation between the strain in a tissue and the degree of pain a patient experiences. We hypothesize that these two quantities are more or less proportional. This question introduces the difficulty of quantifying the strain and the amount of pain, the latter of which is particularly subjective.
- At a certain point in a tumor's life, it becomes so large that the surrounding tissue collapses. When this happens, tumor cells can enter the circulatory system and form colonies at other places in the body. This process is called metastasis. A possible way to formalize this 'breaking point' is to consider the work done by the tumor boundary on the surrounding tissue. Metastasis occurs whenever this energy reaches a threshold value. The finite element model can be easily adapted to compute the energy, but additional research is necessary to verify the validity of this approach.
- In Section 4.5.1 we present the argument that the solution to the morpho-poroelastic system approaches the solution to the morphoelastic system as  $k$  tends to infinity. In the proof of this claim some assumptions are made without being proved. Most notably, it is assumed that  $\mathbf{v}$  stays bounded. While this seems like a reasonable assumption since  $\mathbf{v}$  is an  $H^1$ -function, it should be proved in future work to make the argument more rigorous.
- A possible way to expand the model is to include additional nutrients. For example, it is known that next to oxygen, tumor growth is also dependent on the presence of glucose. By including glucose transport in the model the results will become more realistic. On the other hand, one could argue that oxygen and glucose are both transported by the same intercellular fluid, and that the growth rate depends on both nutrients in a similar fashion. Thus, we predict the model output will not be drastically different.
- An obvious improvement to the model would be to use higher order elements. On one hand, this will lead to a more accurate solution, but on the other hand the systems will grow in size. We have shown in Chapter 3 that Taylor-Hood elements lead to a slightly more stable solution, even without stabilization. Additionally, it could be beneficial to investigate the order of the error made by the model. Since we use linear basis functions, according to the theory of finite element methods the error should be of order  $h^2$ . However, the use of iterative methods, timestepping and possibly stabilization methods might have a negative impact on the error.
- In the numerical experiments from Section 5.4 we mostly use parameters that are close to unity. For real world applications realistic values of these parameters should be chosen. However, correct values

are often unobtainable, in which case the tumor growth model can only be used for approximate, qualitative predictions.

- A more realistic model can be constructed by considering the three-dimensional equations, which means a three-dimensional mesh has to be used. An intermediate step would be to consider a ‘pseudo’ three-dimensional model, in which we make use of rotational symmetry. In this case the model equations have to be adapted.
- In Appendix B we use delta functions to exactly invert the Laplace matrix in one dimension. A similar approach is used to approximate the inverse of the discrete Reaction-Laplacian operator in one dimension. We are interested to see if this method of finding (approximate) inverses is also applicable in two dimensions. It is expected to be more difficult because the Green’s function of the Laplacian contains a singularity.

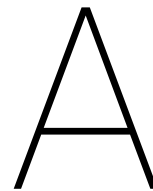
## 6.2 Suggestions for Improving Model Performance

The current implementation should be regarded as a proof of concept rather than a well-optimized numerical model. The simulations performed in Section 5.4 take only a few minutes, but the meshes are quite coarse. If the model is to be converted to three dimensions, the systems will grow from size  $\mathcal{O}(n^2)$  to  $\mathcal{O}(n^3)$ , where  $n$  is the number of meshpoints. Using the current implementation the conversion from two to three dimensions will not be computationally feasible. There are numerous ways to improve its performance. We list a number of suggestions.

- To get more accurate results for growing domains remeshing could be applied. In the current implementation elements are deformed as a result of moving grid nodes. By generating a new, higher quality mesh every few timesteps the accuracy can be improved. The current mesh generating library, `dmsh`, is rather slow. Thus, for the purpose of remeshing it should be replaced. A possible candidate is the library `MeshPy`, which provides a Python interface to the two-dimensional mesh generator `Triangle`. This package is less user-friendly but much faster. Another downside of `dmsh` is that it does not support three-dimensional meshes.
- To improve the computation time the model could be translated to C++ code. The most drastic speed-up will be seen in the matrix construction step. In the current model, all matrices are constructed by looping over every element and adding the corresponding contributions to the relevant matrix entries. Loops in Python are known to be slow compared to C or C++ loops. We suspect that the speed-up in the solving step is much smaller, because the libraries `NumPy` and `SciPy` already use C code ‘under the hood’.
- Another advantage of C++ is the ability to easily adapt the code to run on a GPU. Both the matrix construction step and (sparse) system solving step can be performed in parallel, potentially leading to drastic speed-ups.

