



**HAL**  
open science

# Advanced polyhedral discretization methods for poromechanical modelling

Michele Botti

► **To cite this version:**

Michele Botti. Advanced polyhedral discretization methods for poromechanical modelling. General Mathematics [math.GM]. Université Montpellier, 2018. English. NNT : 2018MONT041 . tel-01871074v3

**HAL Id: tel-01871074**

**<https://theses.hal.science/tel-01871074v3>**

Submitted on 11 Apr 2019

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

**THÈSE POUR OBTENIR LE GRADE DE DOCTEUR  
DE L'UNIVERSITE DE MONTPELLIER**

**En Mathématiques et Modélisation**

**École doctorale : Information, Structures, Systèmes**

**Unité de recherche : Institut Montpellierain Alexander Grothendieck**

**Advanced polyhedral discretization methods for  
poromechanical modelling**

**Présentée par Michele BOTTI**

**Le 27 Novembre 2018**

**Sous la direction de Daniele DI PIETRO**

**Devant le jury composé de**

**BEIRÃO DA VEIGA Lourenço, Professeur des universités, Università di Milano-Bicocca**

**BOFFI Daniele, Professeur des universités, Università di Pavia**

**DI PIETRO Daniele, Professeur des universités, Université de Montpellier**

**GALLOUËT Thierry, Professeur des universités, Institut de Mathématiques de Marseille**

**SOCHALA Pierre, Ingénieur, BRGM**

**VOHRALÍK Martin, Directeur de recherche, INRIA Paris**

**Rapporteur**

**Examineur**

**Directeur**

**Rapporteur**

**Co-encadrant**

**Président du jury**



**UNIVERSITÉ  
DE MONTPELLIER**



---

## Acknowledgements

---

Sono profondamente grato a Daniele A. Di Pietro per essere stato un eccezionale direttore di tesi durante l'intero svolgersi del dottorato. Ammiro la tua passione e dedizione per la ricerca e sono onorato di aver lavorato insieme a questa tesi. Ti ringrazio per avermi sempre manifestato fiducia ed avermi incoraggiato durante tutto il percorso. La tua disponibilità e il tuo supporto sono stati di grandissimo aiuto per raggiungere al meglio questo traguardo.

I would like to thank my second supervisor, Pierre Sochala, for welcoming me at BRGM within the best conditions. Thank you for your kindness, your listening, your passion and for all the time we have spent together. It has been a great pleasure to work with you.

I would also like to express my sincere gratitude to Lourenço Beirão da Veiga and Thierry Gallouët, who have accepted the heavy task of being referees, for the interest they have shown in my work, and their pertinent suggestions. I am also very grateful to Daniele Boffi and Martin Vohralik for having accepted to be members of the PhD committee.

I would like to thank all the members of the SPU research team that I met during the 9 months at the BRGM. I am particularly grateful to André Burnol, Farid Smaï, Fernanda M.L. Veloso, and Hideo Aochi for the fruitful discussions about poromechanics and reservoir modelling. I also express my special thanks to Olivier Le Maître for his collaboration in the results of the last chapter of this manuscript. I really appreciate your great intuitions and clear explanations.

At the IMAG, I would like to thank the researchers and professors with whom I had the opportunity to exchange and work with, especially Berardo Ruffini, Michele Bolognesi, Matthieu Alfaro, Philippe Roche and Vanessa Lleras. I am also grateful to the administrative staff, including Bernadette Lacan, Geneviève Carrière, and Carmela Madonia for their efficiency and patience. Many thanks also to the PhD and post-doctoral students met during the years in Montpellier, especially Abel, Alexandre, Francesco, Jocelyn, Joubine, Mario, Nestor, Paul, Paul-Marie, Rita, Robert, Rodrigo, and Tiffany. I am also particularly grateful to Alice, Fatima, Roberta for the good moments we shared in our office. Finally, I would like to express a huge thank to Florent for being a true friend during these three years, not only for the fantastic time we had while traveling around the world, but also for being close to me when I needed it most.

Questi tre anni di dottorato mi hanno permesso di incontrare tante altre persone stupende che desidero ringraziare perché, a modo loro, hanno fortemente contribuito alla realizzazione di questa tesi. Ringrazio Florent, Alex, Tim, Paul e Alec per avermi accolto al mio arrivo all'IMAG ed avermi fatto sentire sin da subito a casa. Un ringraziamento speciale va anche a tutti i coinquilini dell'appartamento di *Rue Rondelet*: Gregoire, Jean-Florent, Claire, Swan

e Pierre-Alexis. Condividere le cene, le serate, i weekend con voi ha riempito di gioia la mia esperienza a Montpellier.

Desidero ringraziare gli amici dello *Sborgomeo*: Rumi, Stiv, Max, Maffo, Zenk, Andre e Susi; che sono sempre al mio fianco. I momenti condivisi con voi li porterò nel cuore. Un grande grazie anche agli *Sugarcandy Mountains*: Cano, Jaya, Oscar, J, Zenk e Zoia. I concerti e le prove insieme sono per me un fondamentale momento di amicizia e divertimento. Sono felice che il nostro progetto musicale continui. Scusandomi con chi non ho menzionato, ringrazio anche tutti coloro che, nonostante la lontananza, mi sono stati vicini manifestandomi il loro affetto e la loro amicizia.

Sarò eternamente riconoscente ai miei genitori per avermi fatto sentire il loro amore anche a centinaia di chilometri di distanza. Vi sono grato per tutto ciò che avete fatto e continuate a fare per me, per gli insegnamenti ed i valori che mi avete trasmesso. Grazie mamma per il tuo sesto senso capace di sentire quando qualche nuvolone grigio fa capolino nel mio cielo anche quando cerco di nascondere. Grazie papà per il tuo cuore immenso e per avermi insegnato la gioia del dare agli altri. Grazie Mamie per esserti sempre preoccupata per me. Grazie Stefano per la tua vicinanza, la tua energia e la voglia di passare dei bei momenti insieme. Grazie Lorenzo per tutte le ore trascorse a confrontarci sulle nostre teorie e riflessioni. Sono davvero felice di condividere con te la passione per la matematica. Ringrazio immensamente Anna e Martine per il loro supporto e per esserci per celebrare questo importante traguardo.

Concludo ringraziando la persona con cui più di tutte ho condiviso le soddisfazioni e le fatiche di questo dottorato. Grazie Claudia per esserci sempre stata con il tuo prezioso sostegno e per tutta la bellezza e l'amore che mi hai regalato in questi anni.





---

# Résumé

---

Dans cette thèse, on s'intéresse à de nouveaux schémas de discrétisation afin de résoudre les équations couplées de la poroélasticité et nous présentons des résultats analytiques et numériques concernant des problèmes issus de la poromécanique. Nous proposons de résoudre ces problèmes en utilisant les méthodes Hybrid High-Order (HHO), une nouvelle classe de méthodes de discrétisation polyédriques d'ordre arbitraire. Cette thèse a été conjointement financée par le Bureau de Recherches Géologiques et Minières (BRGM) et le LabEx NUMEV. Le couplage entre l'écoulement souterrain et la déformation géomécanique est un sujet de recherche crucial pour les deux institutions de cofinancement.

## Contexte et motivations

On prend en considération des matériaux constitués d'un squelette solide et d'un réseau de pores, connectés entre eux et permettant le passage d'un fluide interstitiel. Le comportement de ce fluide peut influencer sur celui du squelette et réciproquement. De nombreux matériaux présentent cette caractéristique : les matériaux minéraux comme les sols et les roches, des matériaux organiques comme le bois, les os ou le cerveau et les matériaux industriels comme certains joints d'étanchéité. Dans cette thèse, l'accent porte sur la modélisation des sols et des roches. Le milieu poreux est vu comme la superposition de deux milieux continus : le squelette solide déformable et la phase fluide. Le milieu poreux possède une cinématique déterminée par le squelette et les sollicitations considérées peuvent être de type mécanique, hydraulique et thermique. Les sollicitations mécaniques peuvent être causées par des contraintes imposées, telles que le creusement d'ouvrages, ou bien par des contraintes intrinsèques, telles que la consolidation du sol par gravité. Les sollicitations hydrauliques sont liées à l'écoulement du fluide à travers le milieu, tandis que les sollicitations thermiques peuvent être dues à la présence de sources de chaleur. Lorsque les trois phénomènes sont traités de manière couplée, le cadre est celui dit de la thermo-hydro-mécanique. Dans cette thèse, on s'intéresse avant tout au modèle hydro-mécanique saturé obtenu en supposant que les variations de température soient négligeables.

L'intérêt pour les mécanismes couplés de diffusion-déformation était initialement motivé par le problème de la consolidation [27, 189], à savoir, le tassement progressif du sol dû à l'extraction des fluides. Aujourd'hui la théorie de la poroélasticité présente un intérêt pour les scientifiques et les ingénieurs en raison de son potentiel d'applications dans la mécanique des sols, l'industrie pétrolière, la géomécanique environnementale et la biomécanique. En simulation de réservoir, le couplage mécanique-écoulement joue un rôle important pour l'étude des problèmes de compaction et subsidence induits par la mise en production de



réservoirs peu consolidés, pour la stabilité des puits, ou encore la fracturation hydraulique. La non prise en compte de ce couplage peut aussi conduire à de mauvaises prédictions de la production. Le couplage poroélastique est par ailleurs crucial pour l'étude des risques liés à l'injection et au stockage du CO<sub>2</sub> dans les aquifères salins, comme la fuite du CO<sub>2</sub> par le puits ou bien la réactivation mécanique des failles. La surpression entraînée par l'injection peut aussi modifier les contraintes appliquées sur la roche couverture et provoquer une fracturation qui remettrait en cause l'intégrité du stockage.

Dans cette thèse, nous considérons principalement trois problèmes liées à la modélisation de la géomécanique : le problème couplé de Biot linéaire dans le Chapitre 2, le comportement non-linéaire en mécanique et poromécanique dans les chapitres 3 et 4 et la poroélasticité avec coefficients aléatoires dans le Chapitre 5. Comme noté dans [68], de nombreux résultats expérimentaux suggèrent que la réponse volumétrique de la roche poreuse au changement de la pression totale est en réalité non-linéaire. Cette non-linéarité est généralement associée à la fermeture/ouverture de fissures ou, dans les roches très poreuses, à l'effondrement progressif des pores. Les études physiques du comportement mécanique non-linéaire des milieux poreux [22, 28] ont été motivées par la nécessité de quantifier l'effet de la diminution de la pression interstitielle lors de l'épuisement d'un gisement de pétrole ou de gaz sur le volume de la matrice rocheuse.

En simulation numérique, la prise en compte des incertitudes dans les données d'entrée (telles que les constantes physiques, les conditions limites et initiales, le forçage externe et la géométrie) est un problème crucial, en particulier dans l'analyse des risques. L'évaluation des risques dans les applications de poroélasticité reste un défi majeur pour un large éventail d'applications de gestion des ressources essentielles. Les préoccupations récentes de la population concernant la sismicité induite et la contamination des eaux souterraines soulignent la nécessité de quantifier la probabilité d'événements nocifs associés à l'écoulement et à la déformation de la subsurface. La prédiction des contraintes critiques qui peuvent compromettre l'intégrité des roches couvertures et la stabilité des puits dans les réservoirs de stockage présente un intérêt particulier. La clé de la réussite de l'évaluation des risques dans les applications poroélastiques est la capacité de prédire la probabilité d'événements critiques en se basant sur une approximation empirique des paramètres matériels et de la variabilité du terme source. La quantification d'incertitude fournit des barres d'erreur numériques qui facilitent la comparaison avec l'observation expérimentale et l'évaluation des modèles physiques. De plus, elle permet d'identifier les paramètres incertains qui doivent être mesurés avec plus de précision car ils ont un impact plus significatif sur la solution. Parmi les différentes techniques conçues pour la propagation et la quantification de l'incertitude dans les modèles numériques, dans cette thèse nous considérons les méthodes spectrales stochastiques. L'idée centrale de ces méthodes est la décomposition de quantités aléatoires sur des bases d'approximation appropriées telles que les expansions en Polynômes de Chaos [112, 137].

## **Discrétisation**

Ces dernières années, un effort important a été consacré au développement et à l'analyse de méthodes numériques capables de traiter des maillages plus généraux que les maillages simpliciaux ou cartésiens. Dans le contexte de la poroélasticité, cette nécessité est motivée par

la présence des couches géologiques et des fractures dans le milieu poreux. La discrétisation du domaine peut également inclure des éléments dégénérés (comme dans la région proche du puits dans la modélisation des réservoirs) et des interfaces non-conformes apparaissant lors de la modélisation de l'érosion et de la formation de failles. Dans le contexte de la mécanique des structures, les méthodes supportant les maillages polyédriques généraux peuvent être utiles pour plusieurs raisons, notamment l'utilisation de nœuds suspendus pour les problèmes de contact et d'interface et la robustesse par rapport aux distorsions et fractures. De plus, les méthodes polyédriques permettent de simplifier les procédures de raffinement et de déraffinement du maillage utilisées pour l'adaptation.

Dans cette thèse, on s'intéresse à deux familles de méthodes de discrétisation polyédriques d'ordre arbitraire : les méthodes HHO et Galerkin discontinues (dG). L'utilisation de méthodes d'ordre élevé peut accélérer la convergence en présence de solutions régulières, ou lorsqu'elles sont associées au raffinement local du maillage. De plus, la construction de problèmes discrets à l'aide des discrétisations HHO et dG est valable en dimension d'espace arbitraire permettant d'envisager une mise en œuvre indépendante de la dimension. Les deux méthodes sont non-conformes dans le contexte de formulations primales de problèmes elliptiques, car aucune condition de continuité n'est imposée entre les éléments voisins en dG et entre les faces voisines en HHO. Nous nous concentrons sur la méthode HHO pour la discrétisation de la mécanique et sur la méthode dG pour la discrétisation de l'écoulement.

Les méthodes HHO, introduites dans [78] et [74], reposent sur des formulations primales de problèmes elliptiques et aboutissent à des systèmes symétriques et définis positifs. Les inconnues discrètes sont des polynômes du même ordre sur les éléments et les faces. Ces dernières établissent des connexions inter-éléments aux interfaces et peuvent être utilisées pour imposer des conditions limites essentielles aux faces de bord. Les inconnues de maille sont des variables intermédiaires qui peuvent être éliminées du système global par condensation statique, comme détaillé dans le Chapitre 2. La conception des méthodes HHO se fait en deux étapes : (i) tout d'abord, on reconstruit des opérateurs différentiels discrets basés sur la résolution de problèmes locaux peu coûteux à l'intérieur de chaque élément et ensuite, (ii) on introduit une stabilisation qui relie entre elles les inconnues discrètes d'élément et de face. Les définitions des opérateurs de reconstructions différentiels sont basées sur des contreparties discrètes des intégrations par parties. Même si les problèmes locaux liés aux opérateurs de reconstruction doivent être résolus, les résultats numériques de [74] indiquent que le coût associé devient marginal par rapport au coût de la résolution du système global. De plus, les méthodes HHO sont efficaces, car l'utilisation d'inconnues de face donne des stencils compacts. Les méthodes HHO offrent d'autres avantages, dont la possibilité d'établir des propriétés de conservation de quantités physiques et la robustesse par rapport à l'hétérogénéité des coefficients. La discrétisation HHO étudiée dans le Chapitre 3 s'inspire des travaux récents sur les opérateurs Leray–Lions [69, 70], où les auteurs montrent que la méthode est robuste en ce qui concerne les non-linéarités.

Les méthodes dG peuvent être considérées comme des méthodes éléments finis permettant des discontinuités dans les espaces discrets, ou bien comme des volumes finis dans lesquels la solution approchée est représentée dans chaque élément du maillage par un polynôme au lieu que par une constante. Permettre à la solution d'être continue par morceaux offre une grande flexibilité dans la conception du maillage. En plus d'être adaptées aux maillages

généraux y compris non-conformes, les méthodes dG ont l'avantage de pouvoir monter en ordre très facilement. Leur principal inconvénient est le grand nombre d'inconnues qu'elles engendrent, et donc le coût de résolution des systèmes. Une analyse unifiée des méthodes dG pour les problèmes elliptiques peut être trouvée dans [42] et [76]. La stratégie fondamentale pour approcher le problème de la diffusion hétérogène (tel que le flux darcéen dans les milieux poreux) en utilisant les méthodes dG consiste à pénaliser les sauts d'interface en utilisant une pondération dépendante de la moyenne harmonique des composantes normales du tenseur de diffusion. Cette technique permet de pénaliser de façon optimale. De plus, comme indiqué dans les chapitres 2 et 4 et dans [77], la méthode proposée est robuste en ce qui concerne les variations spatiales du coefficient de perméabilité, avec des constantes dans les estimations d'erreur ayant une dépendance faible du rapport d'hétérogénéité.

## Structure du manuscrit

Par la suite, nous allons résumer brièvement le contenu du manuscrit en mettant l'accent sur les résultats marquants.

Dans le **Chapitre 1** nous présentons les modèles de poroélasticité que nous allons considérer successivement. Une fois ceux-ci introduits, nous faisons un inventaire des difficultés liées à leur approximation numérique, avant de présenter un état de l'art documenté nous permettant de définir les orientations des chapitres suivants.

Dans le **Chapitre 2**, publié dans *SIAM Journal on Scientific Computing* (voir [30]), nous introduisons un nouvel algorithme pour le problème de Biot, basé sur une discrétisation HHO de la mécanique et une discrétisation dG Symmetric Weighted Interior Penalty (SWIP) du flux. La méthode a plusieurs atouts, notamment la validité en deux et trois dimensions, la stabilité inf-sup, le support des maillages polyédriques généraux et l'utilisation de l'ordre d'approximation arbitraire en espace. De plus, le coût de résolution peut être réduit en condensant statiquement un grand sous-ensemble d'inconnues. Notre analyse fournit des estimations de stabilité et d'erreur qui se maintiennent même lorsque le coefficient de stockage spécifique est nul et montre que les constantes dépendent faiblement de l'hétérogénéité du coefficient de perméabilité. Nous discutons les détails de la mise en œuvre de la méthode et fournissons des tests numériques démontrant ses performances. Enfin, nous montrons que le schéma est localement conservateur sur le maillage primal, une propriété souhaitable pour les praticiens et fondamentale pour les estimations a posteriori basées sur des flux équilibrés.

Dans le **Chapitre 3**, publié dans *SIAM Journal on Numerical Analysis* (voir [33]), nous proposons et analysons une nouvelle discrétisation HHO d'une classe de modèles d'élasticité (linéaire et) non-linéaire couramment utilisées en mécanique des solides. La méthode satisfait un principe local de travail virtuel à l'intérieur de chaque élément du maillage, avec des tractions numériques d'interface qui obéissent à la loi d'action et réaction. Une analyse complète couvrant des lois de comportement mécanique très générales est effectuée. En particulier, nous prouvons l'existence d'une solution discrète et nous identifions une hypothèse de monotonie stricte sur la loi contrainte-déformation qui assure l'unicité. La convergence aux solutions de régularité minimale est démontrée en utilisant un argument de compacité. Une estimation d'erreur optimale d'ordre  $h^{k+1}$  de la norme d'énergie est alors prouvée dans les conditions supplémentaires de continuité de Lipschitz et de forte monotonie sur la loi de

comportement. Les performances de la méthode sont largement étudiées sur un panel complet de tests numériques avec des lois contrainte-déformation correspondant à des matériaux réels.

Dans le **Chapitre 4**, tiré de [35], nous construisons et analysons une méthode couplée HHO-dG pour la poroélasticité non-linéaire. Il y a deux différences principales par rapport à la méthode conçue dans le Chapitre 2 pour la version linéaire du problème. Premièrement, le gradient symétrique discret se situe dans l'espace complet des polynômes à valeurs tensorielles, par opposition aux gradients symétriques des polynômes à valeurs vectorielles. Comme montré dans le Chapitre 3, cette modification est nécessaire pour l'analyse de convergence en présence de lois de comportement non-linéaires. Deuxièmement, le terme de droite du problème discret est obtenu en prenant la moyenne en temps de la force de chargement  $f$  et de la source de fluide  $g$  entre deux pas de temps consécutifs. Cette modification nous permet de prouver la stabilité et l'estimation d'erreur sans aucune hypothèse de régularité supplémentaire sur les données. Un autre résultat marquant de ce chapitre est une nouvelle preuve de l'inégalité de Korn discrète, ne nécessitant pas d'hypothèse géométrique particulière sur le maillage.