# Bibliography

- [1] G. AGUILAR, F. GASPAR, F. LISBONA, AND C. RODRIGO, *Numerical stabilization of biot's consolidation model by a perturbation on the flow equation*, International journal for numerical methods in engineering, 75 (2008), pp. 1282–1300.
- [2] J. BELL, *Sobolev spaces in one dimension and absolutely continuous functions*, 2015.
- [3] M. BUKAC, I. YOTOV, R. ZAKERZADEH, AND P. ZUNINO, *Effects of Poroelasticity on Fluid-Structure Interaction in Arteries: a Computational Sensitivity Study*, Springer International Publishing, Cham, 2015, pp. 197–220.
- [4] R. M. CORLESS, G. H. GONNET, D. E. HARE, D. J. JEFFREY, AND D. E. KNUTH, *On the lambert w function*, Advances in Computational mathematics, 5 (1996), pp. 329–359.
- [5] S. C. COWIN, *Bone poroelasticity*, Journal of Biomechanics, 32 (1999), pp. 217 – 238.
- [6] G. DZIUK AND C. M. ELLIOTT, *Finite elements on evolving surfaces*, IMA Journal of Numerical Analysis, 27 (2007), pp. 262–292.
- [7] M. FRALDI AND A. R. CAROTENUTO, *Cells competition in tumor growth poroelasticity*, Journal of the Mechanics and Physics of Solids, 112 (2018), pp. 345 – 367.
- [8] C. L. HALL, *Modelling of some biological materials using continuum mechanics*, PhD thesis, Queensland University of Technology, 2008.
- [9] R. JIA, *MATH 418: Honors Real Variable II Lecture Notes*, University of Alberta, 2019, ch. 3: Absolutely Continuous Functions.
- [10] E. KLEIMANN, *Mathematical modelling of burn injuries*, Master's thesis, Delft University of Technology, 10 2018. Supervisor: Fred Vermolen.
- [11] D. KOPPENOL, *Biomedical implications from mathematical models for the simulation of dermal wound healing*, PhD thesis, Delft University of Technology, The address of the publisher, 6 2017.
- [12] E. H. LEE, *Elastic-plastic deformation at finite strains*, Journal of Applied Mechanics, (1969).
- [13] C. RODRIGO, F. GASPAR, X. HU, AND L. T. ZIKATANOV, *Stability and monotonicity for some discretizations of the biot's consolidation model*, Computer Methods in Applied Mechanics and Engineering, 298 (2016), pp. 183–204.
- [14] T. ROOSE, P. A. NETTI, L. L. MUNN, Y. BOUCHER, AND R. K. JAIN, *Solid stress generated by spheroid growth estimated using a linear poroelasticity model*, Microvascular research, 66 (2003), pp. 204–212.
- [15] N. SCHLÖMER, *dmsgh github*. <https://github.com/nschloe/dmsgh>. Accessed: 3/7/2020.
- [16] D. SMITS, *Morphoelastic models for burn contraction*, Master's thesis, Delft University of Technology, 9 2019. Supervisor: Fred Vermolen.
- [17] F. VERMOLEN AND A. SEGAL, *On an integration rule for products of barycentric coordinates of simplexes in  $\mathbb{R}^n$* , Journal of Computational and Applied Mathematics, 330 (2018), pp. 289–294.
- [18] C. VUIK, F. VERMOLEN, M. VAN GIJZEN, AND M. VUIK, *Numerical Methods for Ordinary Differential Equations*, Delft Academic Press, 2 ed., 2016.
- [19] R. A. WEINBERG AND R. A. WEINBERG, *The biology of cancer*, Garland science, 2013.



## Element Integrals

In this chapter, various integrals that appear in the Galerkin equations are computed over a triangular element. We consider the element  $T$  with vertices  $\mathbf{x}_1, \mathbf{x}_2$  and  $\mathbf{x}_3$ . We use linear basis functions and triangular elements, and write

$$\varphi_i(x, y) = \alpha_i x + \beta_i y + \gamma_i. \quad (\text{A.1})$$

for  $i \in \{1, 2, 3\}$ . Thus, whenever the partial derivatives of a basis function appear, we substitute

$$\frac{\partial \varphi_i}{\partial x} = \alpha_i, \quad \text{and} \quad \frac{\partial \varphi_i}{\partial y} = \beta_i. \quad (\text{A.2})$$

The  $\alpha_i$  and  $\beta_i$  can be computed using the relation  $\varphi_i(\mathbf{x}_j) = 1$ , leading to the linear system

$$\begin{pmatrix} x_1 & y_1 & 1 \\ x_2 & y_2 & 1 \\ x_3 & y_3 & 1 \end{pmatrix} \begin{pmatrix} \alpha_i \\ \beta_i \\ \gamma_i \end{pmatrix} = \mathbf{e}_i, \quad (\text{A.3})$$

where  $\mathbf{e}_i$  contains a 1 in its  $i$ 'th position and zeroes otherwise. The area of the triangle  $T$  can also be expressed in terms of the matrix above:

$$|T| = \frac{1}{2} \det \begin{pmatrix} x_1 & y_1 & 1 \\ x_2 & y_2 & 1 \\ x_3 & y_3 & 1 \end{pmatrix}. \quad (\text{A.4})$$

This follows from a well-known theorem by Holand & Bell, see for example [17]. In this chapter we will give the element matrices and vectors corresponding to matrices and vectors found in the main text of this work. For example, let us consider a mass matrix, often defined as

$$M_{ij} = \int_{\Omega} \varphi_i \varphi_j \, d\Omega. \quad (\text{A.5})$$

To construct this matrix, we look at the contribution each element in the triangulation makes per entry:

$$M_{ij} = \sum_{T \in \mathcal{T}} \int_T \varphi_i \varphi_j \, d\Omega. \quad (\text{A.6})$$

Note that we can also view the summand in the right-hand side as an  $n \times n$  matrix. However, since the linear basis functions  $\varphi_i$  and  $\varphi_j$  are only non-zero on elements containing grid node  $i$  or  $j$  respectively, this matrix will have 9 non-zero entries. These entries correspond to the integrals in which both grid node  $i$  and  $j$  are vertices of  $T$ . When constructing  $M$ , it is therefore sufficient to compute for each element  $T$  the  $3 \times 3$  matrix

$$(M_T)_{ij} = \int_T \varphi_{i,T} \varphi_{j,T} \, d\Omega, \quad (\text{A.7})$$

for  $i, j \in \{1, 2, 3\}$ . Here we use  $(i, T)$  to denote the  $i$ 'th vertex of  $T$ . In the following sections, we simply write  $i$  instead of  $(i, T)$ , since it is clear from the context that only one element and its three vertices are considered.

## A.1 Pure Elasticity

We give the element matrices and vectors corresponding to the objects found in system (2.25) from Section 2.2.2. The  $3 \times 3$  element block matrices  $S_T^{11}$ ,  $S_T^{12}$ ,  $S_T^{21}$  and  $S_T^{22}$  are given by

$$(S_T^{11})_{ij} = \int_T (\mu + \lambda) \frac{\partial \varphi_i}{\partial x} \frac{\partial \varphi_j}{\partial x} + \frac{1}{2} \mu \frac{\partial \varphi_i}{\partial y} \frac{\partial \varphi_j}{\partial y} \, d\Omega = |T| \left( (\mu + \lambda) \alpha_i \alpha_j + \frac{1}{2} \mu \beta_i \beta_j \right), \quad (\text{A.8})$$