Le **Chapitre 5** contient des perspectives sur la solution numérique du problème de Biot avec coefficients poroélastiques aléatoires dans le contexte de la quantification des incertitudes. Il recueille une partie des travaux en cours réalisés lors du stage au BRGM qui s'est déroulé de janvier à septembre 2018. L'incertitude est modélisée par un ensemble fini de paramètres avec une distribution de probabilité prescrite. Une attention particulière est accordée pour assurer que la paramétrisation des coefficients soit physiquement admissible. Nous présentons la formulation faible du système d'équations différentielles stochastiques et établissons sa bonne position. Nous discutons ensuite de l'approximation du problème par des techniques non-intrusives basées sur le développement de solutions sur des bases de Polynômes de Chaos. La procédure de projection spectrale considérée permet de réduire le problème stochastique à un ensemble fini de simulations déterministes paramétriques discrétisées par le schéma HHO-dG du Chapitre 2. Nous étudions numériquement la convergence de l'erreur par rapport au niveau de la grille dans l'espace de probabilité. Enfin, nous effectuons une analyse de sensibilité pour évaluer la propagation de l'incertitude en entrée sur les champs de déplacement et pression dans un cas test d'injection et un problème de traction.

En **Annexe A**, publié dans les actes de la conférence *Finite Volumes for Complex Applications VIII* (voir [34]), nous présentons une variante du schéma HHO-dG pour le problème de la poroélasticité non-linéaire étudié dans le Chapitre 4. En particulier cette annexe fournit des tests numériques démontrant la convergence de la méthode en présence de lois de comportement non-linéaires de type Hencky-Mises.

L'**Annexe B** présente des travaux réalisés en marge de la ligne directrice de ce manuscrit. Elle est tirée de [36] publié dans *Computational Methods in Applied Mathematics*. Nous considérons les problèmes hyperélastiques et leurs solutions numériques en utilisant des algorithmes de discrétisation par éléments finis conformes. Pour ces problèmes, nous présentons des reconstructions du tenseur de contraintes conformes en  $H(\text{div})$ , équilibrées et faiblement symétriques, obtenues à partir de problèmes locaux à l'aide des espaces d'éléments finis Arnold–Falk–Winther. Les reconstructions sont indépendantes de la loi de comportement

mécanique. Sur la base de ces reconstructions du tenseur des contraintes, nous obtenons une estimation d'erreur à posteriori en distinguant les estimations d'erreur de discrétisation, de linéarisation et de quadrature, et proposons un algorithme adaptatif équilibrant ces différentes sources d'erreur. Nous confirmons l'efficacité de l'estimation par un test numérique avec une solution analytique. Nous appliquons ensuite l'algorithme adaptatif à un test davantage axé sur l'application.





---

# Contents

---

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Poroelasticity . . . . .	1
1.2	Discretization . . . . .	9
1.3	Uncertainty quantification . . . . .	12
1.4	Plan of the manuscript . . . . .	13
<b>2</b>	<b>The Biot problem</b>	<b>17</b>
2.1	Introduction . . . . .	18
2.2	Discretization . . . . .	20
2.3	Stability analysis . . . . .	27
2.4	Error analysis . . . . .	31
2.5	Implementation . . . . .	37
2.6	Numerical tests . . . . .	40
2.7	Appendix: Flux formulation . . . . .	44
<b>3</b>	<b>Nonlinear elasticity</b>	<b>49</b>
3.1	Introduction . . . . .	50
3.2	Setting and examples . . . . .	52
3.3	Notation and basic results . . . . .	54
3.4	The Hybrid High-Order method . . . . .	55
3.5	Analysis . . . . .	58
3.6	Local principle of virtual work and law of action and reaction . . . . .	68
3.7	Numerical results . . . . .	70
3.8	Appendix: Technical results . . . . .	74
<b>4</b>	<b>Nonlinear poroelasticity</b>	<b>83</b>
4.1	Introduction . . . . .	84
4.2	Continuous setting . . . . .	86
4.3	Discrete setting . . . . .	89
4.4	Discretization . . . . .	93
4.5	Stability and well-posedness . . . . .	98
4.6	Convergence analysis . . . . .	103



---

<b>5</b>	<b>Poroelasticity with uncertain coefficients</b>	<b>111</b>
5.1	Introduction . . . . .	112
5.2	The Biot model . . . . .	113
5.3	Probabilistic framework . . . . .	116
5.4	The Biot problem with random coefficients . . . . .	120
5.5	Discrete setting . . . . .	126
5.6	Test cases . . . . .	131
<b>A</b>	<b>A nonconforming high-order method for nonlinear poroelasticity</b>	<b>141</b>
A.1	Introduction . . . . .	141
A.2	Mesh and notation . . . . .	142
A.3	Discretization . . . . .	143
A.4	Formulation of the method . . . . .	145
A.5	Numerical results . . . . .	146
<b>B</b>	<b>A posteriori error estimation for nonlinear elasticity</b>	<b>147</b>
B.1	Introduction . . . . .	147
B.2	Setting . . . . .	149
B.3	Equilibrated stress reconstruction . . . . .	150
B.4	A posteriori error estimate and adaptive algorithm . . . . .	154
B.5	Numerical results . . . . .	158
	<b>List of Figures</b>	<b>165</b>
	<b>Bibliography</b>	<b>169</b>





# Chapter 1

---

## Introduction

---

### Contents

---

<b>1.1 Poroelasticity</b> . . . . .	<b>1</b>
1.1.1 Context . . . . .	2
1.1.2 Governing equations . . . . .	3
1.1.3 The Biot model . . . . .	5
1.1.4 Nonlinear poroelasticity . . . . .	7
<b>1.2 Discretization</b> . . . . .	<b>9</b>
1.2.1 Numerical issues . . . . .	9
1.2.2 Polyhedral element methods . . . . .	10
<b>1.3 Uncertainty quantification</b> . . . . .	<b>12</b>
<b>1.4 Plan of the manuscript</b> . . . . .	<b>13</b>

---

In this manuscript we focus on novel discretization schemes for solving the coupled equations of poroelasticity and we present analytical and numerical results for poromechanics problems relevant to geoscience applications. We propose to solve these problems using Hybrid High-Order (HHO) methods, a new class of nonconforming high-order methods supporting general polyhedral meshes.

## 1.1 Poroelasticity

Porous media are solid materials comprising a great number of interconnected pores allowing for fluid flow through the medium. The presence of a moving fluid influences the mechanical response of the solid matrix. Two mechanisms play a central role in this interaction between the mechanical behavior and the fluid dynamics: (i) an increase of the pore pressure induces a dilatation of the rock in response to the added stress, and (ii) a compression of the porous skeleton leads to a rise of pore pressure, if the change of the mechanical state is fast relative to the fluid flow rate. These two coupled deformation-diffusion phenomena lie at the heart

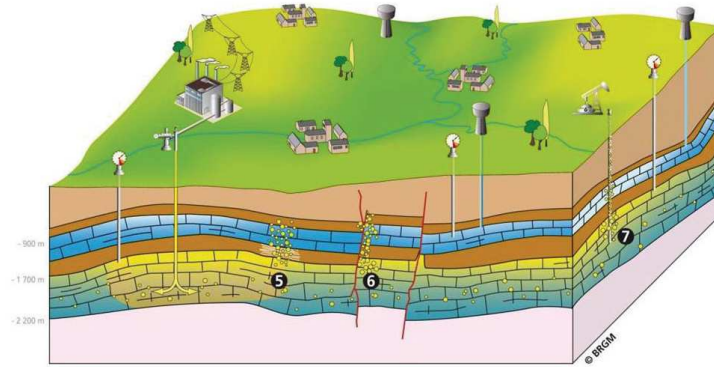


Figure 1.1: Geologic storage of carbon dioxide in saline aquifers

of the theory of poroelasticity. In accordance with these key phenomena, the fluid-filled porous medium acts in a time-dependent manner. Indeed, if the fluid pressure caused by the deformation is allowed to dissipate through a Darcian mass transport, further deformation of the rock progressively appear. At the same time, the induced poroelastic stresses will, in turn, respond back to the fluid pressure field.

### 1.1.1 Context

This Ph.D. thesis was conjointly founded by the Bureau de recherches géologiques et minières (BRGM) and LabEx NUMEV. The coupling between subsurface flow and geomechanical deformation is a crucial research topic for both cofunding institutions. The main goals of BRGM are to implement decision-support tools designed to anticipate and prevent subsurface risks and to establish safety criteria for human geological activities. There is, therefore, a particular interest in poromechanical modelling, that is critical in the assessment of the environmental impacts of groundwater use and exploitation of shale gas reserves. In particular, seismicity induced by fluid injection and withdrawal has emerged as a central element of the scientific discussion around subsurface technologies that tap into water and energy resources.

Interest in the coupled diffusion-deformation mechanisms was initially motivated by the problem of consolidation, namely the progressive settlement of a soil due to fluid extraction. The earliest theory modeling the effects of the pore fluid on the deformation of soils was developed in the pioneering work of Terzaghi [189], who proposed a model for consolidation accounting for the fluid-to-solid coupling only. In this case, the problem can be decoupled and solved in two stages. This kind of theory can successfully model some of the poroelastic processes in the case of highly compressible fluids such as air. However, when one deals with slightly compressible (or incompressible) fluids, the solid-to-fluid coupling cannot be neglected since the changes in the stress can significantly influence the pore pressure. The first detailed mathematical theory of poroelasticity incorporating both the basic phenomena outlined above, was formulated by Biot [27]. The model proposed by Biot was subsequently re-derived via homogenization [9, 45] and mixture theory [146, 147], what placed the Biot theory on a rigorous basis.

Poroelasticity has since been explored in a large number of geomechanical applications,

such as: subsidence due to fluid withdrawal, reservoir impoundment, tensile failure induced by pressurization of a borehole, waste disposal, earthquake triggering due to pressure induced faults slip, and injection-production cycles in geothermal fields. Recently, coupled flow and geomechanics has also gained attention due to its role in the long-term geologic storage of carbon dioxide in saline aquifers (cf. Figure 1.1), which is widely regarded as a promising technology to help mitigate the climate change by reducing anthropogenic CO<sub>2</sub> emissions into the atmosphere. Injection of CO<sub>2</sub> requires a compression of the ambient groundwater and an overpressurization of the aquifer, which could fracture the caprock, trigger seismicity, or activate faults. BRGM plays a key role in the research on CO<sub>2</sub> storage in geological layers, especially in deep aquifers, and develops tools to understand and prevent the onset of the previously mentioned side effects. Other examples of poroelastic structures include cartilage, skin, bone, the myocardium, the brain, and the lungs. Consequently, notable contributions have also been made in a diverse range of biomechanics applications [6, 15, 44, 187, 193].

### 1.1.2 Governing equations

We use a classical continuum representation in which the fluid and the solid skeleton are viewed as overlapping continua [18, 61]. The poroelasticity system, describing the fluid flow in a deformable saturated porous medium, consists of one equation expressing the momentum balance and one expressing the fluid mass conservation law. At the macroscopic scale, they are derived in the work of Biot [27] and Terzaghi [190]. In what follows, the elastic structure  $\Omega \subset \mathbb{R}^d$ ,  $d \in \{2, 3\}$ , forms a porous and permeable matrix saturated by a slightly compressible and viscous fluid which diffuses through it. We assume that the material is isotropic and the conditions are isothermal. In many important applications, such as geothermal energy extraction, nuclear waste disposal, and carbon storage, temperature plays a significant role and must therefore be included in the model. Thermo-poroelastic models are derived in [43, 61, 107] and consists in an additional coupling between deformation, heat diffusion, and fluid flow. Since the pressure and temperature play a similar role in the deformation of the body, in what follows we only focus on the hydromechanical coupling. The displacement of the structure is denoted by  $\mathbf{u} := \{u_1, \dots, u_d\}$  and the fluid pore pressure by  $p$ .

The momentum conservation equation is similar to that found in the context of elasticity, the exception being the addition of a fluid pressure term. Letting  $V$  be a fixed, arbitrary open subset of  $\Omega$ , the momentum of the corresponding portion of the matrix is given by  $\int_V \frac{\partial \mathbf{u}}{\partial t} dV$ . The forces acting on  $V$  consist in the traction forces applied by the complement of  $V$  across its boundary  $\partial V$  and the volume-distributed external forces. Then, for a volumetric load  $\mathbf{f}$  and total stress tensor  $\tilde{\sigma}$ , the momentum balance equation reads

$$\frac{\partial}{\partial t} \int_V \frac{\partial \mathbf{u}}{\partial t} dV = \int_{\partial V} \tilde{\sigma} \mathbf{n} dS + \int_V \mathbf{f} dV,$$

where  $\mathbf{n}$  is the outward normal of  $\partial V$ . Owing to the divergence theorem, this gives

$$\int_V \frac{\partial^2 \mathbf{u}}{\partial t^2} dV = \int_V \nabla \cdot \tilde{\sigma} dV + \int_V \mathbf{f} dV.$$

In the classical poroelasticity model, since the deformation of the material is usually much slower than the flow rate, the inertial effects are considered negligible. This quasi-static assumption consists in ignoring the second time-derivative for displacements and, since  $V$  was chosen arbitrarily, it follows that

$$-\nabla \cdot \tilde{\boldsymbol{\sigma}} = \mathbf{f} \quad \text{in } \Omega. \quad (1.1)$$

Turning to the mass conservation equation, the variables of interest are the fluid content  $\eta$  (fluid mass per unit bulk volume of porous medium), the flux  $\mathbf{w}$  (fluid mass flow rate per unit area and time), and the volumetric fluid source  $g$ . Taking  $V \subset \Omega$  as before, the rate at which fluid moves across the boundary  $\partial V$  is given by  $\int_{\partial V} \mathbf{w} \cdot \mathbf{n} \, dS$ . Then the conservation of mass for isothermal single-phase flow of a slightly compressible fluid takes the integral form

$$\frac{\partial}{\partial t} \int_V \eta \, dV + \int_{\partial V} \mathbf{w} \cdot \mathbf{n} \, dS = \int_V g \, dV.$$

The divergence theorem applied to the second term of the left-hand side of the above equation and the fact that  $V$  was chosen arbitrarily lead to

$$\frac{\partial \eta}{\partial t} + \nabla \cdot \mathbf{w} = g. \quad (1.2)$$

The coupling between the mechanical behavior of the matrix and the fluid dynamics is realized by constitutive relations relating the total stress  $\tilde{\boldsymbol{\sigma}}$ , the flux  $\mathbf{w}$ , and the fluid content  $\eta$  to the primary variables  $\mathbf{u}$  and  $p$ . The total stress must account for the usual material stress, as in solid mechanics, and for the fluid pressure. Owing to the Terzaghi decomposition of the total stress [190], one has

$$\tilde{\boldsymbol{\sigma}} = \boldsymbol{\sigma} - \alpha p \mathbf{I}_d, \quad (1.3)$$

where  $\boldsymbol{\sigma}$  is referred to as the *effective stress tensor* and measures the material properties of the medium. The Biot–Willis coefficient  $0 < \alpha \leq 1$ , defined as the ratio of the volume change of the fluid over the volume change of the medium, accounts for the pressure-deformation interaction, and  $\mathbf{I}_d$  denotes the  $d$ -dimensional identity matrix. The effective stress tensor depends on the deformation according to the stress-strain relation  $\boldsymbol{\sigma} = \boldsymbol{\sigma}(\nabla_s \mathbf{u})$ , where the symmetric gradient operator measures the strain in the case of small deformations and is defined as

$$\nabla_s \mathbf{u} := \frac{\nabla \mathbf{u} + (\nabla \mathbf{u})^T}{2}, \quad \text{with } (\nabla \mathbf{u})_{ij} := \frac{\partial u_i}{\partial x_j}, \quad \text{for } 1 \leq i, j \leq d.$$

The standard assumption of Darcy's law for porous media holds for the flux:

$$\mathbf{w} = -\frac{\mathbb{K} \nabla p}{\mu_f}, \quad (1.4)$$

where  $\mathbb{K}$  is the tensor-valued intrinsic permeability field and  $\mu_f$  is the fluid viscosity. The third constitutive assumption links the change in the fluid content  $\eta$  with the changes in the

fluid pressure  $p$  and the material volume, which is locally measured by  $\nabla \cdot \mathbf{u}$ . More specifically, following [61, 68], we set

$$\eta = c_0 p + \alpha \nabla \cdot \mathbf{u}, \quad (1.5)$$

where the Biot–Willis coefficient  $\alpha$  quantifies the amount of fluid that can be forced into the medium by a variation of pore volume for a constant fluid pressure, while the constrained specific storage coefficient  $c_0 \geq 0$  measures the amount of fluid that can be forced into the medium by pressure increments due to the compressibility of the structure. The case of a solid matrix with incompressible grains corresponds to the limit value  $c_0 = 0$ .

### 1.1.3 The Biot model

The Biot model of linear poroelasticity is valid under the following assumptions: (i) small relative variations of porosity with respect to the equilibrium value  $\varphi \in [0, 1]$ , (ii) small relative variations of the fluid density, and (iii) infinitesimal strain theory. Moreover, the material is assumed to be isotropic, linearly elastic, and homogeneous and, as a consequence, its mechanical behavior is described by Hooke’s stress-strain law

$$\boldsymbol{\sigma}(\nabla_s \mathbf{u}) := 2\mu \nabla_s \mathbf{u} + \lambda(\nabla \cdot \mathbf{u}) \mathbf{I}_d, \quad (1.6)$$

the Lamé’s parameters  $\mu > 0$  and  $\lambda > 0$  correspond to the dilatation and shear moduli, respectively. We recall that the shear modulus  $\mu$  remains bounded and bounded away from 0, namely  $0 < \underline{\mu} \leq \mu \leq \bar{\mu} \in \mathbb{R}$ , whereas  $\lambda$  can take unboundedly large values in the case of a quasi-incompressible material ( $\lambda \rightarrow \infty$ ). We note that the condition  $\lambda > 0$  is stronger than the one required to have the coercivity of the elasticity operator in Section 2.2.2. Indeed,  $\lambda \geq -\frac{2}{d}\mu$  is sufficient to ensure the positivity of the bulk modulus and the well-posedness of the weak formulation (5.18) (see also Lemma 5.4 and [178]).

The Biot model is derived by inserting the constitutive relations (1.3), (1.4), and (1.5) into (1.1) and (1.2). Additional details on the model derivation and investigations on the poromechanical coefficient are given in Section 5.2. Hence, for a given bounded connected domain  $\Omega \subset \mathbb{R}^d$ ,  $d \in \{2, 3\}$ , with boundary  $\partial\Omega$  and outward normal  $\mathbf{n}$ , a finite time  $t_F > 0$ , a volumetric load  $\mathbf{f}$ , and a fluid source  $g$ , the corresponding problem consists in finding a vector-valued displacement field  $\mathbf{u} : \Omega \times (0, t_F] \rightarrow \mathbb{R}^d$  and a scalar-valued pore pressure field  $p : \Omega \times (0, t_F] \rightarrow \mathbb{R}$  solution of

$$-\nabla \cdot \boldsymbol{\sigma}(\nabla_s \mathbf{u}) + \alpha \nabla p = \mathbf{f} \quad \text{in } \Omega \times (0, t_F], \quad (1.7a)$$

$$c_0 d_t p + \alpha d_t(\nabla \cdot \mathbf{u}) - \nabla \cdot (\boldsymbol{\kappa} \nabla p) = g \quad \text{in } \Omega \times (0, t_F], \quad (1.7b)$$

where  $d_t$  denotes the time-derivative operator and  $\boldsymbol{\kappa} := \frac{\mathbb{K}}{\mu_f}$  the tensor-valued fluid mobility field. There are two distinct sets of boundary conditions, one corresponding to the deformation and one corresponding to the flow, and an initial condition on the fluid content to be added



to the above equations to close the model for the Biot problem:

$$\mathbf{u} = \mathbf{u}_D \quad \text{on } \Gamma_D \times (0, t_F], \quad (1.7c)$$

$$(\boldsymbol{\sigma}(\nabla_s \mathbf{u}) + \alpha p \mathbf{I}_d) \mathbf{n} = \mathbf{t}_N \quad \text{on } \Gamma_N \times (0, t_F], \quad (1.7d)$$

$$p = p_d \quad \text{on } \Gamma_d \times (0, t_F], \quad (1.7e)$$

$$\boldsymbol{\kappa} \nabla p \cdot \mathbf{n} = w_n \quad \text{on } \Gamma_n \times (0, t_F], \quad (1.7f)$$

$$(c_0 p + \alpha \nabla \cdot \mathbf{u})(\cdot, 0) = \eta_0 \quad \text{in } \Omega, \quad (1.7g)$$

where  $\Gamma_D, \Gamma_N, \Gamma_d$ , and  $\Gamma_n$  are subsets of the boundary  $\partial\Omega$  such that  $\Gamma_D \cup \Gamma_N = \Gamma_d \cup \Gamma_n = \partial\Omega$ ,  $\Gamma_D \cap \Gamma_N \neq \emptyset$  and  $\Gamma_d \cap \Gamma_n \neq \emptyset$ . If  $\Gamma_D = \emptyset$ , owing to the Neumann condition (1.7d), we need to impose a compatibility condition on the average of the forcing term  $\mathbf{f}$  and the boundary traction  $\mathbf{t}_N$ , as well as to prescribe the rigid-body motions of the medium. Thus, in this case, we have  $\Gamma_N = \partial\Omega$  and we set

$$\int_{\Omega} \mathbf{f}(\cdot, t) = \int_{\partial\Omega} \mathbf{t}_N(\cdot, t), \quad \int_{\Omega} \mathbf{u}(\cdot, t) = \mathbf{u}_{\text{av}}, \quad \text{and} \quad \int_{\Omega} \nabla \times \mathbf{u}(\cdot, t) = \mathbf{u}_{\text{rot}} \quad \forall t \in (0, t_F], \quad (1.7h)$$

where  $\nabla \times$  denotes the curl operator. For the sake of simplicity, we exclude the case of  $\Gamma_D \neq \emptyset$  with zero  $(d-1)$ -dimensional Hausdorff measure. We only mention that, if the domain is clamped on a single point  $\mathbf{x}_D \in \partial\Omega$  there is no need to fix the translations (by imposing the zero-average constraint on  $\mathbf{u}$ ), but it is still necessary to prescribe the rotations. Similarly, in the case  $\Gamma_d = \Gamma_n = \emptyset$  and  $c_0 = 0$ , we observe that, owing to (1.7b) and the inhomogeneous boundary conditions (1.7c) and (1.7f), we need the following compatibility condition relating the fluid source  $g$ , the normal flux  $w_n$ , and the Dirichlet datum  $\mathbf{u}_D$ , and the following constraint on the average of  $p$ :

$$\int_{\Omega} g(\cdot, t) = \int_{\partial\Omega} w_n(\cdot, t) + \alpha \text{d}_t(\mathbf{u}_D \cdot \mathbf{n})(\cdot, t), \quad \text{and} \quad \int_{\Omega} p(\cdot, t) = p_{\text{av}} \quad \forall t \in (0, t_F]. \quad (1.7i)$$

In the context of geomechanics, a reasonable choice for the rigid-body motions of the medium and the pressure average in (1.7h) and (1.7i) can be  $\mathbf{u}_{\text{av}} = \mathbf{0}$ ,  $\mathbf{u}_{\text{rot}} = \mathbf{0}$ , and  $p_{\text{av}}$  equal to the hydrostatic pressure measured at the barycenter of  $\Omega$ , respectively. Finally, in the case  $\Gamma_N = \emptyset$  and  $c_0 = 0$ , if looking for a strong (space-time continuous up to the boundary) solution of problem (1.7), one should also: (i) assume that the boundary data in (1.7c) and (1.7e) admit continuous traces at  $\partial\Omega \times \{t = 0\}$ , and (ii) require the average of the initial fluid content to be compatible with the Dirichlet datum  $\mathbf{u}_D$ , namely

$$\int_{\Omega} \eta_0 = \alpha \int_{\Omega} \nabla \cdot \mathbf{u}(\cdot, 0) = \int_{\partial\Omega} (\mathbf{u}_D \cdot \mathbf{n})(\cdot, 0). \quad (1.7j)$$

The symmetric mobility tensor  $\boldsymbol{\kappa} : \Omega \rightarrow \mathbb{R}_{\text{sym}}^{d \times d}$  is assumed to be uniformly bounded and uniformly elliptic; that means there exist positive constants  $\underline{\kappa}$  and  $\bar{\kappa}$  such that, for all  $\mathbf{x} \in \Omega$  and all vector  $\boldsymbol{\xi} \in \mathbb{R}^d$ ,

$$\underline{\kappa} \|\boldsymbol{\xi}\|^2 \leq \boldsymbol{\xi}^T \boldsymbol{\kappa}(\mathbf{x}) \boldsymbol{\xi} \leq \bar{\kappa} \|\boldsymbol{\xi}\|^2.$$

Parameter	Description	Unit
$\mu \in [\underline{\mu}, \bar{\mu}]$	Shear modulus	Pa
$\lambda \in (0, +\infty)$	Dilatation modulus	Pa
$\boldsymbol{\kappa} \in \mathbb{R}_{\text{sym}}^{d \times d}$	Uniformly elliptic mobility tensor	$\text{m}^2 \text{s}^{-1} \text{Pa}^{-1}$
$\varphi \in [0, 1]$	Reference porosity value	–
$\alpha \in (\varphi, 1]$	Biot–Willis coefficient	–
$c_0 \geq 0$	Constrained specific storage coefficient	$\text{Pa}^{-1}$

Table 1.1: Summary of physical parameters

The above hypothesis allows the permeability tensor  $\boldsymbol{\kappa}$  to be discontinuous in  $\Omega$  even if, in practice, it has often more regularity than just being uniformly bounded. Henceforth, it is assumed that there is a partition  $P_\Omega := \{\Omega_i\}_{1 \leq i \leq N_\Omega}$  of  $\Omega$  such that the restriction of  $\boldsymbol{\kappa}$  to each  $\Omega_i$  is constant. In geomechanics applications, the partition  $P_\Omega$  typically results from the partitioning of the porous medium into various geological layers.

In some applications, such as consolidation processes, the fluid is considered to be incompressible and the solid matrix to have very low sensitiveness with respect to pressure increments. Hence, the constrained specific storage coefficient  $c_0$  can be very small and sometimes vanishing. Since, in this case, the volume change of the solid grains composing the matrix is neglected, the volume of fluid only depends on the variations of pore volume. This model is referred to as *Biot's consolidation model*. From a numerical point of view, as we will detail further, the correct approximation of the consolidation problem is more involved than the one of the linear poroelasticity problem with  $c_0$  bounded away from 0. The well-posedness of the canonical two-field weak formulation of problem (1.7) with displacement and pressure as variables was carried out in [181, 205]. Three-field and four-field formulations, obtained by taking the Darcy flux and the total stress as independent variables, have also been analyzed and can be found in several studies, *e.g.* [131, 167, 204].

### 1.1.4 Nonlinear poroelasticity

As noted in [68, Section 3.2], many experimental results suggest that the volumetric response of porous rock to the change of total pressure is actually nonlinear. The nonlinear behavior is generally associated with the closing/opening of crack-like pores, but in very porous and weak rocks it is mainly caused by progressive pore collapse. Investigations of the nonlinear deformation of porous media have been motivated by the need to quantify the effect of the pore pressure decline during depletion of an oil or gas reservoir on the volume of the rock matrix. For a physical investigation of the nonlinear mechanical behavior of porous solids we refer the reader to [22, 28]. Another possible source of nonlinearity in the poroelasticity system is due to the dependence of the hydraulic mobility on the stress and fluid pressure. In fact, some porous media exhibit a significant difference in conductivity once the material deformations start to occur. The analysis of the nonlinear problem obtained by considering  $\boldsymbol{\kappa}$  depending on  $\nabla \cdot \mathbf{u}$  is provided in [48], while the case of a negative exponential dependence on the volumetric part of total stress is examined in [145]. A model for finite-strain poroelasticity

with deformation dependent permeability field has been proposed and analyzed in [23, 90].

In this manuscript we focus on nonlinearities of the stress-strain relation  $\sigma : \Omega \times \mathbb{R}_{\text{sym}}^{d \times d} \rightarrow \mathbb{R}_{\text{sym}}^{d \times d}$  in the small deformation regime. In particular, in Chapter 4 we study the poroelasticity problem obtained by replacing (1.6) with the hyperelastic nonlinear laws considered in Chapter 3. *Hyperelasticity* is a type of constitutive model for ideally elastic materials in which the stress is determined by the current state of deformation by deriving a stored energy density function  $\Psi : \mathbb{R}_{\text{sym}}^{d \times d} \rightarrow \mathbb{R}$ , namely

$$\sigma(\tau) := \frac{\partial \Psi(\tau)}{\partial \tau}.$$

We next discuss a number of meaningful examples.

**Example 1.1** (Hencky–Mises model). The nonlinear Hencky–Mises model of [105, 157] corresponds to the stored energy density function

$$\Psi_{\text{hm}}(\tau) := \frac{\alpha}{2} \text{tr}(\tau)^2 + \Phi(\text{dev}(\tau)), \quad (1.8)$$

where  $\text{dev} : \mathbb{R}_{\text{sym}}^{d \times d} \rightarrow \mathbb{R}$  defined by  $\text{dev}(\tau) = \text{tr}(\tau^2) - \frac{1}{d} \text{tr}(\tau)^2$  is the deviatoric operator. Here,  $\alpha \in (0, +\infty)$  and  $\Phi : [0, +\infty) \rightarrow \mathbb{R}$  is a function of class  $C^2$  satisfying, for some positive constants  $C_1, C_2$ , and  $C_3$ ,

$$C_1 \leq \Phi'(\rho) < \alpha, \quad |\rho \Phi''(\rho)| \leq C_2 \quad \text{and} \quad \Phi'(\rho) + 2\rho \Phi''(\rho) \geq C_3 \quad \forall \rho \in [0, +\infty). \quad (1.9)$$

We observe that taking  $\alpha = \lambda + \frac{2}{d}\mu$  and  $\Phi(\rho) = \mu\rho$  in (1.8) leads to the linear case (1.6). Deriving the energy density function (1.8) yields

$$\sigma(\tau) = \tilde{\lambda}(\text{dev}(\tau)) \text{tr}(\tau) \mathbf{I}_d + 2\tilde{\mu}(\text{dev}(\tau))\tau, \quad (1.10)$$

with nonlinear Lamé functions  $\tilde{\mu}(\rho) := \Phi'(\rho)$  and  $\tilde{\lambda}(\rho) := \alpha - \Phi'(\rho)$ .

In the previous example the nonlinearity of the model only depends on the deviatoric part of the strain. In the following model it depends on the term  $\tau : \mathfrak{C}\tau$ .

**Example 1.2** (An isotropic reversible damage model). The isotropic reversible damage model of [50] can also be interpreted in the framework of hyperelasticity by setting up the energy density function as

$$\Psi_{\text{dam}}(\tau) := \frac{(1 - D(\tau))}{2} \tau : \mathfrak{C}\tau + \Phi(D(\tau)), \quad (1.11)$$

where  $D : \mathbb{R}_{\text{sym}}^{d \times d} \rightarrow [0, 1]$  is the scalar damage function and  $\mathfrak{C}$  is a fourth-order symmetric and uniformly elliptic tensor, namely, for some positive constants  $C_*$  and  $C^*$ , it holds

$$C_* \|\tau\|_{d \times d}^2 \leq \mathfrak{C}\tau : \tau \leq C^* \|\tau\|_{d \times d}^2, \quad \forall \tau \in \mathbb{R}^{d \times d}. \quad (1.12)$$

The function  $\Phi : [0, 1] \rightarrow \mathbb{R}$  defines the relation between  $\tau$  and  $D$  by  $\frac{\partial \phi}{\partial D} = \frac{1}{2} \tau : \mathfrak{C}\tau$ . The resulting stress-strain relation reads

$$\sigma(\tau) = (1 - D(\tau))\mathfrak{C}\tau. \quad (1.13)$$

Another generalization of the linear stress-strain relation is obtained by adding second-order terms to (1.6).

**Example 1.3** (A second-order elasticity model). In the second-order isotropic elasticity model of [66, 124] the stress-strain relation is

$$\boldsymbol{\sigma}(\boldsymbol{\tau}) = \lambda \operatorname{tr}(\boldsymbol{\tau})\mathbf{I}_d + 2\mu\boldsymbol{\tau} + B \operatorname{tr}((\boldsymbol{\tau})^2)\mathbf{I}_d + 2B \operatorname{tr}(\boldsymbol{\tau})\boldsymbol{\tau} + C \operatorname{tr}(\boldsymbol{\tau})^2\mathbf{I}_d + A(\boldsymbol{\tau})^2, \quad (1.14)$$

where  $\lambda$  and  $\mu$  are the standard Lamé parameter, and  $A, B, C \in \mathbb{R}$  are the second-order moduli. We note that the hyperelastic stored energy function associated to the stress-strain relation (1.14) has third-order terms and in general is not polyconvex, *i.e.* convex with respect to each of the strain invariants.

## 1.2 Discretization

In order to discretize the Biot and nonlinear poroelasticity problems presented in the previous sections, we consider a coupling of the Hybrid High-Order (HHO) and discontinuous Galerkin (dG) method. Concerning the time discretization, we consider backward differentiation formulas (especially the implicit Euler method), which are the simplest and most widely used methods in the literature. We show that this construction yields an inf-sup stable hydro-mechanical coupling, which is a crucial property to ensure robustness for small time steps combined with small permeabilities. Additionally, the analysis presented in Chapter 2 and 4 allows to derive error estimates that are robust in the limit  $c_0 \rightarrow 0$  which, as noted in Section 1.2.1 and in [164, Section 5.2], is expected to mitigate the problem of the nonphysical pressure oscillations which can sometimes arise in the numerical simulation of poroelastic systems. Before discussing the main features of the discretization methods considered in the next chapters, we give an overview of the numerical issues related to the discretization of poroelasticity problems.

### 1.2.1 Numerical issues

Several difficulties have to be accounted for in the design of a discretization methods for problem (1.7). These issues have three origins: the discretization of the elasticity operator, the (possibly saddle-point) coupling between the flow and the mechanics, and the rough variations of the permeability coefficient that the Darcy operator has to accomodate.

First, we remark that the elasticity operator has to be carefully engineered in order to ensure stability expressed by a discrete version of Korn's inequality. Conforming finite elements naturally yields coercive discretizations, but this is not necessarily the case when considering nonconforming approximations. For instance, it is well known that the Crouzeix–Raviart space (spanned by piecewise affine functions that are continuous at the midpoint of mesh interfaces) does not fulfill a discrete Korn's inequality. Another numerical issue concerning the discretization of elasticity problems is the so-called locking phenomenon. When the dilatation modulus  $\lambda$  in (1.6) is very large, which corresponds to a quasi-incompressible material, results of poor quality can be obtained. More specifically, it can be observed that

the material deforms as if it were much stiffer. In other words, it appears to lock. The key point to ensure that a numerical method is locking-free is to establish an error estimate with a multiplicative constant not blowing up in the limit  $\lambda \rightarrow \infty$ , *i.e.* being able to prove uniform convergence with respect to  $\lambda$ .

The stability of the saddle-point mechanics-flow coupling is closely related to the elasticity locking phenomenon. Indeed, for both of them, the difficulty lies in the approximation of the divergence operator. In the linear elasticity context, locking can be handled by projecting the divergence operator onto a discrete (pressure) space which satisfies an inf-sup condition when coupled with the displacement approximation space. In the context of a poroelastic displacement-pressure coupling, stability can thus be obtained by considering a discrete pressure belonging to that space, which is in fact equivalent to projecting the discrete divergence operator onto the pressure space in the coupling term. From a mathematical point of view, the inf-sup condition yields an estimate in the  $L^\infty((0, t_F); L^2(\Omega))$ -norm on the discrete pore pressure which is independent of  $\underline{\kappa}^{-1}$ . It is of some importance to note that inf-sup stability is not strictly needed in the compressible case (*i.e.* with  $c_0 > 0$ ). As a matter of fact, the presence of the term  $c_0 d_t p$  in the left-hand side of the fluid balance Equation (1.7b) directly yields a discrete  $L^\infty((0, t_F); L^2(\Omega))$  estimate on the pore pressure which does not depend neither on  $\underline{\kappa}^{-1}$  nor on  $\lambda$ . An investigation of the role of the inf-sup condition in the context of finite element discretizations of linear poroelasticity can be found, *e.g.*, in [154, 155] and in [165, 166].

The problem of spurious spatial oscillations of the pore pressure is actually more involved than a simple saddle-point coupling issue. The difficulty comes from the fact that, in very early times (or when the permeability is low), the pressure is quasi- $L^2(\Omega)$  as the diffusion term gives an almost vanishing contribution. However, as soon as  $t > 0$ , boundary conditions are imposed on the pore pressure hence giving necessarily to this latter a  $H^1(\Omega)$  regularity. When  $c_0 = 0$ , if no discrete inf-sup condition holds, the only control of pressure is given by the diffusion term, which is almost inexistent in early times. Spurious spatial oscillations then arise. If an inf-sup condition holds, then it yields a control of the pressure approximation, hence reducing the oscillations. However, it has been recently pointed out in [176] that, even for discretization methods leading to an inf-sup stable discretization of the Stokes problem in the steady case, pressure oscillations can arise owing to a lack of monotonicity of the discrete operator.

## 1.2.2 Polyhedral element methods

In recent years, a large effort has been devoted to the development and analysis of numerical methods that apply to more general meshes than simplicial or Cartesian ones. In the context of poroelasticity, this necessity arises from the presence of various layers and fractures in the porous medium. The discretization of the domain may also include degenerate elements (as in the near wellbore region in reservoir modelling, *e.g.* Figure 1.2) and nonconforming interfaces (accounting for the presence of cracks or resulting from local mesh refinement). In the context of structural mechanics, discretization methods supporting polyhedral meshes and nonconforming interfaces can be useful for several reasons including, *e.g.*, the use of hanging nodes for contact and interface elasticity problems, and the greater robustness to mesh

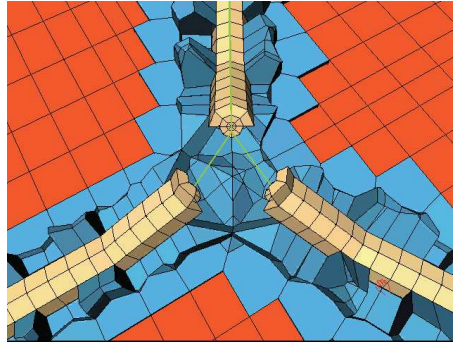


Figure 1.2: Example of polyhedral mesh with degenerate elements in reservoir applications

distorsion and fracture. Moreover, polyhedral methods allow for simple mesh refinement and coarsening procedures for adaptivity.

In this dissertation, we focus on two families of polyhedral discretization methods: the HHO and dG methods. Both are nonconforming in the context of primal formulations of elliptic problems, since no continuity condition is imposed between neighboring mesh elements. They are also high-order methods, because they allow to increase the space approximation order. The use of high-order methods can classically accelerate the convergence in the presence of regular exact solutions or when combined with local mesh refinement. Additionally, the construction of discrete problems using the HHO and dG discretizations is valid for arbitrary space dimension and this enables dimension-independent implementations. In what follows, we consider the HHO method of [33, 74] for the discretization of the (possibly nonlinear) elasticity operator in (1.7a) and the Symmetric Weighted Interior Penalty (SWIP) dG method of [77, 93] for the Darcy operator in (1.7b).

Discontinuous Galerkin methods can be viewed as finite element methods allowing for discontinuities in the discrete trial and test spaces. Localizing test functions to single mesh elements, they can also be viewed as finite volume methods in which the approximate solution is represented on each mesh element by a polynomial function and not only by a constant function. Allowing the approximate solution to be only piecewise continuous offers a substantial amount of flexibility in the mesh design. A unified analysis of dG methods for elliptic problems can be found in [42] and in [76, Chapter 4]. The fundamental strategy to approximate heterogeneous diffusion (such as Darcy flow in porous media) problem using dG methods is to penalize interface and boundary jumps using a diffusion-dependent parameter scaling as the harmonic mean of the normal component of the diffusion tensor. Indeed, using the harmonic mean to penalize will turn out to tune automatically the amount of penalty. Moreover, as pointed out in Chapter 4 and in [77] the proposed method is robust with respect to the spatial variations of the permeability coefficient, with constants in the error estimates having a mild dependance on the heterogeneity ratio.

Hybrid High-Order methods, introduced in [78] and [74], rely on primal formulations of elliptic problems and lead to a symmetric, positive definite system matrices. The methods can be deployed on general polyhedral meshes. The degrees of freedom are polynomials of the same order at mesh elements and faces: face-based discrete unknowns establish inter-elements connections at interfaces and can be used to strongly enforce essential boundary conditions

at boundary faces, element-based discrete unknowns are intermediate variables which can be eliminated from the global system by static condensation, as detailed in [57, Section 2.4] and in Section 2.5. The design of these methods proceeds in two steps: (i) the discrete reconstructions of differential operators hinging on the solution of inexpensive local problems inside each element and (ii) the definition of a least-squares stabilization that weakly enforces the matching of element- and face-based discrete unknowns. The definitions of the differential reconstruction operators are based on discrete counterparts of integrations by parts. Even if local problems related to the reconstruction operators have to be solved, the numerical results of [74] indicate that the associated cost becomes marginal with respect to the cost of solving the global systems as the number of discrete unknowns is increased. Moreover, HHO methods are computationally effective, since the use of face unknowns yields to compact stencils, and they can be efficiently implemented thanks to the possibility of statically condensing a large subset of the unknowns. The HHO discretization studied in Chapter 3 for nonlinear elasticity problems is inspired by recent works on Leray–Lions operators [69, 70], where the authors show that the method is robust with respect to nonlinearities. The robustness of the method with respect to locking has been proved in [74, Section 5].