$$(S_T^{12})_{ij} = \int_T \lambda \frac{\partial \varphi_i}{\partial x} \frac{\partial \varphi_j}{\partial y} + \frac{1}{2} \mu \frac{\partial \varphi_i}{\partial y} \frac{\partial \varphi_j}{\partial x} \, d\Omega = |T| \left( \lambda \alpha_i \beta_j + \frac{1}{2} \mu \beta_i \alpha_j \right), \quad (\text{A.9})$$

$$(S_T^{21})_{ij} = \int_T \frac{1}{2} \mu \frac{\partial \varphi_i}{\partial x} \frac{\partial \varphi_j}{\partial y} + \lambda \frac{\partial \varphi_i}{\partial y} \frac{\partial \varphi_j}{\partial x} \, d\Omega = |T| \left( \frac{1}{2} \mu \alpha_i \beta_j + \lambda \beta_i \alpha_j \right), \quad (\text{A.10})$$

$$(S_T^{22})_{ij} = \int_T \frac{1}{2} \mu \frac{\partial \varphi_i}{\partial x} \frac{\partial \varphi_j}{\partial x} + (\mu + \lambda) \frac{\partial \varphi_i}{\partial y} \frac{\partial \varphi_j}{\partial y} \, d\Omega = |T| \left( \frac{1}{2} \mu \alpha_i \alpha_j + (\mu + \lambda) \beta_i \beta_j \right), \quad (\text{A.11})$$

Using the outer product of two vectors, the element matrices can be written in the following form:

$$S_T^{11} = |T| \left( (\mu + \lambda) \boldsymbol{\alpha} \boldsymbol{\alpha}^\top + \frac{1}{2} \mu \boldsymbol{\beta} \boldsymbol{\beta}^\top \right), \quad (\text{A.12})$$

$$S_T^{12} = |T| \left( \lambda \boldsymbol{\alpha} \boldsymbol{\beta}^\top + \frac{1}{2} \mu \boldsymbol{\beta} \boldsymbol{\alpha}^\top \right), \quad (\text{A.13})$$

$$S_T^{21} = |T| \left( \frac{1}{2} \mu \boldsymbol{\alpha} \boldsymbol{\beta}^\top + \lambda \boldsymbol{\beta} \boldsymbol{\alpha}^\top \right), \quad (\text{A.14})$$

$$S_T^{22} = |T| \left( \frac{1}{2} \mu \boldsymbol{\alpha} \boldsymbol{\alpha}^\top + (\mu + \lambda) \boldsymbol{\beta} \boldsymbol{\beta}^\top \right). \quad (\text{A.15})$$

Where  $\boldsymbol{\alpha}$  and  $\boldsymbol{\beta}$  are the column vectors

$$\boldsymbol{\alpha} = \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{bmatrix}, \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix}. \quad (\text{A.16})$$

Expressions (A.8)-(A.11) can be used to construct (2.30)-(2.33) by summing over all elements. The entries of the right-hand side vectors  $\mathbf{g}$  and  $\boldsymbol{\tau}$  can be approximated on the triangle  $T$  using a first order Newton-Cotes rule:

$$(g_T^1)_i = \int_T \varphi_i g^1 \, d\Omega \approx \frac{|T|}{3} g^1(\mathbf{x}_i), \quad i = 1, 2, 3. \quad (\text{A.17})$$

And similarly for  $(g_T^2)_i$ . Let  $T_B$  be the boundary element (line) connecting the mesh points  $\mathbf{x}_1$  and  $\mathbf{x}_2$ . The boundary element vector component  $(\tau_{T_B}^1)_i$  is approximated by

$$(\tau_{T_B}^1)_i = \int_{T_B} \varphi_i \tau^1 \, d\Gamma \approx \frac{|T_B|}{2} \tau^1(\mathbf{x}_i), \quad i = 1, 2. \quad (\text{A.18})$$

Where  $|T_B|$  is the length of  $T_B$ . The complete vector can be found by summing the (boundary) element vector components over all elements and boundary elements respectively. For example, we have

$$(g^1)_i = \sum_{T \in \mathcal{T}} (g_T^1)_i, \quad \text{and} \quad (\tau^1)_i = \sum_{T_B \in \mathcal{T}_B} (\tau_{T_B}^1)_i, \quad (\text{A.19})$$

where  $\mathcal{T}$  denotes a triangulation of  $\Omega$  and  $\mathcal{T}_B$  is the corresponding partition of  $\Gamma_1$ .

## A.2 Morphoelasticity

We give all element matrices and vectors relative to  $T$  of the objects that make up the system in (2.109) and (2.119). First, consider the mass matrix:

$$M_T = \frac{|T|}{12} \begin{bmatrix} 2 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \end{bmatrix}, \quad (\text{A.20})$$

Next, we consider the blocks that make up the viscous stress matrix. Note that the elastic stress matrix  $S_{el}$  found in sections 2.2 and 2.3 is very similar. Compared to the the following expressions, the parameters  $\mu_v, \lambda_v$  should be replaced by  $\mu_\varepsilon, \lambda_\varepsilon$ .

$$(S_{vis}^{11})_T = |T|(\mu_v + \lambda_v) \begin{bmatrix} \alpha_1\alpha_1 & \alpha_1\alpha_2 & \alpha_1\alpha_3 \\ \alpha_2\alpha_1 & \alpha_2\alpha_2 & \alpha_2\alpha_3 \\ \alpha_3\alpha_1 & \alpha_3\alpha_2 & \alpha_3\alpha_3 \end{bmatrix} + \frac{1}{2}|T|\mu_v \begin{bmatrix} \beta_1\beta_1 & \beta_1\beta_2 & \beta_1\beta_3 \\ \beta_2\beta_1 & \beta_2\beta_2 & \beta_2\beta_3 \\ \beta_3\beta_1 & \beta_3\beta_2 & \beta_3\beta_3 \end{bmatrix}, \quad (\text{A.21})$$

$$(S_{vis}^{12})_T = |T|\lambda_v \begin{bmatrix} \alpha_1\beta_1 & \alpha_1\beta_2 & \alpha_1\beta_3 \\ \alpha_2\beta_1 & \alpha_2\beta_2 & \alpha_2\beta_3 \\ \alpha_3\beta_1 & \alpha_3\beta_2 & \alpha_3\beta_3 \end{bmatrix} + \frac{1}{2}|T|\mu_v \begin{bmatrix} \beta_1\alpha_1 & \beta_1\alpha_2 & \beta_1\alpha_3 \\ \beta_2\alpha_1 & \beta_2\alpha_2 & \beta_2\alpha_3 \\ \beta_3\alpha_1 & \beta_3\alpha_2 & \beta_3\alpha_3 \end{bmatrix}, \quad (\text{A.22})$$

$$(S_{vis}^{21})_T = \frac{1}{2}|T|\mu_v \begin{bmatrix} \alpha_1\beta_1 & \alpha_1\beta_2 & \alpha_1\beta_3 \\ \alpha_2\beta_1 & \alpha_2\beta_2 & \alpha_2\beta_3 \\ \alpha_3\beta_1 & \alpha_3\beta_2 & \alpha_3\beta_3 \end{bmatrix} + |T|\lambda_v \begin{bmatrix} \beta_1\alpha_1 & \beta_1\alpha_2 & \beta_1\alpha_3 \\ \beta_2\alpha_1 & \beta_2\alpha_2 & \beta_2\alpha_3 \\ \beta_3\alpha_1 & \beta_3\alpha_2 & \beta_3\alpha_3 \end{bmatrix}, \quad (\text{A.23})$$

$$(S_{vis}^{22})_T = \frac{1}{2}|T|\mu_v \begin{bmatrix} \alpha_1\alpha_1 & \alpha_1\alpha_2 & \alpha_1\alpha_3 \\ \alpha_2\alpha_1 & \alpha_2\alpha_2 & \alpha_2\alpha_3 \\ \alpha_3\alpha_1 & \alpha_3\alpha_2 & \alpha_3\alpha_3 \end{bmatrix} + |T|(\lambda_v + \mu_v) \begin{bmatrix} \beta_1\beta_1 & \beta_1\beta_2 & \beta_1\beta_3 \\ \beta_2\beta_1 & \beta_2\beta_2 & \beta_2\beta_3 \\ \beta_3\beta_1 & \beta_3\beta_2 & \beta_3\beta_3 \end{bmatrix}, \quad (\text{A.24})$$

The block matrices  $B_x$  and  $B_y$  frequently occur in the system, its element matrices are given by:

$$(B_x)_T = \frac{|T|}{3} \begin{bmatrix} \alpha_1 & \alpha_1 & \alpha_1 \\ \alpha_2 & \alpha_2 & \alpha_2 \\ \alpha_3 & \alpha_3 & \alpha_3 \end{bmatrix}, \quad (B_y)_T = \frac{|T|}{3} \begin{bmatrix} \beta_1 & \beta_1 & \beta_1 \\ \beta_2 & \beta_2 & \beta_2 \\ \beta_3 & \beta_3 & \beta_3 \end{bmatrix}. \quad (\text{A.25})$$

Furthermore, the element vectors are given by:

$$(g^1)_T = \frac{|T|}{3} \begin{bmatrix} g^1(\mathbf{x}_1) \\ g^1(\mathbf{x}_2) \\ g^1(\mathbf{x}_3) \end{bmatrix}, \quad (g^2)_T = \frac{|T|}{3} \begin{bmatrix} g^2(\mathbf{x}_1) \\ g^2(\mathbf{x}_2) \\ g^2(\mathbf{x}_3) \end{bmatrix}, \quad (\text{A.26})$$

$$(G^{11})_T = \frac{|T|}{3} \begin{bmatrix} G^{11}(\mathbf{x}_1) \\ G^{11}(\mathbf{x}_2) \\ G^{11}(\mathbf{x}_3) \end{bmatrix}, \quad (G^{21})_T = \frac{|T|}{3} \begin{bmatrix} G^{12}(\mathbf{x}_1) \\ G^{12}(\mathbf{x}_2) \\ G^{12}(\mathbf{x}_3) \end{bmatrix}, \quad (G^{22})_T = \frac{|T|}{3} \begin{bmatrix} G^{22}(\mathbf{x}_1) \\ G^{22}(\mathbf{x}_2) \\ G^{22}(\mathbf{x}_3) \end{bmatrix}. \quad (\text{A.27})$$

Let us now consider the element matrices involving the non-linear part of (2.109). Note that we can view  $(N^x)_T$  and  $(N^y)_T$  as  $3 \times 3 \times 3$  tensors that we can write as

$$((N^x)_T)_{ijk} = \alpha_k \frac{|T|}{3} (1 + \delta_{ij}), \quad ((N^y)_T)_{ijk} = \beta_k \frac{|T|}{3} (1 + \delta_{ij}). \quad (\text{A.28})$$

Thus, we get the following element matrices:

$$((N^x)_T)_{1..} = \frac{|T|}{3} \begin{bmatrix} 2\alpha_1 & 2\alpha_2 & 2\alpha_3 \\ \alpha_1 & \alpha_2 & \alpha_3 \\ \alpha_1 & \alpha_2 & \alpha_3 \end{bmatrix}, \quad (\text{A.29})$$

$$((N^x)_T)_{2..} = \frac{|T|}{3} \begin{bmatrix} \alpha_1 & \alpha_2 & \alpha_3 \\ 2\alpha_1 & 2\alpha_2 & 2\alpha_3 \\ \alpha_1 & \alpha_2 & \alpha_3 \end{bmatrix}, \quad (\text{A.30})$$

$$((N^x)_T)_{3..} = \frac{|T|}{3} \begin{bmatrix} \alpha_1 & \alpha_2 & \alpha_3 \\ \alpha_1 & \alpha_2 & \alpha_3 \\ 2\alpha_1 & 2\alpha_2 & 2\alpha_3 \end{bmatrix}, \quad (\text{A.31})$$

and

$$((N^y)_T)_{1..} = \frac{|T|}{3} \begin{bmatrix} 2\beta_1 & 2\beta_2 & 2\beta_3 \\ \beta_1 & \beta_2 & \beta_3 \\ \beta_1 & \beta_2 & \beta_3 \end{bmatrix}, \quad (\text{A.32})$$

$$((N^y)_T)_{2..} = \frac{|T|}{3} \begin{bmatrix} \beta_1 & \beta_2 & \beta_3 \\ 2\beta_1 & 2\beta_2 & 2\beta_3 \\ \beta_1 & \beta_2 & \beta_3 \end{bmatrix}, \quad (\text{A.33})$$

$$((N^y)_T)_{3..} = \frac{|T|}{3} \begin{bmatrix} \beta_1 & \beta_2 & \beta_3 \\ \beta_1 & \beta_2 & \beta_3 \\ 2\beta_1 & 2\beta_2 & 2\beta_3 \end{bmatrix}. \quad (\text{A.34})$$