Other locking-free polyhedral schemes have been proposed for the discretization of linear elasticity problems. A non-exhaustive list includes: (i) the Hybridizable discontinuous Galerkin (HDG) method of [185], (ii) the Mimetic Finite Difference (MFD) scheme of [19], (iii) the Virtual Element Method (VEM) of [20], (iv) the generalized Crouzeix–Raviart scheme introduced in [80], and (v) the Weak Galerkin (WG) method of [197]. Concerning the nonlinear version of the problem, discretization methods supporting general meshes have also been considered in [21] and [54], where the authors propose a low-order VEM scheme for small and finite deformations, respectively. The HHO method has been applied for the discretization of finite deformations elasticity problems in [1]. To the best of our knowledge, the existing literature on the approximation of poroelasticity problem on general polyhedral meshes is more scarce. We can cite the vectorial Multi-Point Flux Approximation methods of [160], the Hybrid Finite Volume discretization considered in [138], and the very recent HDG method of [101].

### 1.3 Uncertainty quantification

In numerical simulation, accounting for uncertainties in input data (such as physical constants, boundary and initial conditions, external forcing, and geometry) is a crucial issue, especially in risk analysis, safety, and design. Risk assessment in poroelasticity is a critical challenge for a wide range of vital resource management applications. Recent public concerns over induced seismicity [63, 113] and groundwater contamination [182] underscore the need to quantify the probability of harmful events associated with subsurface flow and deformation. Of particular interest is the prediction of critical stresses which can compromise the caprock integrity [162] and wellbore stability [153] in storage reservoirs. The key to successful risk assessment in these contexts is the ability to predict the probability of critical events based on some empirical approximation of the material parameters and source term variability. For this reason there is a particular need to combine numerical methods for poromechanical

simulation with efficient techniques for quantifying uncertainties.

Uncertainty Quantification (UQ) methods have been developed in the last decades to take into account the effect of random input parameters on the quantity of interest. They enable to obtain information on the model output that is richer than in a deterministic context, since they provide the statistical moments and the probability distribution. This information makes the comparison with experimental observation easier and facilitate the evaluation of the validity of the physical models. Indeed, UQ methods provide confidence intervals in computed predictions and identify the uncertain parameters that should be measured or controlled with more accuracy because they have the most significant impact on the solution. Additionally, they allow to assess the limits of predictability and the level of reliability that can be attached to numerical simulations.

Among the techniques designed for UQ in numerical models, the stochastic spectral methods have received a considerable attention. The core idea of these methods is the decomposition of random quantities on suitable approximation bases such as the Karhunen–Loève [142, 180] or the Polynomial Chaos [112, 137] expansions. The former represents the random fields as a linear combination of an infinite number of uncorrelated random variables, while the latter uses polynomial expansions in terms of independent random variables. Their main interest is that they provide a complete probabilistic description of the uncertain solution.

In Chapter 5, in order to investigate the effect of uncertainty in poroelasticity problems, we rely on a Polynomial Chaos (PC) approach. The fundamental concept on which PC decompositions are based is to regard uncertainty as generating a new dimension and the solution as being dependent on this dimension. A convergent expansion along the new dimension is then sought in terms of a set of orthogonal basis functions, whose coefficients can be used to characterize the uncertainty. The motivation behind PC approaches includes: (i) the suitability to model expressed in terms of partial differential equations, (ii) the ability to deal with situations exhibiting steep nonlinear dependence of the solution on random model data, and (iii) the promise of obtaining efficient and accurate estimates of uncertainty. In addition, the provided information is given in a format that can be readily exploited to probe the dependence of specific observables on particular components of the input data.

## 1.4 Plan of the manuscript

The rest manuscript is organized as follows.

In Chapter 2, published in *SIAM Journal on Scientific Computing* (cf. [30]), we introduce a novel algorithm for the Biot problem based on a HHO discretization of the mechanics and a Symmetric Weighted Interior Penalty dG discretization of the flow. The method has several assets, including, in particular, the validity in two and three space dimensions, inf-sup stability, and the support of general polyhedral meshes, nonmatching interfaces, and arbitrary space approximation order. Additionally, the resolution cost can be reduced by statically condensing a large subset of the unknowns. Our analysis delivers stability and error estimates that hold also when the constrained specific storage coefficient vanishes, and shows that the constants have only a mild dependence on the heterogeneity of the permeability coefficient. We discuss implementation details and provide numerical tests demonstrating



the performance of the method. In particular, we numerically check the robustness of the method with respect to pressure spurious oscillations. Finally, we show that the scheme is locally conservative on the primal mesh, a desirable property for practitioners and key for a posteriori estimates based on equilibrated fluxes.

In **Chapter 3**, published in *SIAM Journal on Numerical Analysis* (cf. [33]), we propose and analyze a novel HHO discretization of a class of (linear and) nonlinear elasticity models in the small deformation regime which are of common use in solid mechanics. The method satisfies a local principle of virtual work inside each mesh element, with interface tractions that obey the law of action and reaction. A complete analysis covering very general stress-strain laws is carried out. In particular, we prove the existence of a discrete solution and we identify a strict monotonicity assumption on the stress-strain law which ensures uniqueness. Convergence to minimal regularity solutions  $\mathbf{u} \in H_0^1(\Omega; \mathbb{R}^d)$  is proved using a compactness argument. An optimal energy-norm error estimate in  $h^{k+1}$  is then proved under the additional conditions of Lipschitz continuity and strong monotonicity on the stress-strain law. The performance of the method is extensively investigated on a complete panel of model problems using stress-strain laws corresponding to real materials. The numerical tests show that the method is robust with respect to strong nonlinearities.

In **Chapter 4**, building on the material of the previous chapters, we construct and analyze a coupled HHO-dG discretization method for the nonlinear poroelasticity problem. The chapter is submitted for publication (see [35] for the preprint version). Compared to the method proposed in Chapter 2 for the linear poroelasticity problem, there are two main differences in the design. First, the discrete symmetric gradient sits in the full space of tensor-valued polynomials, as opposed to symmetric gradients of vector-valued polynomials. As shown in Chapter 3, this modification is required for the convergence analysis in the presence of nonlinear stress-strain laws. Second, the right-hand side of the discrete problem is obtained by taking the average in time of the loading force  $\mathbf{f}$  and fluid source  $g$  between two consecutive time steps. This modification allows us to prove stability and optimal error estimates without any additional time regularity assumptions on data. Moreover, the results holds for both nonzero and vanishing storage coefficients. In this chapter we also give a new simple proof of discrete Korn's inequality, not requiring particular geometrical assumptions on the mesh.

In **Chapter 5**, we consider the numerical solution of the Biot problem with random poroelastic coefficients in the context of uncertainty quantification. The chapter collects part of the ongoing work carried out during the internship at the BRGM that took place from January to September 2018. The uncertainty is modelled with a finite set of parameters with prescribed probability distribution. We present the variational formulation of the stochastic partial differential system and establish its well-posedness. We then discuss the approximation of the parameter-dependent problem by non-intrusive techniques based on Polynomial Chaos decompositions. We specifically focus on sparse spectral projection methods which essentially amount to performing an ensemble of deterministic model simulations to estimate the expansion coefficients. We numerically investigate the convergence of the probability error of the PC approximation with respect to the level of the sparse grid. Finally, we perform a sensitivity analysis to assess the propagation of the input uncertainty on the solutions

considering an injection test and a traction problem.

In **Appendix A**, published in the conference book *Finite Volumes for Complex Applications VIII* (cf. [34]), we present a variant of the HHO-dG algorithm for the quasi-static nonlinear poroelasticity problem studied in Chapter 4. In particular, this appendix provides numerical tests demonstrating the convergence of the method considering the Hencky-Mises non-linear behavior law.

**Appendix B** is an excerpt of a complementary work that has been published in *Computational Methods in Applied Mathematics* (cf. [36]). I decided not to include the full version to give a more consistent structure to the manuscript. Therein we consider hyperelastic problems and their numerical solution using a conforming finite element discretization and iterative linearization algorithms. For these problems, we present equilibrated, weakly symmetric,  $H(\text{div})$ -conforming stress tensor reconstructions, obtained from local problems on patches around vertices using the Arnold–Falk–Winther finite element spaces. We distinguish two stress reconstructions, one for the discrete stress and one representing the error of the iterative linearization algorithm. The reconstructions are independent of the mechanical behavior law. Based on these stress tensor reconstructions, we derive an a posteriori error estimate distinguishing the discretization, linearization, and quadrature error estimates, and propose an adaptive algorithm balancing these different error sources. We prove the efficiency of the estimate, and confirm it on a numerical test with analytical solution. We then apply the adaptive algorithm to a more application-oriented test.



# Chapter 2

---

## The Biot problem

---

This chapter has been published in the following peer-reviewed journal (see [30]):

**SIAM Journal on Scientific Computing,**  
Volume 38, Issue 3, 2016, Pages A1508–A1537.

### Contents

---

<b>2.1</b>	<b>Introduction</b>	<b>18</b>
<b>2.2</b>	<b>Discretization</b>	<b>20</b>
2.2.1	Mesh and notation	20
2.2.2	Linear elasticity operator	22
2.2.3	Darcy operator	24
2.2.4	Hydro-mechanical coupling	25
2.2.5	Formulation of the method	26
<b>2.3</b>	<b>Stability analysis</b>	<b>27</b>
<b>2.4</b>	<b>Error analysis</b>	<b>31</b>
2.4.1	Projection	31
2.4.2	Error equations	33
2.4.3	Convergence	33
<b>2.5</b>	<b>Implementation</b>	<b>37</b>
<b>2.6</b>	<b>Numerical tests</b>	<b>40</b>
2.6.1	Convergence	40
2.6.2	Barry and Mercer’s test case	42
<b>2.7</b>	<b>Appendix: Flux formulation</b>	<b>44</b>

---

## 2.1 Introduction

We consider in this chapter the quasi-static Biot's consolidation problem describing Darcian flow in a deformable saturated porous medium. Our original motivation comes from applications in geosciences, where the support of general polyhedral meshes is crucial, e.g., to handle nonconforming interfaces arising from local mesh adaptation or Voronoi elements in the near wellbore region when modelling petroleum extraction. Let  $\Omega \subset \mathbb{R}^d$ ,  $1 \leq d \leq 3$ , denote a bounded connected polyhedral domain with boundary  $\partial\Omega$  and outward normal  $\mathbf{n}$ . For a given finite time  $t_F > 0$ , volumetric load  $\mathbf{f}$ , fluid source  $g$ , the Biot problem consists in finding a vector-valued displacement field  $\mathbf{u}$  and a scalar-valued pore pressure field  $p$  solution of

$$-\nabla \cdot \boldsymbol{\sigma}(\mathbf{u}) + \alpha \nabla p = \mathbf{f} \quad \text{in } \Omega \times (0, t_F), \quad (2.1a)$$

$$c_0 d_t p + \nabla \cdot (\alpha d_t \mathbf{u}) - \nabla \cdot (\boldsymbol{\kappa} \nabla p) = g \quad \text{in } \Omega \times (0, t_F), \quad (2.1b)$$

where  $c_0 \geq 0$  and  $\alpha > 0$  are real numbers corresponding to the constrained specific storage and Biot–Willis coefficients, respectively,  $\boldsymbol{\kappa}$  is a real-valued permeability field such that  $\underline{\kappa} \leq \boldsymbol{\kappa} \leq \bar{\kappa}$  a.e. in  $\Omega$  for given real numbers  $0 < \underline{\kappa} \leq \bar{\kappa}$ , and the Cauchy stress tensor is given by

$$\boldsymbol{\sigma}(\mathbf{u}) := 2\mu \nabla_s \mathbf{u} + \lambda \mathbf{I}_d \nabla \cdot \mathbf{u},$$

with real numbers  $\lambda \geq 0$  and  $\mu > 0$  corresponding to Lamé's parameters,  $\nabla_s$  denoting the symmetric part of the gradient operator applied to vector-valued fields, and  $\mathbf{I}_d$  denoting the identity matrix of  $\mathbb{R}^{d \times d}$ . Equations (2.1a) and (2.1b) express, respectively, the mechanical equilibrium and the fluid mass balance. We consider, for the sake of simplicity, the following homogeneous boundary conditions:

$$\mathbf{u} = \mathbf{0} \quad \text{on } \partial\Omega \times (0, t_F), \quad (2.1c)$$

$$\boldsymbol{\kappa} \nabla p \cdot \mathbf{n} = 0 \quad \text{on } \partial\Omega \times (0, t_F). \quad (2.1d)$$

Initial conditions are set prescribing  $\mathbf{u}(\cdot, 0) = \mathbf{u}^0$  and, if  $c_0 > 0$ ,  $p(\cdot, 0) = p^0$ . In the incompressible case  $c_0 = 0$ , we also need the following compatibility condition on  $g$ :

$$\int_{\Omega} g(\cdot, t) = 0 \quad \forall t \in (0, t_F), \quad (2.1e)$$

as well as the following zero-average constraint on  $p$ :

$$\int_{\Omega} p(\cdot, t) = 0 \quad \forall t \in (0, t_F). \quad (2.1f)$$

For the derivation of the Biot model we refer to the seminal work of Terzaghi [190] and Biot [27, 29]. A theoretical study of problem (2.1) can be found in [181]. For the precise regularity assumptions on the data and on the solution under which our a priori bounds and convergence estimates are derived, we refer to Lemma 2.7 and Theorem 2.12, respectively.

A few simplifications are made to keep the exposition as simple as possible while still retaining all the principal difficulties. For the Biot–Willis coefficient we take

$$\alpha = 1,$$

an assumption often made in practice. For the scalar-valued permeability  $\kappa$ , we assume that it is piecewise constant on a partition  $P_\Omega$  of  $\Omega$  into bounded open polyhedra. The treatment of more general permeability coefficients can be done following the ideas of [77]. Also, more general boundary conditions than (2.1c)–(2.1d) can be considered up to minor modifications.

Our focus is here on a novel space discretization for the Biot problem (standard choices are made for the time discretization). Several difficulties have to be accounted for in the design of the space discretization of problem (2.1): in the context of nonconforming methods, the linear elasticity operator has to be carefully engineered to ensure stability expressed by a discrete counterpart of the Korn’s inequality; the Darcy operator has to accommodate rough variations of the permeability coefficient; the choice of discrete spaces for the displacement and the pressure must satisfy an inf-sup condition to contribute reducing spurious pressure oscillations for small time steps combined with small permeabilities when  $c_0 = 0$ . An investigation of the role of the inf-sup condition in the context of finite element discretizations can be found, e.g., in Murad and Loula [154, 155]. A recent work of Rodrigo, Gaspar, Hu, and Zikatanov [176] has pointed out that, even for discretization methods leading to an inf-sup stable discretization of the Stokes problem in the steady case, pressure oscillations can arise owing to a lack of monotonicity. Therein, the authors suggest that stabilizing is possible by adding to the mass balance equation an artificial diffusion term with coefficient proportional to  $h^2/\tau$  (with  $h$  and  $\tau$  denoting, respectively, the spatial and temporal meshsizes). However, computing the exact amount of stabilization required is in general feasible only in 1 space dimension.

Several space discretization methods for the Biot problem have been considered in the literature. Finite element discretizations are discussed, e.g., in the monograph of Lewis and Schrefler [140]; cf. also references therein. A finite volume discretization for the three-dimensional Biot problem with discontinuous physical coefficients is considered by Naumovich [156]. In [165, 166], Phillips and Wheeler propose and analyze an algorithm that models displacements with continuous elements and the flow with a mixed method. In [167], the same authors also propose a different method where displacements are instead approximated using discontinuous Galerkin methods. In [200], Wheeler, Xue and Yotov study the coupling of multipoint flux discretization for the flow with a discontinuous Galerkin discretization of the displacements. While certainly effective on matching simplicial meshes, discontinuous Galerkin discretizations of the displacements usually do not allow to prove inf-sup stability on general polyhedral meshes.

In this chapter, we propose a novel space discretization of problem (2.1) where the linear elasticity operator is discretized using the Hybrid High-Order (HHO) method of [74] (c.f. also [71, 73, 78]), while the flow relies on the Symmetric Weighted Interior Penalty (SWIP) discontinuous Galerkin method of [77], see also [76, Chapter 4]. The proposed method has several assets: (i) It delivers an inf-sup stable discretization on general meshes including, e.g., polyhedral elements and nonmatching interfaces; (ii) it allows to increase the space approximation order to accelerate convergence in the presence of (locally) regular solutions;

(iii) it is locally conservative on the primal mesh, a desirable property for practitioners and key for a posteriori estimates based on equilibrated fluxes; (iv) it is robust with respect to the spatial variations of the permeability coefficient, with constants in the error estimates that depend on the square root of the heterogeneity ratio; (v) it is (relatively) inexpensive: at the lowest order, after static condensation of element unknowns for the displacement, we have 4 (resp. 9) unknowns per face for the displacements + 3 (resp. 4) unknowns per element for the pore pressure in 2d (resp. 3d). Finally, the proposed construction is valid for arbitrary space dimension, a feature which can be exploited in practice to conceive dimension-independent implementations.

The material is organized as follows. In Section 2.2, we introduce the discrete setting and formulate the method. In Section 2.3, we derive a priori bounds on the exact solution for regular-in-time volumetric load and mass source. The convergence analysis of the method is carried out in Section 2.4. Implementation details are discussed in Section 2.5, while numerical tests proposed in Section 2.6. Finally, in Appendix 2.7, we investigate the local conservation properties of the method by identifying computable conservative normal tractions and mass fluxes.

## 2.2 Discretization

In this section we introduce the assumptions on the mesh, define the discrete counterparts of the elasticity and Darcy operators and of the hydro-mechanical coupling terms, and formulate the discretization method.

### 2.2.1 Mesh and notation

Denote by  $\mathcal{H} \subset \mathbb{R}_*^+$  a countable set of meshsizes having 0 as its unique accumulation point. Following [76, Chapter 1], we consider  $h$ -refined spatial mesh sequences  $(\mathcal{T}_h)_{h \in \mathcal{H}}$  where, for all  $h \in \mathcal{H}$ ,  $\mathcal{T}_h$  is a finite collection of nonempty disjoint open polyhedral elements  $T$  such that  $\overline{\Omega} = \bigcup_{T \in \mathcal{T}_h} \overline{T}$  and  $h = \max_{T \in \mathcal{T}_h} h_T$  with  $h_T$  standing for the diameter of the element  $T$ . We assume that mesh regularity holds in the following sense: For all  $h \in \mathcal{H}$ ,  $\mathcal{T}_h$  admits a matching simplicial submesh  $\mathfrak{T}_h$  and there exists a real number  $\varrho > 0$  independent of  $h$  such that, for all  $h \in \mathcal{H}$ , (i) for all simplex  $S \in \mathfrak{T}_h$  of diameter  $h_S$  and inradius  $r_S$ ,  $\varrho h_S \leq r_S$  and (ii) for all  $T \in \mathcal{T}_h$ , and all  $S \in \mathfrak{T}_h$  such that  $S \subset T$ ,  $\varrho h_T \leq h_S$ . A mesh sequence with this property is called regular. It is worth emphasizing that the simplicial submesh  $\mathfrak{T}_h$  is just an analysis tool, and it is not used in the actual construction of the discretization method. These assumptions are essentially analogous to those made in the context of other recent methods supporting general meshes; cf., e.g., [20, Section 2.2] for the Virtual Element method. For a collection of useful geometric and functional inequalities that hold on regular mesh sequences we refer to [76, Chapter 1] and [69].

*Remark 2.1* (Face degeneration). The above regularity assumptions on the mesh imply that the diameter of the mesh faces is uniformly comparable to that of the cell(s) they belong to, i.e., face degeneration is not allowed. Face degeneration has been considered, on the other hand, in [47] in the context of interior penalty discontinuous Galerkin methods. One

could expect that this framework could be used herein while adapting accordingly the penalty strategy (2.11) and (2.19). This point lies out of the scope of the present work and will be inspected in the future.

To avoid dealing with jumps of the permeability inside elements, we additionally assume that, for all  $h \in \mathcal{H}$ ,  $\mathcal{T}_h$  is compatible with the known partition  $P_\Omega$  on which the diffusion coefficient  $\kappa$  is piecewise constant, so that jumps can only occur at interfaces.

We define a face  $F$  as a hyperplanar closed connected subset of  $\overline{\Omega}$  with positive  $(d-1)$ -dimensional Hausdorff measure and such that (i) either there exist  $T_1, T_2 \in \mathcal{T}_h$  such that  $F \subset \partial T_1 \cap \partial T_2$  (with  $\partial T_i$  denoting the boundary of  $T_i$ ) and  $F$  is called an interface or (ii) there exists  $T \in \mathcal{T}_h$  such that  $F \subset \partial T \cap \partial \Omega$  and  $F$  is called a boundary face. Interfaces are collected in the set  $\mathcal{F}_h^i$ , boundary faces in  $\mathcal{F}_h^b$ , and we let  $\mathcal{F}_h := \mathcal{F}_h^i \cup \mathcal{F}_h^b$ . The diameter of a face  $F \in \mathcal{F}_h$  is denoted by  $h_F$ . For all  $T \in \mathcal{T}_h$ ,  $\mathcal{F}_T := \{F \in \mathcal{F}_h : F \subset \partial T\}$  denotes the set of faces contained in  $\partial T$  and, for all  $F \in \mathcal{F}_T$ ,  $\mathbf{n}_{TF}$  is the unit normal to  $F$  pointing out of  $T$ . For a regular mesh sequence, the maximum number of faces in  $\mathcal{F}_T$  can be bounded by an integer  $N_\partial$  uniformly in  $h$ . For each interface  $F \in \mathcal{F}_h^i$ , we fix once and for all the ordering for the elements  $T_1, T_2 \in \mathcal{T}_h$  such that  $F \subset \partial T_1 \cap \partial T_2$  and we let  $\mathbf{n}_F := \mathbf{n}_{T_1, F}$ . For a boundary face, we simply take  $\mathbf{n}_F = \mathbf{n}$ , the outward unit normal to  $\Omega$ .