Finally, let  $T_B$  be the boundary element connecting vertices  $\mathbf{x}_1$  and  $\mathbf{x}_2$ , then:

$$(\tau^1)_{T_B} = \frac{|T_B|}{2} \begin{bmatrix} \tau^1(\mathbf{x}_1) \\ \tau^1(\mathbf{x}_2) \end{bmatrix}, \quad (\tau^2)_{T_B} = \frac{|T_B|}{2} \begin{bmatrix} \tau^2(\mathbf{x}_1) \\ \tau^2(\mathbf{x}_2) \end{bmatrix}, \quad (\text{A.35})$$

where  $|T_B|$  denotes the length of  $T_B$ .

### A.3 Bubble Functions

We give the element matrices and vector corresponding to bubble matrices  $S_b$ ,  $D_b$  and bubble vector  $\mathbf{g}_b$  from Section 3.3.2. Let  $\varphi_{b,k}$  be the bubble function corresponding to element  $T_k$ , where  $k \in \{1, \dots, n_T\}$ .

$$(S_b^{11})_{kk} = (\mu + \lambda) \int_{T_k} \frac{\partial \varphi_{b,k}}{\partial x} \frac{\partial \varphi_{b,k}}{\partial x} d\Omega + \frac{1}{2} \mu \int_{T_k} \frac{\partial \varphi_{b,k}}{\partial y} \frac{\partial \varphi_{b,k}}{\partial y} d\Omega \quad (\text{A.36})$$

$$(S_b^{12})_{kk} = \lambda \int_{T_k} \frac{\partial \varphi_{b,k}}{\partial x} \frac{\partial \varphi_{b,k}}{\partial y} d\Omega + \frac{1}{2} \mu \int_{T_k} \frac{\partial \varphi_{b,k}}{\partial y} \frac{\partial \varphi_{b,k}}{\partial x} d\Omega \quad (\text{A.37})$$

$$(S_b^{21})_{kk} = \frac{1}{2} \mu \int_{T_k} \frac{\partial \varphi_{b,k}}{\partial x} \frac{\partial \varphi_{b,k}}{\partial y} d\Omega + \lambda \int_{T_k} \frac{\partial \varphi_{b,k}}{\partial y} \frac{\partial \varphi_{b,k}}{\partial x} d\Omega \quad (\text{A.38})$$

$$(S_b^{22})_{kk} = \frac{1}{2} \mu \int_{T_k} \frac{\partial \varphi_{b,k}}{\partial x} \frac{\partial \varphi_{b,k}}{\partial x} d\Omega + (\mu + \lambda) \int_{T_k} \frac{\partial \varphi_{b,k}}{\partial y} \frac{\partial \varphi_{b,k}}{\partial y} d\Omega \quad (\text{A.39})$$

Note that all the  $S_b$  blocks are diagonal matrices, because a bubble-function is only non-zero on its corresponding element. This is also why can integrate over  $T_k$  instead of  $\Omega_h$ . The three distinct integrals in the above formulas are given by:

$$\int_{T_k} \frac{\partial \varphi_{b,k}}{\partial x} \frac{\partial \varphi_{b,k}}{\partial x} d\Omega = \frac{81|T_k|}{10} [\alpha_1^2 + \alpha_1\alpha_2 + \alpha_2^2 + \alpha_2\alpha_3 + \alpha_3^2 + \alpha_1\alpha_3] \quad (\text{A.40})$$

$$\int_{T_k} \frac{\partial \varphi_{b,k}}{\partial x} \frac{\partial \varphi_{b,k}}{\partial y} d\Omega = \frac{81|T_k|}{20} [2\alpha_1\beta_1 + \alpha_1\beta_2 + \alpha_1\beta_3 + \alpha_2\beta_1 + 2\alpha_2\beta_2 + \alpha_2\beta_3 + \alpha_3\beta_1 + \alpha_3\beta_2 + 2\alpha_3\beta_3] \quad (\text{A.41})$$

$$\int_{T_k} \frac{\partial \varphi_{b,k}}{\partial y} \frac{\partial \varphi_{b,k}}{\partial y} d\Omega = \frac{81|T_k|}{10} [\beta_1^2 + \beta_1\beta_2 + \beta_2^2 + \beta_2\beta_3 + \beta_3^2 + \beta_1\beta_3] \quad (\text{A.42})$$

The matrix  $D_b$  is also subdivided into blocks:  $D_b = (D_b^1 \quad D_b^2)$ . They are given component-wise by:

$$(D_b^1)_{kj} = \int_{\Omega_h} \frac{\partial \varphi_{b,k}}{\partial x} \psi_j d\Omega, \quad (D_b^2)_{kj} = \int_{\Omega_h} \frac{\partial \varphi_{b,k}}{\partial y} \psi_j d\Omega, \quad (\text{A.43})$$

where  $k \in \{1, \dots, n_T\}$  and  $j \in \{1, \dots, n\}$  and  $\psi_j$  is a linear basis function. The corresponding element matrices are given by:

$$(D_b^1)_T = \frac{9|T|}{20} \begin{bmatrix} \alpha_1 + 2\alpha_2 + 2\alpha_3 \\ 2\alpha_1 + \alpha_2 + 2\alpha_3 \\ 2\alpha_1 + 2\alpha_2 + \alpha_3 \end{bmatrix}, \quad (D_b^2)_T = \frac{9|T|}{20} \begin{bmatrix} \beta_1 + 2\beta_2 + 2\beta_3 \\ 2\beta_1 + \beta_2 + 2\beta_3 \\ 2\beta_1 + 2\beta_2 + \beta_3 \end{bmatrix}. \quad (\text{A.44})$$