For integers  $l \geq 0$  and  $s \geq 1$ , we denote by  $\mathbb{P}_d^l(\mathcal{T}_h)$  the space of fully discontinuous piecewise polynomial functions of total degree  $\leq l$  on  $\mathcal{T}_h$  and by  $H^s(\mathcal{T}_h)$  the space of functions in  $L^2(\Omega)$  that lie in  $H^s(T)$  for all  $T \in \mathcal{T}_h$ . The notation  $H^s(P_\Omega)$  will also be used with obvious meaning. Under the mesh regularity assumptions detailed above, using [76, Lemma 1.40] together with the results of [88], one can prove that there exists a real number  $C_{\text{app}}$  depending on  $\varrho$  and  $l$ , but independent of  $h$ , such that, denoting by  $\pi_h^l$  the  $L^2$ -orthogonal projector on  $\mathbb{P}_d^l(\mathcal{T}_h)$ , the following holds: For all  $s \in \{1, \dots, l+1\}$  and all  $v \in H^s(\mathcal{T}_h)$ ,

$$|v - \pi_h^l v|_{H^m(\mathcal{T}_h)} \leq C_{\text{app}} h^{s-m} |v|_{H^s(\mathcal{T}_h)} \quad \forall m \in \{0, \dots, s-1\}. \quad (2.2)$$

For an integer  $l \geq 0$ , we consider the space

$$C^l(V) := C^l([0, t_F]; V),$$

spanned by  $V$ -valued functions that are  $l$  times continuously differentiable in the time interval  $[0, t_F]$ . The space  $C^0(V)$  is a Banach space when equipped with the norm  $\|\varphi\|_{C^0(V)} := \max_{t \in [0, t_F]} \|\varphi(t)\|_V$ , and the space  $C^l(V)$  is a Banach space when equipped with the norm  $\|\varphi\|_{C^l(V)} := \max_{0 \leq m \leq l} \|d_t^m \varphi\|_{C^0(V)}$ . For the time discretization, we consider a uniform mesh of the time interval  $(0, t_F)$  of step  $\tau := t_F/N$  with  $N \in \mathbb{N}^*$ , and introduce the discrete times  $t^n := n\tau$  for all  $0 \leq n \leq N$ . For any  $\varphi \in C^l(V)$ , we set  $\varphi^n := \varphi(t^n) \in V$ , and we introduce the backward differencing operator  $\delta_t$  such that, for all  $1 \leq n \leq N$ ,

$$\delta_t \varphi^n := \frac{\varphi^n - \varphi^{n-1}}{\tau} \in V. \quad (2.3)$$

In what follows, for  $X \subset \overline{\Omega}$ , we respectively denote by  $(\cdot, \cdot)_X$  and  $\|\cdot\|_X$  the standard inner product and norm in  $L^2(X)$ , with the convention that the subscript is omitted whenever  $X = \Omega$ . The same notation is used in the vector- and tensor-valued cases. For the sake of



brevity, throughout the chapter, we will often use the notation  $a \lesssim b$  for the inequality  $a \leq Cb$  with generic constant  $C > 0$  independent of  $h$ ,  $\tau$ ,  $c_0$ ,  $\lambda$ ,  $\mu$ , and  $\kappa$ , but possibly depending on  $\varrho$  and the polynomial degree  $k$ . We will name generic constants only in statements or when this helps to follow the proofs.

## 2.2.2 Linear elasticity operator

The discretization of the linear elasticity operator is based on the Hybrid High-Order method of [74]. Let a polynomial degree  $k \geq 1$  be fixed. The degrees of freedom (DOFs) for the displacement are collected in the space

$$\underline{U}_h^k := \left\{ \prod_{T \in \mathcal{T}_h} \mathbb{P}_d^k(T)^d \right\} \times \left\{ \prod_{F \in \mathcal{F}_h} \mathbb{P}_{d-1}^k(F)^d \right\}.$$

For a generic collection of DOFs in  $\underline{U}_h^k$  we use the notation  $\underline{v}_h := ((v_T)_{T \in \mathcal{T}_h}, (v_F)_{F \in \mathcal{F}_h})$ . We also denote by  $v_h$  (not underlined) the function of  $\mathbb{P}_d^k(\mathcal{T}_h)^d$  such that  $v_h|_T = v_T$  for all  $T \in \mathcal{T}_h$ . The restrictions of  $\underline{U}_h^k$  and  $\underline{v}_h$  to an element  $T$  are denoted by  $\underline{U}_T^k$  and  $\underline{v}_T = (v_T, (v_F)_{F \in \mathcal{F}_T})$ , respectively. For further use, we define the reduction map  $\underline{I}_h^k : H^1(\Omega)^d \rightarrow \underline{U}_h^k$  such that, for all  $v \in H^1(\Omega)^d$ ,

$$\underline{I}_h^k v = ((\pi_T^k v)_{T \in \mathcal{T}_h}, (\pi_F^k v)_{F \in \mathcal{F}_h}), \quad (2.4)$$

where  $\pi_T^k$  and  $\pi_F^k$  denote the  $L^2$ -orthogonal projectors on  $\mathbb{P}_d^k(T)$  and  $\mathbb{P}_{d-1}^k(F)$ , respectively. For all  $T \in \mathcal{T}_h$ , the reduction map on  $\underline{U}_T^k$  obtained by a restriction of  $\underline{I}_h^k$  is denoted by  $\underline{I}_T^k$ .

For all  $T \in \mathcal{T}_h$ , we obtain a high-order polynomial reconstruction  $\mathbf{r}_T^{k+1} : \underline{U}_T^k \rightarrow \mathbb{P}_d^{k+1}(T)^d$  of the displacement field by solving the following local pure traction problem: For a given local collection of DOFs  $\underline{v}_T = (v_T, (v_F)_{F \in \mathcal{F}_T}) \in \underline{U}_T^k$ , find  $\mathbf{r}_T^{k+1} \underline{v}_T \in \mathbb{P}_d^{k+1}(T)^d$  such that

$$(\nabla_s \mathbf{r}_T^{k+1} \underline{v}_T, \nabla_s \mathbf{w})_T = (\nabla_s v_T, \nabla_s \mathbf{w})_T + \sum_{F \in \mathcal{F}_T} (v_F - v_T, \nabla_s \mathbf{w} \mathbf{n}_{TF})_F \quad \forall \mathbf{w} \in \mathbb{P}_d^{k+1}(T)^d. \quad (2.5)$$

In order to uniquely define the solution to (2.5), we prescribe the conditions  $\int_T \mathbf{r}_T^{k+1} \underline{v}_T = \int_T v_T$  and  $\int_T \nabla_{ss} \mathbf{r}_T^{k+1} \underline{v}_T = \sum_{F \in \mathcal{F}_T} \int_F \frac{1}{2} (v_F \otimes \mathbf{n}_{TF} - \mathbf{n}_{TF} \otimes v_F)$ , where  $\nabla_{ss}$  denotes the skew-symmetric part of the gradient operator. We also define the global reconstruction of the displacement  $\mathbf{r}_h^{k+1} : \underline{U}_h^k \rightarrow \mathbb{P}_d^{k+1}(\mathcal{T}_h)^d$  such that, for all  $\underline{v}_h \in \underline{U}_h^k$ ,

$$(\mathbf{r}_h^{k+1} \underline{v}_h)|_T = \mathbf{r}_T^{k+1} \underline{v}_T \quad \forall T \in \mathcal{T}_h.$$

The following approximation property is proved in [74, Lemma 2]: For all  $v \in H^1(\Omega)^d \cap H^{k+2}(P_\Omega)^d$ ,

$$\|\nabla_s(\mathbf{r}_h^{k+1} \underline{I}_h^k v - v)\| \lesssim h^{k+1} \|v\|_{H^{k+2}(P_\Omega)^d}. \quad (2.6)$$

We next introduce the discrete divergence operator  $D_T^k : \underline{U}_T^k \rightarrow \mathbb{P}_d^k(T)$  such that, for all

$q \in \mathbb{P}_d^k(T)$

$$(D_T^k \underline{\mathbf{v}}_T, q)_T = (\nabla \cdot \mathbf{v}_T, q)_T + \sum_{F \in \mathcal{F}_T} (\mathbf{v}_F - \mathbf{v}_T, q \mathbf{n}_{TF})_F \quad (2.7a)$$

$$= -(\mathbf{v}_T, \nabla q)_T + \sum_{F \in \mathcal{F}_T} (\mathbf{v}_F, q \mathbf{n}_{TF})_F, \quad (2.7b)$$

where we have used integration by parts to pass to the second line. The divergence operator satisfies the following commuting property: For all  $T \in \mathcal{T}_h$  and all  $\mathbf{v} \in H^1(T)^d$ ,

$$D_T^k I_T^k \mathbf{v} = \pi_T^k(\nabla \cdot \mathbf{v}). \quad (2.8)$$

The local contribution to the discrete linear elasticity operator is expressed by the bilinear form  $a_T$  on  $\underline{\mathbf{U}}_T^k \times \underline{\mathbf{U}}_T^k$  such that, for all  $\underline{\mathbf{w}}_T, \underline{\mathbf{v}}_T \in \underline{\mathbf{U}}_T^k$ ,

$$a_T(\underline{\mathbf{w}}_T, \underline{\mathbf{v}}_T) := 2\mu \left\{ (\nabla_s \mathbf{r}_T^{k+1} \underline{\mathbf{w}}_T, \nabla_s \mathbf{r}_T^{k+1} \underline{\mathbf{v}}_T)_T + s_T(\underline{\mathbf{w}}_T, \underline{\mathbf{v}}_T) \right\} + \lambda (D_T^k \underline{\mathbf{w}}_T, D_T^k \underline{\mathbf{v}}_T)_T, \quad (2.9)$$

where the stabilization bilinear form  $s_T$  is such that

$$s_T(\underline{\mathbf{w}}_T, \underline{\mathbf{v}}_T) := \sum_{F \in \mathcal{F}_T} h_F^{-1} (\Delta_{TF}^k \underline{\mathbf{w}}_T, \Delta_{TF}^k \underline{\mathbf{v}}_T)_F, \quad (2.10)$$

with face-based residual such that, for all  $\underline{\mathbf{w}}_T \in \underline{\mathbf{U}}_T^k$ ,

$$\Delta_{TF}^k \underline{\mathbf{w}}_T := (\pi_F^k \mathbf{r}_T^{k+1} \underline{\mathbf{w}}_T - \mathbf{w}_F) - (\pi_T^k \mathbf{r}_T^{k+1} \underline{\mathbf{w}}_T - \mathbf{w}_T).$$

The global bilinear form  $a_h$  on  $\underline{\mathbf{U}}_h^k \times \underline{\mathbf{U}}_h^k$  is assembled element-wise from local contributions:

$$a_h(\underline{\mathbf{w}}_h, \underline{\mathbf{v}}_h) := \sum_{T \in \mathcal{T}_h} a_T(\underline{\mathbf{w}}_T, \underline{\mathbf{v}}_T). \quad (2.11)$$

To account for the zero-displacement boundary condition (2.1c), we consider the subspace

$$\underline{\mathbf{U}}_{h,D}^k := \left\{ \underline{\mathbf{v}}_h = ((\mathbf{v}_T)_{T \in \mathcal{T}_h}, (\mathbf{v}_F)_{F \in \mathcal{F}_h}) \in \underline{\mathbf{U}}_h^k : \mathbf{v}_F \equiv \mathbf{0} \quad \forall F \in \mathcal{F}_h^b \right\}. \quad (2.12)$$

Define on  $\underline{\mathbf{U}}_h^k$  the discrete strain seminorm

$$\|\underline{\mathbf{v}}_h\|_{\epsilon,h}^2 := \sum_{T \in \mathcal{T}_h} \|\underline{\mathbf{v}}_T\|_{\epsilon,T}^2, \quad \|\underline{\mathbf{v}}_h\|_{a,h}^2 := \|\nabla_s \mathbf{v}_T\|_T^2 + \sum_{F \in \mathcal{F}_T} h_F^{-1} \|\mathbf{v}_F - \mathbf{v}_T\|_F^2. \quad (2.13)$$

It can be proved that  $\|\cdot\|_{\epsilon,h}$  defines a norm on  $\underline{\mathbf{U}}_{h,D}^k$ . Moreover, using [74, Corollary 6], one has the following coercivity and boundedness result for  $a_h$ :

$$\eta^{-1}(2\mu) \|\underline{\mathbf{v}}_h\|_{\epsilon,h}^2 \leq \|\underline{\mathbf{v}}_h\|_{a,h}^2 := a_h(\underline{\mathbf{v}}_h, \underline{\mathbf{v}}_h) \leq \eta(2\mu + d\lambda) \|\underline{\mathbf{v}}_h\|_{\epsilon,h}^2, \quad (2.14)$$

where  $\eta > 0$  is a real number independent of  $h$ ,  $\tau$  and the physical coefficients. Additionally, we know from [74, Theorem 8] that, for all  $\mathbf{w} \in H_0^1(\Omega)^d \cap H^{k+2}(P_\Omega)^d$  such that  $\nabla \cdot \mathbf{w} \in H^{k+1}(P_\Omega)$  and all  $\underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_{h,D}^k$ , the following consistency result holds:

$$\left| a_h(I_h^k \mathbf{w}, \underline{\mathbf{v}}_h) + (\nabla \cdot \boldsymbol{\sigma}(\mathbf{w}), \mathbf{v}_h) \right| \lesssim h^{k+1} \left( 2\mu \|\mathbf{w}\|_{H^{k+2}(P_\Omega)^d} + \lambda \|\nabla \cdot \mathbf{w}\|_{H^{k+1}(P_\Omega)} \right) \|\underline{\mathbf{v}}_h\|_{\epsilon,h}. \quad (2.15)$$

To close this section, we give the following discrete counterpart of Korn's inequality (see 4.3 for the proof).

**Proposition 2.2** (Discrete Korn's inequality). *There is a real number  $C_K > 0$  depending on  $\varrho$  and on  $k$  but independent of  $h$  such that, for all  $\underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_{h,\mathbb{D}}^k$ , recalling that  $\mathbf{v}_h \in \mathbb{P}_d^k(\mathcal{T}_h)^d$  denotes the broken polynomial function such that  $\mathbf{v}_h|_T = \mathbf{v}_T$  for all  $T \in \mathcal{T}_h$ ,*

$$\|\mathbf{v}_h\| \leq C_K d_\Omega \|\underline{\mathbf{v}}_h\|_{\epsilon,h}, \quad (2.16)$$

where  $d_\Omega$  denotes the diameter of  $\Omega$ .

### 2.2.3 Darcy operator

The discretization of the Darcy operator is based on the Symmetric Weighted Interior Penalty method of [77], cf. also [76, Section 4.5]. At each time step, the discrete pore pressure is sought in the broken polynomial space

$$P_h^k := \begin{cases} \mathbb{P}_d^k(\mathcal{T}_h) & \text{if } c_0 > 0, \\ \mathbb{P}_{d,0}^k(\mathcal{T}_h) & \text{if } c_0 = 0, \end{cases} \quad (2.17)$$

where we have introduced the zero-average subspace  $\mathbb{P}_{d,0}^k(\mathcal{T}_h) := \{q_h \in \mathbb{P}_d^k(\mathcal{T}_h) : (q_h, 1) = 0\}$ . For all  $F \in \mathcal{F}_h^i$ , we define the jump and (weighted) average operators such that, for all  $\varphi \in H^1(\mathcal{T}_h)$ , denoting by  $\varphi_T$  and  $\kappa_T$  the restrictions of  $\varphi$  and  $\kappa$  to  $T \in \mathcal{T}_h$ , respectively,

$$[\varphi]_F := \varphi_{T_1} - \varphi_{T_2}, \quad \{\varphi\}_F := \omega_{T_1} \varphi_{T_1} + \omega_{T_2} \varphi_{T_2}, \quad (2.18)$$

where  $\omega_{T_1} = 1 - \omega_{T_2} := \frac{\kappa_{T_2}}{(\kappa_{T_1} + \kappa_{T_2})}$ . Denoting by  $\nabla_h$  the broken gradient on  $H^1(\mathcal{T}_h)$  and letting, for all  $F \in \mathcal{F}_h^i$ ,  $\lambda_{\kappa,F} := \frac{2\kappa_{T_1}\kappa_{T_2}}{(\kappa_{T_1} + \kappa_{T_2})}$ , we define the bilinear form  $c_h$  on  $P_h^k \times P_h^k$  such that, for all  $q_h, r_h \in P_h^k$ ,

$$\begin{aligned} c_h(r_h, q_h) &:= (\kappa \nabla_h r_h, \nabla_h q_h) + \sum_{F \in \mathcal{F}_h^i} \frac{\varsigma \lambda_{\kappa,F}}{h_F} ([r_h]_F, [q_h]_F)_F \\ &\quad - \sum_{F \in \mathcal{F}_h^i} ((\kappa \nabla_h r_h)_F \cdot \mathbf{n}_F, [q_h]_F)_F + ([r_h]_F, \{\kappa \nabla_h q_h\}_F \cdot \mathbf{n}_F)_F, \end{aligned} \quad (2.19)$$

where  $\varsigma > 0$  is a user-defined penalty parameter. The fact that the boundary terms only appear on internal faces in (2.19) reflects the Neumann boundary condition (2.1d). From this point on, we will assume that  $\varsigma > C_{\text{tr}}^2 N_\partial$  with  $C_{\text{tr}}$  denoting the constant from the discrete trace inequality [76, Eq. (1.37)], which ensures that the bilinear form  $c_h$  is coercive (in the numerical tests of Section 2.6, we took  $\varsigma = (N_\partial + 0.1)k^2$ ). Since the bilinear form  $c_h$  is also symmetric, it defines a seminorm on  $P_h^k$ , denoted hereafter by  $\|\cdot\|_{c,h}$  (the map  $\|\cdot\|_{c,h}$  is in fact a norm on  $\mathbb{P}_{d,0}^k(\mathcal{T}_h)$ ).

*Remark 2.3* (Alternative stabilization). To get rid of the dependence of the lower threshold of  $\varsigma$  on  $C_{\text{tr}}$ , one can resort to the BR2 stabilization; c.f. [17] and also [76, Section 5.3.2]. In passing, this stabilization could also contribute to handle face degeneration since the penalty parameter no longer depends on the inverse of the face diameter (cf. Remark 2.1). This topic will make the object of future investigations.

The following known results will be needed in the analysis. Let

$$P_* := \left\{ r \in H^1(\Omega) \cap H^2(P_\Omega) : \boldsymbol{\kappa} \nabla r \cdot \mathbf{n} = 0 \text{ on } \partial\Omega \right\}, \quad P_{*h}^k := P_* + P_h^k.$$

Extending the bilinear form  $c_h$  to  $P_{*h}^k \times P_{*h}^k$ , the following consistency result can be proved adapting the arguments of [76, Chapter 4] to account for the homogeneous Neumann boundary condition (2.1d):

$$\forall r \in P_*, \quad -(\nabla \cdot (\boldsymbol{\kappa} \nabla r), q) = c_h(r, q) \quad \forall q \in P_{*h}. \quad (2.20)$$

Assuming, additionally, that  $r \in H^{k+2}(P_\Omega)$ , as a consequence of [76, Lemma 5.52] together with the optimal approximation properties (2.2) of  $\pi_h^k$  on regular mesh sequences one has,

$$\sup_{q_h \in \mathbb{P}_{d,0}^k(\mathcal{T}_h) \setminus \{0\}} \frac{c_h(r - \pi_h^k r, q_h)}{\|q_h\|_{c,h}} \lesssim \bar{\kappa}^{1/2} h^k \|r\|_{H^{k+1}(P_\Omega)}.$$

### 2.2.4 Hydro-mechanical coupling

The hydro-mechanical coupling is realized by means of the bilinear form  $b_h$  on  $\underline{U}_h^k \times \mathbb{P}_d^k(\mathcal{T}_h)$  such that, for all  $\underline{\mathbf{v}}_h \in \underline{U}_h^k$  and all  $q_h \in \mathbb{P}_d^k(\mathcal{T}_h)$ ,

$$b_h(\underline{\mathbf{v}}_h, q_h) := \sum_{T \in \mathcal{T}_h} b_T(\underline{\mathbf{v}}_T, q_{h|T}), \quad b_T(\underline{\mathbf{v}}_T, q_{h|T}) := -(D_T^k \underline{\mathbf{v}}_T, q_{h|T})_T, \quad (2.21)$$

where  $D_T^k$  is the discrete divergence operator defined by (2.7a). A simple verification shows that, for all  $\underline{\mathbf{v}}_h \in \underline{U}_h^k$  and all  $q_h \in \mathbb{P}_d^k(\mathcal{T}_h)$ ,

$$b_h(\underline{\mathbf{v}}_h, q_h) \lesssim \|\underline{\mathbf{v}}_h\|_{\epsilon,h} \|q_h\|. \quad (2.22)$$

Additionally, using the definition (2.7a) of  $D_T^k$  and (2.12) of  $\underline{U}_{h,D}^k$ , it can be proved that, for all  $\underline{\mathbf{v}}_h \in \underline{U}_{h,D}^k$ , it holds ( $\chi_\Omega$  denotes here the characteristic function of  $\Omega$ ),

$$b_h(\underline{\mathbf{v}}_h, \chi_\Omega) = 0. \quad (2.23)$$

The following inf-sup condition expresses the stability of the hydro-mechanical coupling:

**Lemma 2.4** (Inf-sup condition for  $b_h$ ). *There is a real number  $\beta$  depending on  $\Omega$ ,  $\varrho$  and  $k$  but independent of  $h$  such that, for all  $q_h \in \mathbb{P}_{d,0}^k(\mathcal{T}_h)$ ,*

$$\|q_h\| \leq \beta \sup_{\underline{\mathbf{v}}_h \in \underline{U}_{h,D}^k \setminus \{0\}} \frac{b_h(\underline{\mathbf{v}}_h, q_h)}{\|\underline{\mathbf{v}}_h\|_{\epsilon,h}}. \quad (2.24)$$

*Proof.* Let  $q_h \in \mathbb{P}_{d,0}^k(\mathcal{T}_h)$ . Classically [31], there is  $\mathbf{v}_{q_h} \in H_0^1(\Omega)^d$  such that  $\nabla \cdot \mathbf{v}_{q_h} = q_h$  and  $\|\mathbf{v}_{q_h}\|_{H^1(\Omega)^d} \lesssim \|q_h\|$ . Let  $T \in \mathcal{T}_h$ . Using the  $H^1$ -stability of the  $L^2$ -orthogonal projector (cf., e.g., [69, Corollary 3.7]), it is inferred that

$$\|\nabla_s \pi_T^k \mathbf{v}_{q_h}\|_T \leq \|\nabla \mathbf{v}_{q_h}\|_T.$$

Moreover, for all  $F \in \mathcal{F}_T$ , using the boundedness of  $\pi_F^k$  and the continuous trace inequality of [76, Lemma 1.49] followed by a local Poincaré's inequality for the zero-average function  $(\pi_T^k \mathbf{v}_{q_h} - \mathbf{v}_{q_h})$ , we have

$$h_F^{-1/2} \|\pi_F^k(\pi_T^k \mathbf{v}_{q_h} - \mathbf{v}_{q_h})\|_F \leq h_F^{-1/2} \|\pi_T^k \mathbf{v}_{q_h} - \mathbf{v}_{q_h}\|_F \lesssim \|\nabla \mathbf{v}_{q_h}\|_T.$$

As a result, recalling the definition (2.4) of the local reduction map  $\underline{I}_T^k$  and (2.13) of the strain norm  $\|\cdot\|_{\epsilon, T}$ , it follows that  $\|\underline{I}_T^k \mathbf{v}_{q_h}\|_{\epsilon, T} \lesssim \|\mathbf{v}_{q_h}\|_{H^1(T)^d}$ . Squaring and summing over  $T \in \mathcal{T}_h$  the latter inequality, we get

$$\|\underline{I}_h^k \mathbf{v}_{q_h}\|_{\epsilon, h} \lesssim \|\mathbf{v}_{q_h}\|_{H^1(\Omega)^d} \lesssim \|q_h\|. \quad (2.25)$$

Using (2.25), the commuting property (2.8), and denoting by  $\mathbf{S}$  the supremum in (2.24), one has

$$\|q_h\|^2 = (\nabla \cdot \mathbf{v}_{q_h}, q_h) = \sum_{T \in \mathcal{T}_h} (D_T^k \underline{I}_T^k \mathbf{v}_{q_h}, q_h)_T = -b_h(\underline{I}_h^k \mathbf{v}_{q_h}, q_h) \leq \mathbf{S} \|\underline{I}_h^k \mathbf{v}_{q_h}\|_{\epsilon, h} \lesssim \mathbf{S} \|q_h\|. \quad \square$$

## 2.2.5 Formulation of the method

For all  $1 \leq n \leq N$ , the discrete solution  $(\underline{\mathbf{u}}_h^n, p_h^n) \in \underline{\mathbf{U}}_{h, D}^k \times P_h^k$  at time  $t^n$  is such that, for all  $(\underline{\mathbf{v}}_h, q_h) \in \underline{\mathbf{U}}_{h, D}^k \times \mathbb{P}_d^k(\mathcal{T}_h)$ ,

$$a_h(\underline{\mathbf{u}}_h^n, \underline{\mathbf{v}}_h) + b_h(\underline{\mathbf{v}}_h, p_h^n) = l_h^n(\underline{\mathbf{v}}_h), \quad (2.26a)$$

$$(c_0 \delta_t p_h^n, q_h) - b_h(\delta_t \underline{\mathbf{u}}_h^n, q_h) + c_h(p_h^n, q_h) = (g^n, q_h), \quad (2.26b)$$

where the linear form  $l_h^n$  on  $\underline{\mathbf{U}}_h^k$  is defined as

$$l_h^n(\underline{\mathbf{v}}_h) := (\mathbf{f}^n, \mathbf{v}_h) = \sum_{T \in \mathcal{T}_h} (\mathbf{f}^n, \mathbf{v}_T)_T. \quad (2.27)$$

In petroleum engineering, the usual way to enforce the initial condition is to compute a displacement from an initial (usually hydrostatic) pressure distribution. For a given scalar-valued initial pressure field  $p^0 \in L^2(\Omega)$ , we let  $\widehat{p}_h^0 := \pi_h^k p^0$  and set  $\underline{\mathbf{u}}_h^0 = \widehat{\underline{\mathbf{u}}}_h^0$  with  $\widehat{\underline{\mathbf{u}}}_h^0 \in \underline{\mathbf{U}}_{h, D}^k$  unique solution of

$$a_h(\widehat{\underline{\mathbf{u}}}_h^0, \underline{\mathbf{v}}_h) = l_h^0(\underline{\mathbf{v}}_h) - b_h(\underline{\mathbf{v}}_h, \widehat{p}_h^0) \quad \forall \underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_{h, D}^k. \quad (2.28)$$

If  $c_0 = 0$ , the value of  $\widehat{p}_h^0$  is only needed to enforce the initial condition on the displacement while, if  $c_0 > 0$ , we also set  $p_h^0 = \widehat{p}_h^0$  to initialize the discrete pressure.

*Remark 2.5* (Discrete compatibility condition for  $c_0 = 0$ ). Also when  $c_0 = 0$  it is possible to take the test function  $q_h$  in (2.26b) in the full space  $\mathbb{P}_d^k(\mathcal{T}_h)$  instead of the zero-average subspace  $\mathbb{P}_{d,0}^k(\mathcal{T}_h)$ , since the compatibility condition is verified at the discrete level. To check it, it suffices to let  $q_h = \chi_\Omega$  in (2.26b), observe that the right-hand side is equal to zero since  $g^n$  has zero average on  $\Omega$  (cf. (2.1e)), and use the definition (2.19) of  $c_h$  together with (2.23) to prove that the left-hand side also vanishes. This remark is crucial to ensure the local conservation properties of the method detailed in Section 2.7.

## 2.3 Stability analysis

In this section we study the stability of problem (2.26) and prove its well-posedness. We recall the following discrete Gronwall's inequality, which is a minor variation of [121, Lemma 5.1].

**Lemma 2.6** (Discrete Gronwall's inequality). *Let an integer  $N$  and reals  $\delta, G > 0$ , and  $K \geq 0$  be given, and let  $(a^n)_{0 \leq n \leq N}$ ,  $(b^n)_{0 \leq n \leq N}$ , and  $(\gamma^n)_{0 \leq n \leq N}$  denote three sequences of nonnegative real numbers such that, for all  $0 \leq n \leq N$*

$$a^n + \delta \sum_{m=0}^n b^m + K \leq \delta \sum_{m=0}^n \gamma^m a^m + G.$$

Then, if  $\gamma^m \delta < 1$  for all  $0 \leq m \leq N$ , letting  $\varsigma^m := (1 - \gamma^m \delta)^{-1}$ , it holds, for all  $0 \leq n \leq N$ ,

$$a^n + \delta \sum_{m=0}^n b^m + K \leq \exp\left(\delta \sum_{m=0}^n \varsigma^m \gamma^m\right) \times G. \quad (2.29)$$

**Lemma 2.7** (A priori bounds). *Assume  $f \in C^1(L^2(\Omega)^d)$  and  $g \in C^0(L^2(\Omega))$ , and let  $(\underline{\mathbf{u}}_h^0, p_h^0) = (\widehat{\underline{\mathbf{u}}}_h^0, \widehat{p}_h^0)$  with  $(\widehat{\underline{\mathbf{u}}}_h^0, \widehat{p}_h^0)$  defined as in Section 2.2.5. For all  $1 \leq n \leq N$ , denote by  $(\underline{\mathbf{u}}_h^n, p_h^n)$  the solution to (2.26). Then, for  $\tau$  small enough, it holds that*

$$\begin{aligned} & \|\underline{\mathbf{u}}_h^N\|_{a,h}^2 + \|c_0^{1/2} p_h^N\|^2 + \frac{1}{2\mu + d\lambda} \|p_h^N - \bar{p}_h^N\|^2 + \sum_{n=1}^N \tau \|p_h^n\|_{c,h}^2 \lesssim \\ & \left(\frac{1}{2\mu} + c_0\right) \|p^0\|^2 + \frac{d_\Omega^2}{2\mu} \|f\|_{C^1(L^2(\Omega)^d)}^2 + (2\mu + d\lambda) t_F^2 \|g\|_{C^0(L^2(\Omega))}^2 + \frac{t_F^2}{c_0} \|\bar{g}\|_{C^0(L^2(\Omega))}^2, \end{aligned} \quad (2.30)$$

with the convention that  $c_0^{-1} \|\bar{g}\|_{C^0(L^2(\Omega))}^2 = 0$  if  $c_0 = 0$  and, for  $0 \leq n \leq N$ ,  $\bar{p}_h^n := (p_h^n, 1)$ .

*Remark 2.8* (Well-posedness). Owing to linearity, the well-posedness of (2.26) is an immediate consequence of Lemma 2.7.

*Remark 2.9* (A priori bound for  $c_0 = 0$ ). When  $c_0 = 0$ , the choice (2.17) of the discrete space for the pressure ensures that  $\bar{p}_h^n = 0$  for all  $0 \leq n \leq N$ . Thus, the third term in the left-hand side of (2.30) yields an estimate on  $\|p_h^N\|^2$ , and the a priori bound reads

$$\begin{aligned} & \|\underline{\mathbf{u}}_h^N\|_{a,h}^2 + \frac{1}{2\mu + d\lambda} \|p_h^N\|^2 + \sum_{n=1}^N \tau \|p_h^n\|_{c,h}^2 \lesssim \\ & (2\mu)^{-1} \left(d_\Omega^2 \|f\|_{C^1(L^2(\Omega)^d)}^2 + \|p^0\|^2\right) + (2\mu + d\lambda) t_F^2 \|g\|_{C^0(L^2(\Omega))}^2. \end{aligned}$$

The convention  $c_0^{-1} \|\bar{g}\|_{C^0(L^2(\Omega))}^2 = 0$  if  $c_0 = 0$  is justified since the term  $\mathfrak{T}_2$  in point (4) of the following proof vanishes in this case thanks to the compatibility condition (2.1e).

*Proof of Lemma 2.7.* Throughout the proof,  $C_i$  with  $i \in \mathbb{N}^*$  will denote a generic positive constant independent of  $h$ ,  $\tau$ , and of the physical parameters  $c_0$ ,  $\lambda$ ,  $\mu$ , and  $\kappa$ .

(1) *Estimate of  $\|p_h^n - \bar{p}_h^n\|$ .* Using the inf-sup condition (2.24) followed by (2.23) to infer that  $b_h(\underline{\mathbf{v}}_h, \bar{p}_h^n) = 0$ , the mechanical equilibrium equation (2.26a), and the second inequality in (2.14), for all  $1 \leq n \leq N$  we get

$$\begin{aligned} \|p_h^n - \bar{p}_h^n\| &\leq \beta \sup_{\underline{\mathbf{v}}_h \in \underline{U}_{h,D}^k \setminus \{\mathbf{0}\}} \frac{b_h(\underline{\mathbf{v}}_h, p_h^n - \bar{p}_h^n)}{\|\underline{\mathbf{v}}_h\|_{\epsilon,h}} = \beta \sup_{\underline{\mathbf{v}}_h \in \underline{U}_{h,D}^k \setminus \{\mathbf{0}\}} \frac{b_h(\underline{\mathbf{v}}_h, p_h^n)}{\|\underline{\mathbf{v}}_h\|_{\epsilon,h}} \\ &= \beta \sup_{\underline{\mathbf{v}}_h \in \underline{U}_{h,D}^k \setminus \{\mathbf{0}\}} \frac{l_h^n(\underline{\mathbf{v}}_h) - a_h(\underline{\mathbf{u}}_h^n, \underline{\mathbf{v}}_h)}{\|\underline{\mathbf{v}}_h\|_{\epsilon,h}} \leq C_1^{1/2} \left( d_\Omega \|f^n\| + (2\mu + d\lambda)^{1/2} \|\underline{\mathbf{u}}_h^n\|_{a,h} \right), \end{aligned}$$

where we have set, for the sake of brevity,  $C_1^{1/2} := \beta \max(C_K, \eta)$ . This implies, in particular,

$$\|p_h^n - \bar{p}_h^n\|^2 \leq 2C_1 \left( d_\Omega^2 \|f^n\|^2 + (2\mu + d\lambda) \|\underline{\mathbf{u}}_h^n\|_{a,h}^2 \right) \quad (2.31)$$

(2) *Energy balance.* Adding (2.26a) with  $\underline{\mathbf{v}}_h = \tau \delta_t \underline{\mathbf{u}}_h^n$  to (2.26b) with  $q_h = \tau p_h^n$ , and summing the resulting equation over  $1 \leq n \leq N$ , it is inferred

$$\sum_{n=1}^N \tau a_h(\underline{\mathbf{u}}_h^n, \delta_t \underline{\mathbf{u}}_h^n) + \sum_{n=1}^N \tau (c_0 \delta_t p_h^n, p_h^n) + \sum_{n=1}^N \tau \|p_h^n\|_{c,h}^2 = \sum_{n=1}^N \tau l_h^n(\delta_t \underline{\mathbf{u}}_h^n) + \sum_{n=1}^N \tau (g^n, p_h^n). \quad (2.32)$$

We denote by  $\mathcal{L}$  and  $\mathcal{R}$  the left- and right-hand side of (2.32) and proceed to find suitable lower and upper bounds, respectively.

(3) *Lower bound for  $\mathcal{L}$ .* Using twice the formula

$$2x(x - y) = x^2 + (x - y)^2 - y^2, \quad (2.33)$$

and telescoping out the appropriate summands, the first two terms in the left-hand side of (2.32) can be rewritten as, respectively,

$$\begin{aligned} \sum_{n=1}^N \tau a_h(\underline{\mathbf{u}}_h^n, \delta_t \underline{\mathbf{u}}_h^n) &= \frac{1}{2} \|\underline{\mathbf{u}}_h^N\|_{a,h}^2 + \frac{1}{2} \sum_{n=1}^N \tau^2 \|\delta_t \underline{\mathbf{u}}_h^n\|_{a,h}^2 - \frac{1}{2} \|\underline{\mathbf{u}}_h^0\|_{a,h}^2, \\ \sum_{n=1}^N \tau (c_0 \delta_t p_h^n, p_h^n) &= \frac{1}{2} \|c_0^{1/2} p_h^N\|^2 + \frac{1}{2} \sum_{n=1}^N \tau^2 \|c_0^{1/2} \delta_t p_h^n\|^2 - \frac{1}{2} \|c_0^{1/2} p_h^0\|^2. \end{aligned}$$

Using the above relation together with (2.31) and  $\|f^N\| \leq \|f\|_{C^1(L^2(\Omega)^d)}$ , it is inferred that

$$\begin{aligned} &\frac{1}{4} \|\underline{\mathbf{u}}_h^N\|_{a,h}^2 - \frac{1}{2} \|\underline{\mathbf{u}}_h^0\|_{a,h}^2 + \frac{1}{2} \|c_0^{1/2} p_h^N\|^2 - \frac{1}{2} \|c_0^{1/2} p_h^0\|^2 \\ &+ \frac{1}{8C_1(2\mu + d\lambda)} \|p_h^N - \bar{p}_h^N\|^2 + \sum_{n=1}^N \tau \|p_h^n\|_{c,h}^2 \leq \mathcal{L} + \frac{d_\Omega^2}{4(2\mu + d\lambda)} \|f\|_{C^1(L^2(\Omega)^d)}^2. \end{aligned} \quad (2.34)$$

(4) *Upper bound for  $\mathcal{R}$ .* For the first term in the right-hand side of (2.32), discrete integration by parts in time yields

$$\sum_{n=1}^N \tau l_h^n (\delta_t \underline{\mathbf{u}}_h^n) = (\mathbf{f}^N, \underline{\mathbf{u}}_h^N) - (\mathbf{f}^0, \underline{\mathbf{u}}_h^0) - \sum_{n=1}^N \tau (\delta_t \mathbf{f}^n, \underline{\mathbf{u}}_h^{n-1}),$$

hence, using the Cauchy–Schwarz inequality, the discrete Korn’s inequality followed by (2.14) to estimate  $\|\underline{\mathbf{u}}_h^n\|^2 \leq \frac{C_2 d_\Omega^2}{\mu} \|\underline{\mathbf{u}}_h^n\|_{a,h}^2$  for all  $1 \leq n \leq N$  (with  $C_2 := C_K^2 \eta / 2$ ), and Young’s inequality, one has

$$\begin{aligned} \left| \sum_{n=1}^N \tau l_h^n (\delta_t \underline{\mathbf{u}}_h^n) \right| &\leq \frac{1}{8} \left( \|\underline{\mathbf{u}}_h^N\|_{a,h}^2 + \|\underline{\mathbf{u}}_h^0\|_{a,h}^2 + \frac{1}{2t_F} \sum_{n=1}^N \tau \|\underline{\mathbf{u}}_h^{n-1}\|_{a,h}^2 \right) \\ &\quad + \frac{C_2 d_\Omega^2}{\mu} \left( \|\mathbf{f}^N\|^2 + \|\mathbf{f}^0\|^2 + 2t_F \sum_{n=1}^N \tau \|\delta_t \mathbf{f}^n\|^2 \right) \\ &\leq \frac{1}{8} \left( \|\underline{\mathbf{u}}_h^N\|_{a,h}^2 + \|\underline{\mathbf{u}}_h^0\|_{a,h}^2 + \frac{1}{2t_F} \sum_{n=0}^N \tau \|\underline{\mathbf{u}}_h^n\|_{a,h}^2 \right) + \frac{C_2 C_3 d_\Omega^2}{\mu} \|\mathbf{f}\|_{C^1(L^2(\Omega)^d)}^2, \end{aligned} \quad (2.35)$$

where we have used the classical bound  $\|\mathbf{f}^N\|^2 + \|\mathbf{f}^0\|^2 + 2t_F \sum_{n=1}^N \tau \|\delta_t \mathbf{f}^n\|^2 \leq C_3 \|\mathbf{f}\|_{C^1(L^2(\Omega)^d)}^2$  to conclude. We proceed to estimate the second term in the right-hand side of (2.32) by splitting it into two contributions as follows (here,  $\bar{g}^n := (g^n, 1)$ ):

$$\sum_{n=1}^N \tau (g^n, p_h^n) = \sum_{n=1}^N \tau (g^n, p_h^n - \bar{p}_h^n) + \sum_{n=1}^N \tau (\bar{g}^n, p_h^n) := \mathfrak{I}_1 + \mathfrak{I}_2.$$