We also write  $\mathbf{g}_b^m = ((\mathbf{g}_b^1)^m, (\mathbf{g}_b^2)^m)^\top$ , they are given by:

$$(\mathbf{g}_b^1)_{T_k}^m = \int_{T_k} \varphi_{b,k} (g^1)^m d\Omega, \quad (\mathbf{g}_b^2)_{T_k}^m = \int_{T_k} \varphi_{b,k} (g^2)^m d\Omega. \quad (\text{A.45})$$

The integral is approximated using Newton-Cotes quadrature, where we use the barycentric coordinates of  $T_k$  to linearly approximate  $g^1$  and  $g^2$ . Holand & Bell's theorem can then be used to exactly compute the resulting integrals, such that we obtain:

$$(\mathbf{g}_b^1)_{T_k}^m \approx \frac{3|T_k|}{20} (g^1(\mathbf{x}_1, t^m) + g^1(\mathbf{x}_2, t^m) + g^1(\mathbf{x}_3, t^m)), \quad (\text{A.46})$$

$$(\mathbf{g}_b^2)_{T_k}^m \approx \frac{3|T_k|}{20} (g^2(\mathbf{x}_1, t^m) + g^2(\mathbf{x}_2, t^m) + g^2(\mathbf{x}_3, t^m)). \quad (\text{A.47})$$

## A.4 Tumor Growth Model

Although most integrals found in the tumor growth model have already been calculated in previous sections, a handful of new ones are introduced. Note that the last integral in (5.56) is a kind of mass matrix with an additional basis function, we denote it by:

$$((\mathcal{M}_\ell)_T)_{ij} = \int_T \varphi_{\ell,T} \varphi_i \varphi_j \, d\Omega. \quad (\text{A.48})$$

Its element matrices are given by:

$$(\mathcal{M}_1)_T = \frac{|T|}{60} \begin{bmatrix} 6 & 2 & 2 \\ 2 & 2 & 1 \\ 2 & 1 & 2 \end{bmatrix}, \quad (\mathcal{M}_2)_T = \frac{|T|}{60} \begin{bmatrix} 2 & 2 & 1 \\ 2 & 6 & 2 \\ 1 & 2 & 2 \end{bmatrix}, \quad (\mathcal{M}_3)_T = \frac{|T|}{60} \begin{bmatrix} 2 & 1 & 2 \\ 1 & 2 & 2 \\ 2 & 2 & 6 \end{bmatrix}. \quad (\text{A.49})$$

The above integral also occurs in (5.66). The other integral in this equation is similar, but contains yet another basis function. We denote it by

$$((\mathcal{M}_{k\ell})_T)_{ij} = \int_T \varphi_i \varphi_j \varphi_{k,T} \varphi_{\ell,T} \, d\Omega. \quad (\text{A.50})$$

Its nine corresponding element matrices are given by

$$\begin{aligned} (\mathcal{M}_{11})_T &= \frac{|T|}{180} \begin{bmatrix} 12 & 3 & 3 \\ 3 & 2 & 1 \\ 3 & 1 & 2 \end{bmatrix} & (\mathcal{M}_{12})_T &= \frac{|T|}{180} \begin{bmatrix} 3 & 2 & 1 \\ 2 & 3 & 1 \\ 1 & 1 & 1 \end{bmatrix} & (\mathcal{M}_{13})_T &= \frac{|T|}{180} \begin{bmatrix} 3 & 1 & 2 \\ 1 & 1 & 1 \\ 2 & 1 & 3 \end{bmatrix} \\ (\mathcal{M}_{21})_T &= \frac{|T|}{180} \begin{bmatrix} 3 & 2 & 1 \\ 2 & 3 & 1 \\ 1 & 1 & 1 \end{bmatrix} & (\mathcal{M}_{22})_T &= \frac{|T|}{180} \begin{bmatrix} 2 & 3 & 1 \\ 3 & 12 & 3 \\ 1 & 3 & 2 \end{bmatrix} & (\mathcal{M}_{23})_T &= \frac{|T|}{180} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 3 & 2 \\ 1 & 2 & 3 \end{bmatrix} \\ (\mathcal{M}_{31})_T &= \frac{|T|}{180} \begin{bmatrix} 3 & 1 & 2 \\ 1 & 1 & 1 \\ 2 & 1 & 3 \end{bmatrix} & (\mathcal{M}_{32})_T &= \frac{|T|}{180} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 3 & 2 \\ 1 & 2 & 3 \end{bmatrix} & (\mathcal{M}_{33})_T &= \frac{|T|}{180} \begin{bmatrix} 2 & 1 & 3 \\ 1 & 2 & 3 \\ 3 & 3 & 12 \end{bmatrix} \end{aligned}$$

# B

## Matrix Inversion by Analyzing Delta Problems

### B.1 Inverting the Laplace Matrix in 1D

In section 4.3.3 we computed an approximation to the inverse of a discrete differential operator. It turns out that for the Laplace operator, this approximation is exact. We analyze the following equation:

$$-u'' = \delta(x - a), \quad x \in (0, 1) \quad (\text{B.1})$$

where  $a \in (0, 1)$ . A number of different sets of boundary conditions can be imposed. We will first consider homogeneous Dirichlet boundary conditions:

$$u(0) = u(1) = 0. \quad (\text{B.2})$$

Later we will also consider different sets of boundary conditions. The exact (weak) solution to (B.1)-(B.2) is given by:

$$u(x) = (1 - a)x - (x - a)H(x - a), \quad (\text{B.3})$$

where  $H(\cdot)$  is the Heaviside step-function. The exact solution can be verified by integrating the differential equation twice, and substituting the boundary conditions to obtain the values of the integration constants. Another option would be to use the Laplace transform. Note that  $u$  is not a classical solution, but we do have  $u \in H^1(0, 1)$ . We build a finite element model to approximate the solution to (B.1). First, consider the weak form:

Find  $u \in H_0^1(0, 1)$  such that

$$\int_0^1 u' \varphi' \, dx = \varphi(a), \quad \forall \varphi \in H_0^1(0, 1), \quad (\text{B.4})$$

where  $H_0^1(0, 1) = \{u \in H^1(0, 1) : u(0) = u(1) = 0\}$ .