Using the Cauchy–Schwarz inequality, the bound  $\sum_{n=1}^N \tau \|g^n\|^2 \leq t_F \|g\|_{C^0(L^2(\Omega))}^2$  together with (2.31) and Young’s inequality, it is inferred that

$$\begin{aligned} |\mathfrak{I}_1| &\leq \left\{ \sum_{n=1}^N \tau \|g^n\|^2 \right\}^{1/2} \times \left\{ \sum_{n=1}^N \tau \|p_h^n - \bar{p}_h^n\|^2 \right\}^{1/2} \\ &\leq t_F \|g\|_{C^0(L^2(\Omega))} \times \left\{ \frac{2C_1}{t_F} \sum_{n=1}^N \tau \left( d_\Omega^2 \|\mathbf{f}^n\|^2 + (2\mu + d\lambda) \|\underline{\mathbf{u}}_h^n\|_{a,h}^2 \right) \right\}^{1/2} \\ &\leq 8C_1 t_F^2 (2\mu + d\lambda) \|g\|_{C^0(L^2(\Omega))}^2 + \frac{d_\Omega^2}{16(2\mu + d\lambda)} \|\mathbf{f}\|_{C^1(L^2(\Omega)^d)}^2 + \frac{1}{16t_F} \sum_{n=1}^N \tau \|\underline{\mathbf{u}}_h^n\|_{a,h}^2. \end{aligned} \quad (2.36)$$

Owing the compatibility condition (2.1e),  $\mathfrak{I}_2 = 0$  if  $c_0 = 0$ . If  $c_0 > 0$ , using the Cauchy–



Schwarz and Young's inequalities, we have

$$\begin{aligned} |\mathfrak{I}_2| &\leq \left\{ t_F \sum_{n=1}^N \tau c_0^{-1} \|\bar{g}^n\|^2 \right\}^{1/2} \times \left\{ t_F^{-1} \sum_{n=1}^N \tau \|c_0^{1/2} p_h^n\|^2 \right\}^{1/2} \\ &\leq \frac{t_F^2}{2c_0} \|\bar{g}\|_{C^0(L^2(\Omega))}^2 + \frac{1}{2t_F} \sum_{n=1}^N \tau \|c_0^{1/2} p_h^n\|^2. \end{aligned} \quad (2.37)$$

Using (2.35), (2.36), and (2.37), we infer

$$\begin{aligned} \mathcal{R} &\leq \frac{1}{8} \left( \|\underline{\mathbf{u}}_h^N\|_{a,h}^2 + t_F^{-1} \sum_{n=0}^N \tau \|\underline{\mathbf{u}}_h^n\|_{a,h}^2 + \|\underline{\mathbf{u}}_h^0\|_{a,h}^2 \right) + \frac{1}{2t_F} \sum_{n=1}^N \tau \|c_0^{1/2} p_h^n\|^2 + \frac{t_F^2}{2c_0} \|\bar{g}\|_{C^0(L^2(\Omega))}^2 \\ &\quad + 8C_1 t_F^2 (2\mu + d\lambda) \|g\|_{C^0(L^2(\Omega))}^2 + \left( \frac{1}{16(2\mu + d\lambda)} + \frac{C_2 C_3}{\mu} \right) d_\Omega^2 \|\mathbf{f}\|_{C^1(L^2(\Omega)^d)}^2. \end{aligned} \quad (2.38)$$

(5) *Conclusion.* Using (2.34), the fact that  $\mathcal{L} = \mathcal{R}$  owing to (2.32), and (2.38), it is inferred that

$$\begin{aligned} \|\underline{\mathbf{u}}_h^N\|_{a,h}^2 + 4\|c_0^{1/2} p_h^N\|^2 + \frac{1}{(2\mu + d\lambda)} \|p_h^N - \bar{p}_h^N\|^2 + 8 \sum_{n=1}^N \tau \|p_h^n\|_{c,h}^2 &\leq \\ \frac{C_4}{t_F} \sum_{n=0}^N \tau \|\underline{\mathbf{u}}_h^n\|_{a,h}^2 + \frac{C_4}{t_F} \sum_{n=1}^N \tau 4\|c_0^{1/2} p_h^n\|^2 + G, \end{aligned}$$

where  $C_4 := \max(1, C_1)$  while, observing that  $\|c_0^{1/2} p_h^0\| \leq \|c_0^{1/2} p^0\|$  since  $\pi_h^k$  is a bounded operator, and that it follows from (2.39) below that  $\|\underline{\mathbf{u}}_h^0\|_{a,h}^2 \leq C_5 (2\mu)^{-1} (d_\Omega^2 \|\mathbf{f}^0\|^2 + \|p^0\|^2)$ ,

$$\begin{aligned} C_4^{-1} G &:= \frac{5C_5}{2\mu} (d_\Omega^2 \|\mathbf{f}^0\|^2 + \|p^0\|^2) + 4\|c_0^{1/2} p^0\|^2 + \frac{4t_F^2}{c_0} \|\bar{g}\|_{C^0(L^2(\Omega))}^2 \\ &\quad + 64C_1 t_F^2 (2\mu + d\lambda) \|g\|_{C^0(L^2(\Omega))}^2 + \left( \frac{5}{2(2\mu + d\lambda)} + \frac{8C_2 C_3}{\mu} \right) d_\Omega^2 \|\mathbf{f}\|_{C^1(L^2(\Omega)^d)}^2. \end{aligned}$$

Using Gronwall's Lemma 2.6 with  $a^0 := \|\underline{\mathbf{u}}_h^0\|_{a,h}^2$  and  $a^n := \|\underline{\mathbf{u}}_h^n\|_{a,h}^2 + 4\|c_0^{1/2} p_h^n\|^2$  for  $1 \leq n \leq N$ ,  $\delta := \tau$ ,  $b^0 := 0$  and  $b^n := \|p_h^n\|_{c,h}^2$  for  $1 \leq n \leq N$ ,  $K = \frac{1}{(2\mu + d\lambda)} \|p_h^N - \bar{p}_h^N\|^2$ , and  $\gamma^n = \frac{C_4}{t_F}$ , the desired result follows.  $\square$

**Proposition 2.10** (Stability and approximation properties for  $\widehat{\underline{\mathbf{u}}}_h^0$ ). *The initial displacement (2.28) satisfies the following stability condition:*

$$\|\widehat{\underline{\mathbf{u}}}_h^0\|_{a,h} \lesssim (2\mu)^{-1/2} (d_\Omega \|\mathbf{f}^0\| + \|p^0\|). \quad (2.39)$$

Additionally, recalling the global reduction map  $\underline{I}_h^k$  defined by (2.4), and assuming the additional regularity  $p_0 \in H^{k+1}(P_\Omega)$ ,  $\mathbf{u}^0 \in H^{k+2}(P_\Omega)^d$ , and  $\nabla \cdot \mathbf{u}^0 \in H^{k+1}(P_\Omega)$ , it holds

$$(2\mu)^{1/2} \|\widehat{\underline{\mathbf{u}}}_h^0 - \underline{I}_h^k \mathbf{u}^0\|_{a,h} \lesssim h^{k+1} \left( 2\mu \|\mathbf{u}^0\|_{H^{k+2}(P_\Omega)^d} + \lambda \|\nabla \cdot \mathbf{u}^0\|_{H^{k+1}(P_\Omega)} + \rho_\kappa^{1/2} \|p^0\|_{H^{k+1}(P_\Omega)} \right), \quad (2.40)$$

with global heterogeneity ratio  $\rho_\kappa := \bar{\kappa}/\underline{\kappa}$ .

*Proof.* (1) *Proof of (2.39).* Using the first inequality in (2.14) followed by the definition (2.28) of  $\widehat{\underline{\mathbf{u}}}_h^0$ , we have

$$\begin{aligned} \|\widehat{\underline{\mathbf{u}}}_h^0\|_{a,h} &\lesssim \sup_{\underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_{h,D}^k \setminus \{\mathbf{0}\}} \frac{a_h(\widehat{\underline{\mathbf{u}}}_h^0, \underline{\mathbf{v}}_h)}{(2\mu)^{1/2} \|\underline{\mathbf{v}}_h\|_{\epsilon,h}} \\ &= (2\mu)^{-1/2} \sup_{\underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_{h,D}^k \setminus \{\mathbf{0}\}} \frac{l_h^0(\underline{\mathbf{v}}_h) - b_h(\underline{\mathbf{v}}_h, \pi_h^k p^0)}{\|\underline{\mathbf{v}}_h\|_{\epsilon,h}} \lesssim (2\mu)^{-1/2} (d_\Omega \|f^0\| + \|p^0\|), \end{aligned}$$

where to conclude we have used the Cauchy–Schwarz and discrete Korn’s (2.16) inequalities for the first term in the numerator and the continuity (2.22) of  $b_h$  together with the  $L^2(\Omega)$ -stability of  $\pi_h^k$  for the second. (2) *Proof of (2.40).* The proof is analogous to that of point (3) in Lemma 2.11 except that we use the approximation properties (2.2) of  $\pi_h^k$  instead of (2.45). For this reason, elliptic regularity is not needed.  $\square$

## 2.4 Error analysis

In this section we carry out the error analysis of the method.

### 2.4.1 Projection

We consider the error with respect to the sequence of projections  $(\widehat{\underline{\mathbf{u}}}_h^n, \widehat{p}_h^n)_{1 \leq n \leq N}$ , of the exact solution defined as follows: For  $1 \leq n \leq N$ ,  $\widehat{p}_h^n \in P_h^k$  solves

$$c_h(\widehat{p}_h^n, q_h) = c_h(p^n, q_h) \quad \forall q_h \in \mathbb{P}_d^k(\mathcal{T}_h), \quad (2.41a)$$

with the closure condition  $\int_\Omega \widehat{p}_h^n = \int_\Omega p^n$ . Once  $\widehat{p}_h^n$  has been computed,  $\widehat{\underline{\mathbf{u}}}_h^n \in \underline{\mathbf{U}}_{h,D}^k$  solves

$$a_h(\widehat{\underline{\mathbf{u}}}_h^n, \underline{\mathbf{v}}_h) = l_h^n(\underline{\mathbf{v}}_h) - b_h(\underline{\mathbf{v}}_h, \widehat{p}_h^n) \quad \forall \underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_{h,D}^k. \quad (2.41b)$$

The well-posedness of problems (2.41a) and (2.41b) follow, respectively, from the coercivity of  $c_h$  on  $\mathbb{P}_{d,0}^k(\mathcal{T}_h)$  and of  $a_h$  on  $\underline{\mathbf{U}}_{h,D}^k$ . The projection  $(\widehat{\underline{\mathbf{u}}}_h^n, \widehat{p}_h^n)$  is chosen so that a convergence rate of  $(k+1)$  in space analogous to the one derived in [74] can be proved for the  $\|\cdot\|_{a,h}$ -norm of the displacement at final time  $t_F$ . To this purpose, we also need in what follows the following elliptic regularity, which holds, e.g., when  $\Omega$  is convex: There is a real number  $C_{\text{ell}} > 0$  only depending on  $\Omega$  such that, for all  $\psi \in L_0^2(\Omega)$ , with  $L_0^2(\Omega) := \{q \in L^2(\Omega) : (q, 1) = 0\}$ , the unique function  $\zeta \in H^1(\Omega) \cap L_0^2(\Omega)$  solution of the homogeneous Neumann problem

$$-\nabla \cdot (\kappa \nabla \zeta) = \psi \quad \text{in } \Omega, \quad \kappa \nabla \zeta \cdot \mathbf{n} = 0 \quad \text{on } \partial\Omega, \quad (2.42)$$

is such that

$$\|\zeta\|_{H^2(P_\Omega)} \leq C_{\text{ell}} \kappa^{-1/2} \|\psi\|. \quad (2.43)$$

For further insight on the role of the choice (2.41) and of the elliptic regularity assumption we refer to Remark 2.14.

**Lemma 2.11** (Approximation properties for  $(\widehat{\underline{\mathbf{u}}}_h^n, \widehat{p}_h^n)$ ). *Let a time step  $1 \leq n \leq N$  be fixed. Assuming the regularity  $p^n \in H^{k+1}(P_\Omega)$ , it holds*

$$\|\widehat{p}_h^n - p^n\|_{c,h} \lesssim h^k \bar{\kappa}^{1/2} \|p^n\|_{H^{k+1}(P_\Omega)}. \quad (2.44)$$

Moreover, recalling the global reduction map  $\underline{I}_h^k$  defined by (2.4), further assuming the regularity  $\mathbf{u}^n \in H^{k+2}(P_\Omega)^d$ ,  $\nabla \cdot \mathbf{u}^n \in H^{k+1}(P_\Omega)$ , and provided that the elliptic regularity (2.43) holds, one has

$$\|\widehat{p}_h^n - p^n\| \lesssim h^{k+1} \rho_\kappa^{1/2} \|p^n\|_{H^{k+1}(P_\Omega)}, \quad (2.45)$$

$$(2\mu)^{1/2} \|\widehat{\underline{\mathbf{u}}}_h^n - \underline{I}_h^k \mathbf{u}^n\|_{a,h} \lesssim h^{k+1} \left( 2\mu \|\mathbf{u}^n\|_{H^{k+2}(P_\Omega)^d} + \lambda \|\nabla \cdot \mathbf{u}^n\|_{H^{k+1}(P_\Omega)} + \rho_\kappa^{1/2} \|p^n\|_{H^{k+1}(P_\Omega)} \right). \quad (2.46)$$

*Proof.* (1) *Proof of (2.44).* By definition, we have that  $\|\widehat{p}_h^n - p^n\|_{c,h} = \inf_{q_h \in \mathbb{P}_d^k(\mathcal{T}_h)} \|q_h - p^n\|_{c,h}$ . To prove (2.44), it suffices to take  $q_h = \pi_h^k p^n$  in the right-hand side of the previous expression and use the approximation properties (2.2) of  $\pi_h^k$ .

(2) *Proof of (2.45).* Let  $\zeta \in H^1(\Omega)$  solve (2.42) with  $\psi = p^n - \widehat{p}_h^n$ . From the consistency property (2.20), it follows that

$$\|\widehat{p}_h^n - p^n\|^2 = -(\nabla \cdot (\kappa \nabla \zeta), \widehat{p}_h^n - p^n) = c_h(\zeta, \widehat{p}_h^n - p^n) = c_h(\zeta - \pi_h^1 \zeta, \widehat{p}_h^n - p^n).$$

Then, using the Cauchy–Schwarz inequality, the estimate (2.44) together with the approximation properties (2.2) of  $\pi_h^1$ , and elliptic regularity, it is inferred that

$$\begin{aligned} \|\widehat{p}_h^n - p^n\|^2 &= c_h(\zeta - \pi_h^1 \zeta, \widehat{p}_h^n - p^n) \leq \|\zeta - \pi_h^1 \zeta\|_{c,h} \|\widehat{p}_h^n - p^n\|_{c,h} \\ &\lesssim h^{k+1} \bar{\kappa}^{1/2} \|\zeta\|_{H^2(P_\Omega)} \|p^n\|_{H^{k+1}(P_\Omega)} \lesssim h^{k+1} \rho_\kappa^{1/2} \|\widehat{p}_h^n - p^n\| \|p^n\|_{H^{k+1}(P_\Omega)}, \end{aligned}$$

and (2.45) follows.

(3) *Proof of (2.46).* We start by observing that

$$\|\widehat{\underline{\mathbf{u}}}_h^n - \underline{I}_h^k \mathbf{u}^n\|_{a,h} = \sup_{\underline{\mathbf{v}}_h \in \underline{U}_h^k \setminus \{\mathbf{0}\}} \frac{a_h(\widehat{\underline{\mathbf{u}}}_h^n - \underline{I}_h^k \mathbf{u}^n, \underline{\mathbf{v}}_h)}{\|\underline{\mathbf{v}}_h\|_{a,h}} \lesssim \sup_{\underline{\mathbf{v}}_h \in \underline{U}_h^k \setminus \{\mathbf{0}\}} \frac{a_h(\widehat{\underline{\mathbf{u}}}_h^n - \underline{I}_h^k \mathbf{u}^n, \underline{\mathbf{v}}_h)}{(2\mu)^{1/2} \|\underline{\mathbf{v}}_h\|_{\epsilon,h}},$$

where we have used the first inequality in (2.14). Recalling the definition (2.27) of the linear form  $l_h^n$ , the fact that  $\mathbf{f}^n = -\nabla \cdot \boldsymbol{\sigma}(\mathbf{u}) + \nabla p$ , and using (2.41a), it is inferred that

$$\begin{aligned} a_h(\widehat{\underline{\mathbf{u}}}_h^n - \underline{I}_h^k \mathbf{u}^n, \underline{\mathbf{v}}_h) &= l_h^n(\underline{\mathbf{v}}_h^n) - a_h(\underline{I}_h^k \mathbf{u}^n, \underline{\mathbf{v}}_h) - b_h(\underline{\mathbf{v}}_h, \widehat{p}_h^n) \\ &= \{ -a_h(\underline{I}_h^k \mathbf{u}^n, \underline{\mathbf{v}}_h) - (\nabla \cdot \boldsymbol{\sigma}(\mathbf{u}^n), \underline{\mathbf{v}}_h) \} + \{ (\nabla p^n, \underline{\mathbf{v}}_h) - b_h(\underline{\mathbf{v}}_h, \widehat{p}_h^n) \}. \end{aligned} \quad (2.47)$$

Denote by  $\mathfrak{T}_1$  and  $\mathfrak{T}_2$  the terms in braces. Using (2.15), it is readily inferred that

$$|\mathfrak{T}_1| \lesssim h^{k+1} \left( 2\mu \|\mathbf{u}^n\|_{H^{k+2}(P_\Omega)^d} + \lambda \|\nabla \cdot \mathbf{u}^n\|_{H^{k+1}(P_\Omega)} \right) \|\underline{\mathbf{v}}_h\|_{\epsilon,h}. \quad (2.48)$$

For the second term, performing an element-wise integration by parts on  $(\nabla p, \mathbf{v}_h)$  and recalling the definition (2.21) of  $b_h$  and (2.7a) of  $D_T^k$  with  $q = \widehat{p}_h^n$ , it is inferred that

$$\begin{aligned} |\mathfrak{I}_2| &= \left| \sum_{T \in \mathcal{T}_h} \left\{ (\widehat{p}_h^n - p^n, \nabla \cdot \mathbf{v}_T)_T + \sum_{F \in \mathcal{F}_T} (\widehat{p}_h^n - p^n, (\mathbf{v}_F - \mathbf{v}_T) \mathbf{n}_{TF})_F \right\} \right| \\ &\lesssim h^{k+1} \rho_\kappa^{1/2} \|p^n\|_{H^{k+1}(P_\Omega)} \|\mathbf{v}_h\|_{\epsilon, h}, \end{aligned} \quad (2.49)$$

where the conclusion follows from the Cauchy–Schwarz inequality together with (2.45). Plugging (2.48)–(2.49) into (2.47) we obtain (2.46).  $\square$

## 2.4.2 Error equations

We define the discrete error components as follows: For all  $1 \leq n \leq N$ ,

$$\underline{\mathbf{e}}_h^n := \underline{\mathbf{u}}_h^n - \widehat{\underline{\mathbf{u}}}_h^n, \quad \rho_h^n := p_h^n - \widehat{p}_h^n. \quad (2.50)$$

Owing to the choice of the initial condition detailed in Section 2.2.5, the initial error  $(\underline{\mathbf{e}}_h^0, \rho_h^0) := (\underline{\mathbf{u}}_h^0 - \widehat{\underline{\mathbf{u}}}_h^0, p_h^0 - \widehat{p}_h^0)$  is the null element in the product space  $\underline{\mathbf{U}}_{h,D}^k \times P_h^k$ . On the other hand, for all  $1 \leq n \leq N$ ,  $(\underline{\mathbf{e}}_h^n, \rho_h^n)$  solves

$$a_h(\underline{\mathbf{e}}_h^n, \mathbf{v}_h) + b_h(\mathbf{v}_h, \rho_h^n) = 0 \quad \forall \mathbf{v}_h \in \underline{\mathbf{U}}_h^k, \quad (2.51a)$$

$$(c_0 \delta_t \rho_h^n, q_h) - b_h(\delta_t \underline{\mathbf{e}}_h^n, q_h) + c_h(\rho_h^n, q_h) = \mathcal{E}_h^n(q_h), \quad \forall q_h \in P_h^k, \quad (2.51b)$$

with consistency error

$$\mathcal{E}_h^n(q_h) := (g^n, q_h) - (c_0 \delta_t \widehat{p}_h^n, q_h) - c_h(\widehat{p}_h^n, q_h) + b_h(\delta_t \widehat{\underline{\mathbf{u}}}_h^n, q_h).$$

## 2.4.3 Convergence

**Theorem 2.12** (Estimate for the discrete errors). *Let  $(\mathbf{u}, p)$  denote the unique solution to (2.1), for which we assume the regularity*

$$\mathbf{u} \in C^2(H^1(P_\Omega)^d) \cap C^1(H^{k+2}(P_\Omega)^d), \quad p \in C^1(H^{k+1}(P_\Omega)).$$

If  $c_0 > 0$ , we further assume  $p \in C^2(L^2(\Omega))$ . Define, for the sake of brevity, the bounded quantities

$$\begin{aligned} \mathcal{N}_1 &:= (2\mu + d\lambda)^{1/2} \|\mathbf{u}\|_{C^2(H^1(P_\Omega)^d)} + \|c_0^{1/2} p\|_{C^2(L^2(\Omega)^d)}, \\ \mathcal{N}_2 &:= \frac{(2\mu + d\lambda)^{1/2}}{2\mu} \left( 2\mu \|\mathbf{u}\|_{C^1(H^{k+2}(P_\Omega)^d)} + \lambda \|\nabla \cdot \mathbf{u}\|_{C^1(H^{k+1}(P_\Omega))} + \rho_\kappa^{1/2} \|p\|_{C^1(H^{k+1}(P_\Omega))} \right) \\ &\quad + \|c_0^{1/2} p\|_{C^0(H^{k+1}(P_\Omega))}. \end{aligned}$$

Then, assuming the elliptic regularity (2.43), it holds, letting  $\overline{\rho}_h^n := (\rho_h^n, 1)$ ,

$$\|\underline{\mathbf{e}}_h^N\|_{a,h}^2 + \|c_0^{1/2} \rho_h^N\|^2 + \frac{1}{2\mu + d\lambda} \|\rho_h^N - \overline{\rho}_h^N\|^2 + \sum_{n=1}^N \tau \|\rho_h^n\|_{c,h}^2 \lesssim (\tau \mathcal{N}_1 + h^{k+1} \mathcal{N}_2)^2. \quad (2.52)$$

*Remark 2.13* (Pressure estimate for  $c_0 = 0$ ). In the incompressible case  $c_0 = 0$ , the third term in the left-hand side of (2.52) delivers an estimate on the  $L^2$ -norm of the pressure since  $\bar{\rho}_h^N = 0$  (cf. (2.1f)).