Divide  $(0, 1)$  into a uniform mesh consisting of  $n + 1$  elements. Let  $h = 1/(n + 1)$  be the meshwidth and let  $x_i = ih$  be a grid point for  $i = 0, \dots, n + 1$ . Note that since  $u(0) = u(1) = 0$ , the unknowns are located at  $x_1, \dots, x_n$ . Let  $\varphi_i$  be a piecewise linear basis function such that  $\varphi_i(x_j) = \delta_{ij}$ , and let  $u_h$  be the finite element approximation to  $u$ . We write

$$u_h(x) = \sum_{j=1}^n c_j \varphi_j(x). \quad (\text{B.5})$$

Substitute  $u_h$  into the weak form (B.4), and set  $\varphi = \varphi_i$  to obtain

$$\sum_{j=1}^n c_j \int_0^1 \varphi_i' \varphi_j' \, dx = \varphi_i(a), \quad i = 1, \dots, n. \quad (\text{B.6})$$

We let  $a$  coincide with one of the grid points, say  $a = x_k$  for some  $k \in \{1, \dots, n\}$ . Then (B.6) turns into:

$$\sum_{j=1}^n c_j \int_0^1 \varphi'_i \varphi'_j \, dx = \delta_{ik}, \quad i = 1, \dots, n. \quad (\text{B.7})$$

Let

$$S_{ij} = \int_0^1 \varphi'_i \varphi'_j \, dx, \quad (\text{B.8})$$

then the matrix  $S$  is given by

$$S = \frac{1}{h} \begin{pmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 & -1 \\ & & & -1 & 2 \end{pmatrix}. \quad (\text{B.9})$$

**Theorem 12.** *Let the  $n \times n$  matrix  $S$  be given by (B.9), and let  $\mathbf{e}_k$  be the  $k$ 'th standard basis vector in  $\mathbb{R}^n$ , such that  $e_k^j = \delta_{kj}$ . Then the vector  $\mathbf{v}_k \in \mathbb{R}^n$  given by*

$$v_k^j = (1 - x_k)x_j - (x_j - x_k)H(j - k) \quad (\text{B.10})$$

*solves the equation  $S\mathbf{v}_k = \mathbf{e}_k$  for  $k = 1, \dots, n$ . Therefore, the entries of  $S^{-1}$  are defined by (B.10).*

*Proof.* We use Céa's Lemma to prove the theorem. Let  $u$  and  $u_h$  be the exact solution and finite element solution respectively. Assume that  $u_h$  satisfies (B.7). Then by Céa's Lemma we have

$$\|u' - u'_h\|_{L^2(0,1)} \leq \|u' - v'\|_{L^2(0,1)}$$

for all  $v \in \text{Span}\{\varphi_i\}_{i=1}^n$ . Note that  $u$  is linear on each element. Thus, by choosing for  $v$  the linear interpolant of  $u$ , the error becomes zero. Then, by the Poincaré inequality we obtain

$$\|u - u_h\|_{L^2(0,1)} \leq c \|u' - u'_h\|_{L^2(0,1)} \leq 0.$$

Because the finite element error is zero, expression (B.3) can be used to construct the finite element solution. Thus, by setting  $a = x_k$  in (B.3), we find that the solution to  $S\mathbf{v}_k = \mathbf{e}_k$  is given by:

$$v_k^j = (1 - x_k)x_j - (x_j - x_k)H(j - k).$$

Consider the matrix equation

$$S(\mathbf{v}_1 \quad \dots \quad \mathbf{v}_n) = (\mathbf{e}_1 \quad \dots \quad \mathbf{e}_n) = I,$$

then it is clear that

$$S^{-1} = (\mathbf{v}_1 \quad \dots \quad \mathbf{v}_n).$$

□

Since we have found an expression for the inverse of  $S$ , general systems of the form

$$S\mathbf{u} = \mathbf{f} \quad (\text{B.11})$$

can now easily be solved by setting

$$\mathbf{u} = \sum_{k=1}^n f_k \mathbf{v}_k. \quad (\text{B.12})$$

## B.2 Connection to Green's Functions

It is well known that for any linear differential operator  $L$ , the problem

$$Lu(x) = f(x) \quad (\text{B.13})$$

can be solved exactly by using the Green's function corresponding to  $L$ . The Green's function  $G(x, s)$  satisfies the equation

$$LG(x, s) = \delta(x - s). \quad (\text{B.14})$$

Then the solution to (B.13) is given by

$$u(x) = \int G(x, s) f(s) \, ds. \quad (\text{B.15})$$

This can be easily verified by noting that  $L$  is linear and only acts on the variable  $x$ , thus:

$$\begin{aligned} Lu(x) &= L \int G(x, s) f(s) \, ds \\ &= \int LG(x, s) f(s) \, ds \\ &= \int \delta(x - s) f(s) \, ds \\ &= f(x). \end{aligned}$$

Note that in our case  $L = -\frac{d^2}{dx^2}$  with boundary conditions (B.2). The corresponding Green's function is given by (B.3). We rewrite it to fit the notation of this section:

$$G(x, s) = (1 - s)x - (x - s)H(x - s). \quad (\text{B.16})$$

By the previous section, we know that  $G(x, x_k)$  is approximated exactly in a grid point  $x_k$  by the solution to  $S\mathbf{v}_k = \mathbf{e}_k$ , where  $S$  is given in (B.9). This fact can be used to prove that the finite element scheme, as introduced in the previous section, also approximates the solution to arbitrary problems of the form  $-u''(x) = f(x)$  exactly.

**Theorem 13.** *Assume that  $u$  satisfies  $-u''(x) = f(x)$  and the boundary conditions (B.2). Let  $u_h$  be the corresponding finite element solution using linear basis functions on a uniform grid of  $n + 1$  elements. We write*

$$u_h(x) = \sum_{j=1}^n u_j \varphi_j(x). \quad (\text{B.17})$$

*Let the vector  $\mathbf{u}$  contain the coefficients  $u_j$ . Then it satisfies the equation  $S\mathbf{u} = \mathbf{f}$ , where the matrix  $S$  is given by (B.9) and  $\mathbf{f}$  is defined pointwise by*

$$f_i = \int_0^1 f(x) \varphi_i(x) \, dx. \quad (\text{B.18})$$

*Then  $u_h(x_i) = u(x_i)$  in each grid point  $x_i$ .*

For general functions  $f$ , the integrals in the right-hand side of (B.18) are usually approximated using some sort of Newton-Cotes quadrature rule. For the proof to hold, we need the elements of  $\mathbf{f}$  to be exact.

*Proof.* Since  $\mathbf{u}$  satisfies  $S\mathbf{u} = \mathbf{f}$ , by (B.11) and (B.12) we can write

$$\mathbf{u} = \sum_{k=1}^n f_k \mathbf{v}_k.$$

It follows that

$$\begin{aligned} u_h(x_i) &= \sum_{k=1}^n f_k v_k^i \\ &= \sum_{k=1}^n G(x_i, x_k) f_k \\ &= \sum_{k=1}^n G(x_i, x_k) \int_0^1 f(x) \varphi_k(x) \, dx \\ &= \int_0^1 f(x) \sum_{k=1}^n G(x_i, x_k) \varphi_k(x) \, dx \end{aligned}$$



and  $\mathbf{u} = \mathbf{u}_p$ . We then get the following cascading solution:

$$\begin{aligned}
u^i &= \sum_{\alpha_{p-1}} u_{p-1}^{\alpha_{p-1}} v_{\alpha_{p-1}}^i \\
&= \sum_{\alpha_{p-2}} \sum_{\alpha_{p-1}} u_{p-2}^{\alpha_{p-2}} v_{\alpha_{p-2}}^{\alpha_{p-1}} v_{\alpha_{p-1}}^i \\
&\quad \vdots \\
&= \sum_{\alpha_0} \cdots \sum_{\alpha_{p-1}} f^{\alpha_0} v_{\alpha_0}^{\alpha_1} v_{\alpha_1}^{\alpha_2} \cdots v_{\alpha_{p-2}}^{\alpha_{p-1}} v_{\alpha_{p-1}}^i.
\end{aligned} \tag{B.24}$$

From the above expression we can derive the inverse of  $S^p$  element-wise. Indeed, we must have

$$u^i = \sum_{j=1}^n (S^{-p})_{ij} f^j, \tag{B.25}$$

and thus we find

$$(S^{-p})_{ij} = \sum_{\alpha_1} \cdots \sum_{\alpha_{p-1}} v_j^{\alpha_1} v_{\alpha_1}^{\alpha_2} \cdots v_{\alpha_{p-2}}^{\alpha_{p-1}} v_{\alpha_{p-1}}^i. \tag{B.26}$$

By choosing the right boundary conditions and order of the basis functions, the systems  $S^p \mathbf{u} = \mathbf{f}$  correspond to finite element approximations of higher order problems  $-u^{(p)} = f$ . We will not investigate these connections any further at this time.

## B.5 Mixed Boundary Conditions

To illustrate that we can also exactly invert other Laplace matrices, different boundary conditions are considered. We impose the following Robin or mixed boundary conditions:

$$-u'(0) + \alpha_0 u(0) = 0, \quad u'(1) + \alpha_1 u(1) = 0. \tag{B.27}$$

Integrating twice gives the following weak solution to this problem:

$$u(x) = c_0 + c_1 x - H(x-a)(x-a), \tag{B.28}$$

where the integration constants are given by

$$c_0 = \frac{\alpha_1(1-a) + 1}{\alpha_0 + \alpha_0 \alpha_1 + \alpha_1}, \quad c_1 = \frac{\alpha_0 \alpha_1(1-a) + \alpha_0}{\alpha_0 + \alpha_0 \alpha_1 + \alpha_1}.$$

Note that a solution only exists whenever  $\alpha_0 + \alpha_0 \alpha_1 + \alpha_1 \neq 0$ . Due to the different choice of boundary conditions, the corresponding weak form is also different:

Find  $u \in H^1(0, 1)$  such that

$$\int_0^1 u' \varphi' dx + \alpha_0 u(0) \varphi(0) + \alpha_1 u(1) \varphi(1) = \varphi(a), \quad \forall \varphi \in H^1(0, 1). \tag{B.29}$$

We consider again a uniform grid with  $n + 1$  and meshwidth  $h = 1/(n + 1)$ . The unknowns are now located at  $x_0, x_1, \dots, x_{n+1}$ , where  $x_i = ih$ . Also let  $a$  equal one of the grid points;  $a = x_k$  for some  $k \in \{0, \dots, n + 1\}$ . Then the system of Galerkin equations is given by  $S \mathbf{u}_k = \mathbf{e}_k$ , where  $\mathbf{e}_k$  is the  $k$ 'th standard basis vector and  $S$  is the  $(n + 2) \times (n + 2)$  matrix

$$S = \frac{1}{h} \begin{pmatrix} 1 + \alpha_0 h & -1 & & & & \\ -1 & 2 & -1 & & & \\ & & \ddots & \ddots & \ddots & \\ & & & -1 & 2 & -1 \\ & & & & -1 & 1 + \alpha_1 h \end{pmatrix}. \tag{B.30}$$

The proof of Theorem 12 can be repeated almost exactly to show that the inverse of  $S$  is given by the exact solution (B.28) evaluated at the grid points:

$$\begin{aligned} S_{jk}^{-1} &= c_0 + c_1 x_j - (x_j - x_k)H(j - k) \\ &= \frac{(\alpha_1(1 - a) + 1)(1 + \alpha_0 x_j)}{\alpha_0 + \alpha_0 \alpha_1 + \alpha_1} - (x_j - x_k)H(j - k). \end{aligned} \tag{B.31}$$