*Proof of Theorem 2.12.* Throughout the proof,  $C_i$  with  $i \in \mathbb{N}^*$  will denote a generic positive constant independent of  $h$ ,  $\tau$ , and of the physical parameters  $c_0$ ,  $\lambda$ ,  $\mu$ , and  $\kappa$ .

(1) *Basic error estimate.* Using the inf-sup condition (2.24), equation (2.23) followed by (2.51a), and the second inequality in (2.14), it is readily seen that

$$\begin{aligned} \|\rho_h^n - \bar{\rho}_h^n\| &\leq \beta \sup_{\mathbf{v}_h \in \underline{U}_{h,D}^k \setminus \{\mathbf{0}\}} \frac{b_h(\mathbf{v}_h, \rho_h^n - \bar{\rho}_h^n)}{\|\mathbf{v}_h\|_{\epsilon,h}} \\ &= \beta \sup_{\mathbf{v}_h \in \underline{U}_{h,D}^k \setminus \{\mathbf{0}\}} \frac{-a_h(\mathbf{e}_h^n, \mathbf{v}_h)}{\|\mathbf{v}_h\|_{\epsilon,h}} \leq C_1^{1/2} (2\mu + d\lambda)^{1/2} \|\mathbf{e}_h^n\|_{a,h}, \end{aligned}$$

with  $C_1^{1/2} = \beta\eta^{1/2}$ . Adding (2.51a) with  $\mathbf{v}_h = \tau\delta_t \mathbf{e}_h$  to (2.51b) with  $q_h = \tau\rho_h^n$  and summing the resulting equation over  $1 \leq n \leq N$ , it is inferred that

$$\sum_{n=1}^N \tau a_h(\mathbf{e}_h^n, \delta_t \mathbf{e}_h^n) + \sum_{n=1}^N \tau (c_0 \delta_t \rho_h^n, \rho_h^n) + \sum_{n=1}^N \tau \|\rho_h^n\|_{c,h}^2 = \sum_{n=1}^N \tau \mathcal{E}_h^n(\rho_h^n). \quad (2.53)$$

Proceeding as in point (3) of the proof of Lemma 2.7, and recalling that  $(\mathbf{e}_h^0, \rho_h^0) = (\mathbf{0}, 0)$ , we arrive at the following error estimate:

$$\frac{1}{4} \|\mathbf{e}_h^N\|_{a,h}^2 + \frac{1}{4C_1(2\mu + d\lambda)} \|\rho_h^N - \bar{\rho}_h^N\|^2 + \frac{1}{2} \|c_0^{1/2} \rho_h^N\|^2 + \sum_{n=1}^N \tau \|\rho_h^n\|_{c,h}^2 \leq \sum_{n=1}^N \tau \mathcal{E}_h^n(\rho_h^n). \quad (2.54)$$

(2) *Bound of the consistency error.* Using  $g^n = c_0 d_t p^n + \nabla \cdot (d_t \mathbf{u}^n - \kappa \nabla p^n)$ , the consistency property (2.20), and observing that, using the definition (2.19) of  $c_h$ , integration by parts together with the homogeneous displacement boundary condition (2.1c), and (2.23),

$$c_h(p^n - \widehat{p}_h^n, \bar{\rho}_h^n) + (\nabla \cdot (d_t \mathbf{u}^n), \bar{\rho}_h^n) + b_h(\delta_t \widehat{\mathbf{u}}_h^n, \bar{\rho}_h^n) = 0,$$

we can decompose the right-hand side of (2.54) as follows:

$$\begin{aligned} \sum_{n=1}^N \tau \mathcal{E}_h^n(\rho_h^n) &= \sum_{n=1}^N \tau (c_0 (d_t p^n - \delta_t \widehat{p}_h^n), \rho_h^n) + \sum_{n=1}^N \tau c_h(p^n - \widehat{p}_h^n, \rho_h^n - \bar{\rho}_h^n) \\ &\quad + \sum_{n=1}^N \tau \left\{ (\nabla \cdot (d_t \mathbf{u}^n), \rho_h^n - \bar{\rho}_h^n) + b_h(\delta_t \widehat{\mathbf{u}}_h^n, \rho_h^n - \bar{\rho}_h^n) \right\} := \mathfrak{I}_1 + \mathfrak{I}_2 + \mathfrak{I}_3. \end{aligned} \quad (2.55)$$

For the first term, inserting  $\pm \delta_t p^n$  into the first factor and using the Cauchy-Schwarz inequality followed by the approximation properties of  $\widehat{p}_h^0$  (a consequence of (2.2)) and (2.45) of  $\widehat{p}_h^n$ , it

is inferred that

$$\begin{aligned} |\mathfrak{I}_1| &\lesssim \left\{ c_0 \sum_{n=1}^N \tau \left[ \|\mathbf{d}_t p^n - \delta_t p^n\|^2 + \|\delta_t(p^n - \widehat{p}_h^n)\|^2 \right] \right\}^{1/2} \times \left\{ \sum_{n=1}^N \tau \|c_0^{1/2} \rho_h^n\|^2 \right\}^{1/2} \\ &\leq C_2 \left( \tau \mathcal{N}_1 + h^{k+1} \mathcal{N}_2 \right) + \frac{1}{2} \sum_{n=1}^N \tau \|c_0^{1/2} \rho_h^n\|^2. \end{aligned} \quad (2.56)$$

For the second term, the choice (2.41a) of the pressure projection readily yields

$$\mathfrak{I}_2 = 0.$$

For the last term, inserting  $\pm \underline{I}_h^k \mathbf{u}^n$  into the first argument of  $b_h$ , and using the commuting property (2.8) of  $D_T^k$ , it is inferred that

$$\mathfrak{I}_3 = \sum_{n=1}^N \tau \left\{ \sum_{T \in \mathcal{T}_h} \left[ (\nabla \cdot (\mathbf{d}_t \mathbf{u}^n - \delta_t \mathbf{u}^n), \rho_h^n - \bar{\rho}_h^n)_T + (D_T^k \delta_t (\underline{I}_T^k \mathbf{u}^n - \widehat{\underline{\mathbf{u}}}_T^n), \rho_h^n - \bar{\rho}_h^n)_T \right] \right\}.$$

Using the Cauchy–Schwarz inequality, the bound  $\|D_T^k \delta_t (\underline{I}_T^k \mathbf{u}^n - \widehat{\underline{\mathbf{u}}}_T^n)\|_T \lesssim \|\delta_t (\underline{I}_T^k \mathbf{u}^n - \widehat{\underline{\mathbf{u}}}_T^n)\|_{\epsilon, T}$  valid for all  $T \in \mathcal{T}_h$ , and the approximation properties (2.40) and (2.46) of  $\widehat{\underline{\mathbf{u}}}_h^0$  and  $\widehat{\underline{\mathbf{u}}}_h^n$ , respectively, we obtain

$$\begin{aligned} |\mathfrak{I}_3| &\lesssim \left\{ \sum_{n=1}^N \tau \left[ \|\mathbf{d}_t \mathbf{u}^n - \delta_t \mathbf{u}^n\|_{H^1(\Omega)^d}^2 + \|\delta_t (\underline{I}_h^k \mathbf{u}^n - \widehat{\underline{\mathbf{u}}}_h^n)\|_{\epsilon, h}^2 \right] \right\}^{1/2} \times \left\{ \sum_{n=1}^N \tau \|\rho_h^n - \bar{\rho}_h^n\|^2 \right\}^{1/2} \\ &\leq C_3 C_1 \left( \tau \mathcal{N}_1 + h^{k+1} \mathcal{N}_2 \right)^2 + \frac{1}{4C_1(2\mu + d\lambda)} \sum_{n=1}^N \tau \|\rho_h^n - \bar{\rho}_h^n\|^2. \end{aligned} \quad (2.57)$$

Using (2.56)–(2.57) to bound the right-hand side of (2.55), it is inferred

$$\begin{aligned} \|\underline{\mathbf{e}}_h^N\|_{a, h}^2 + \frac{1}{C_1(2\mu + d\lambda)} \|\rho_h^N - \bar{\rho}_h^N\|^2 + 2\|c_0^{1/2} \rho_h^N\|^2 + 4 \sum_{n=1}^N \tau \|\rho_h^n\|_{c, h}^2 \\ \leq \frac{1}{C_1(2\mu + d\lambda)} \sum_{n=1}^N \tau \|\rho_h^n - \bar{\rho}_h^n\|^2 + 2 \sum_{n=1}^N \tau \|c_0^{1/2} \rho_h^n\|^2 + G, \end{aligned}$$

with  $G := 4(C_1 C_3 + C_2) \left( \tau \mathcal{N}_1 + h^{k+1} \mathcal{N}_2 \right)^2$ . The conclusion follows using the discrete Gronwall's inequality (2.29) with  $\delta = \tau$ ,  $K = \|\underline{\mathbf{e}}_h^N\|_{a, h}^2$ ,  $a^0 = 0$  and  $a^n = \frac{1}{C_1(2\mu + d\lambda)} \|\rho_h^n - \bar{\rho}_h^n\|^2 + 2\|c_0^{1/2} \rho_h^n\|^2$  for  $1 \leq n \leq N$ ,  $b^n = 4\|\rho_h^n\|_{c, h}^2$ , and  $\gamma^n = 1$ .  $\square$

*Remark 2.14* (Role of the choice (2.41) and of elliptic regularity). The choice (2.41) for the projection ensures that the term  $\mathfrak{I}_2$  in step (2) of the proof of Theorem 2.12 vanishes. This is a key point to obtain an order of convergence of  $(k + 1)$  in space. For a different choice, say  $\widehat{p}_h^n = \pi_h^k p^n$ , this term would be of order  $k$ , and therefore yield a suboptimal estimate for the terms in the left-hand side of (2.58) below (the estimate (2.59) would not change and remain optimal). This would also be the case if we removed the elliptic regularity assumption (2.43).

*Remark 2.15* (BDF2 time discretization). In some of the numerical test cases of Section 2.6, we have used a BDF2 time discretization, which corresponds to the backward differencing operator

$$\delta_t^{(2)} \varphi^{n+2} := \frac{3\varphi^{n+2} - 4\varphi^{n+1} + \varphi^n}{2\tau},$$

used in place of (2.3). As BDF2 requires two starting values, we perform a first march in time using the backward Euler scheme (another possibility would have been to resort to the second-order Crank–Nicolson scheme). For the BDF2 time discretization, stability estimates similar to those of Lemma 2.7 can be proved with this initialization, while the error can be shown to scale as  $\tau^2 + h^{k+1}$  (compare with (2.52)). The main difference with respect to the present analysis focused on the backward Euler scheme is that formula (2.33) is replaced in the proofs by

$$2x(3x - 4y + z) = x^2 - y^2 + (2x - y)^2 - (2y - z)^2 + (x - 2y + z)^2.$$

The modifications of the proofs are quite classical and are not detailed here for the sake of conciseness (for a pedagogic exposition, one can consult, e.g., [92, Chapter 6]).

**Corollary 2.16** (Convergence). *Under the assumptions of Theorem 2.12, it holds that*

$$\begin{aligned} (2\mu)^{1/2} \|\nabla_{s,h}(\mathbf{r}_h^{k+1} \underline{\mathbf{u}}_h^N - \mathbf{u}^N)\| + \|c_0^{1/2}(p_h^N - p^N)\| + \frac{1}{2\mu + d\lambda} \|(p_h^N - p^N) - (\bar{p}_h^N - \bar{p}^N)\| \\ \lesssim \tau \mathcal{N}_1 + h^{k+1} \mathcal{N}_2 + c_0^{1/2} h^{k+1} \|p^N\|_{H^{k+1}(P_\Omega)}, \end{aligned} \quad (2.58)$$

$$\left\{ \sum_{n=1}^N \tau \|p_h^n - p^n\|_{c,h}^2 \right\}^{1/2} \lesssim \tau \mathcal{N}_1 + h^{k+1} \mathcal{N}_2 + h^k \bar{\kappa}^{-1/2} t_F^{1/2} \|p\|_{C^0(H^{k+1}(P_\Omega))}. \quad (2.59)$$

*Proof.* Using the triangle inequality, recalling the definition (2.50) of  $\underline{\mathbf{e}}_h^N$  and  $\widehat{p}_h^N$  and (2.14) of  $\|\cdot\|_{a,h}$ -norm, it is inferred that

$$\begin{aligned} (2\mu)^{1/2} \|\nabla_{s,h}(\mathbf{r}_h^{k+1} \underline{\mathbf{u}}_h^N - \mathbf{u}^N)\| &\lesssim \|\underline{\mathbf{e}}_h^N\|_{a,h} + (2\mu)^{1/2} \|\nabla_{s,h}(\mathbf{r}_h^{k+1} \widehat{\underline{\mathbf{u}}}_h - \mathbf{r}_h^{k+1} \underline{\mathbf{I}}_h^k \mathbf{u}^N)\| \\ &\quad + (2\mu)^{1/2} \|\nabla_s(\mathbf{r}_h^{k+1} \underline{\mathbf{I}}_h^k \mathbf{u}^N - \mathbf{u}^N)\|, \\ \|p_h^N - p^N - (\bar{p}_h^N - \bar{p}^N)\| &\leq \|\rho_h^N - \bar{\rho}_h^N\| + \|\widehat{p}_h^N - p^N\|, \\ \|c_0^{1/2}(p_h^N - p^N)\| &\leq \|c_0^{1/2} \rho_h^N\| + \|c_0^{1/2}(\widehat{p}_h^N - p^N)\|. \end{aligned}$$

To conclude, use (2.52) to estimate the left-most terms in the right-hand sides of the above equations. Use (2.46) and (2.45), the approximation properties (2.6) of  $\mathbf{r}_h^{k+1} \underline{\mathbf{I}}_h^k$ , respectively, for the right-most terms. This proves (2.58). A similar decomposition of the error yields (2.59).  $\square$

## 2.5 Implementation

In this section we discuss practical aspects including, in particular, static condensation. The implementation is based on the hho platform<sup>1</sup>, which relies on the linear algebra facilities provided by the Eigen3 library [116].

The starting point consists in selecting a basis for each of the polynomial spaces appearing in the construction. Let  $\mathbf{s} = (s_1, \dots, s_d)$  be a  $d$ -dimensional multi-index with the usual notation  $|\mathbf{s}|_1 = \sum_{i=1}^d s_i$ , and let  $\mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d$ . Given  $k \geq 0$  and  $T \in \mathcal{T}_h$ , we denote by  $\mathcal{B}_T^k$  a basis for the polynomial space  $\mathbb{P}_d^k(T)$ . In the numerical experiments of Section 2.6, we have used the set of locally scaled monomials:

$$\mathcal{B}_T^k := \left\{ \left( \frac{\mathbf{x} - \mathbf{x}_T}{h_T} \right)^{\mathbf{s}}, |\mathbf{s}|_1 \leq k \right\}, \quad (2.60)$$

with  $\mathbf{x}_T$  denoting the barycenter of  $T$ . Similarly, for all  $F \in \mathcal{F}_h$ , we denote by  $\mathcal{B}_F^k$  a basis for the polynomial space  $\mathbb{P}_{d-1}^k(F)$  which, in the proposed implementation, is again a set of locally scaled monomials similar to (2.60).

*Remark 2.17* (Choice of the polynomial bases). The choice of the polynomial bases can have a sizeable impact on the conditioning of both the local problems defining the displacement reconstruction  $\mathbf{r}_T^{k+1}$  (cf. (2.5)) and the global problem. This is particularly the case when using high polynomial orders (typically,  $k \geq 7$ ). The scaled monomial basis (2.60) is appropriate when dealing with isotropic elements. In the presence of anisotropic elements, a better choice is to use for each element a local frame aligned with its principal axes of rotation together with normalization factors tailored for each direction. A further improvement, originally investigated in [16] in the context of dG methods, consists in performing a Gram–Schmidt orthonormalization with respect to a suitably selected inner product. In the numerical test cases of Section 2.6, which focus on isotropic meshes and moderate polynomial degrees ( $k \leq 3$ ), the basis (2.60) proved fully satisfactory.

Introducing the vector bases  $\underline{\mathcal{B}}_T^k := (\mathcal{B}_T^k)^d$ ,  $T \in \mathcal{T}_h$ , and  $\underline{\mathcal{B}}_F^k := (\mathcal{B}_F^k)^d$ ,  $F \in \mathcal{F}_h^i$ , a basis  $\underline{\mathcal{U}}_{h,0}^k$  for the space  $\underline{U}_{h,D}^k$  (cf. (2.12)) is given by

$$\underline{\mathcal{U}}_{h,0}^k := \underline{\mathcal{U}}_{\mathcal{T}}^k \times \underline{\mathcal{U}}_{\mathcal{F}}^k, \quad \underline{\mathcal{U}}_{\mathcal{T}}^k := \bigotimes_{T \in \mathcal{T}_h} \mathcal{B}_T^k, \quad \underline{\mathcal{U}}_{\mathcal{F}}^k := \bigotimes_{F \in \mathcal{F}_h^i} \mathcal{B}_F^k,$$

while a basis  $\mathcal{P}_h^k$  for the space  $P_h^k$  (cf. (2.17)) is obtained setting

$$\mathcal{P}_h^k := \bigotimes_{T \in \mathcal{T}_h} \mathcal{B}_T^k.$$

When  $c_0 = 0$ , the zero average constraint in  $P_h^k$  can be accounted for using as a Lagrange multiplier the characteristic function of  $\Omega$ . Notice also that boundary faces have been excluded from the Cartesian product in the definition of  $\underline{\mathcal{U}}_{\mathcal{F}}^k$  to strongly account for boundary

<sup>1</sup>DL15105 Université de Montpellier



conditions. Letting, for the sake of brevity,  $N_n^k := \binom{k+n}{k}$ ,  $n \in \mathbb{N}$ , a simple computation shows that

$$\dim(\mathbf{U}_{\mathcal{T}}^k) = d \operatorname{card}(\mathcal{T}_h) N_d^k, \quad \dim(\mathbf{U}_{\mathcal{F}}^k) = d \operatorname{card}(\mathcal{F}_h^i) N_{d-1}^k, \quad \dim(\mathcal{P}_h^k) = \operatorname{card}(\mathcal{T}_h) N_d^k.$$

The total DOF count thus yields

$$d \operatorname{card}(\mathcal{T}_h) N_d^k + d \operatorname{card}(\mathcal{F}_h^i) N_{d-1}^k + \operatorname{card}(\mathcal{T}_h) N_d^k. \quad (2.61)$$

In what follows, for a given time step  $0 \leq n \leq N$ , we denote by  $\mathbf{U}_{\mathcal{T}}^n$  and  $\mathbf{U}_{\mathcal{F}}^n$  the vectors collecting element-based and face-based displacement DOFs, respectively, and by  $\mathbf{P}^n$  the vector collecting pressure DOFs.

Denote now by  $\mathbf{A}$  and  $\mathbf{B}$ , respectively, the matrices that represent the bilinear forms  $a_h$  (cf. (2.11)) and  $b_h$  (cf. (2.21)) in the selected basis. Distinguishing element-based and face-based displacement DOFs, the matrices  $\mathbf{A}$  and  $\mathbf{B}$  display the following block structure:

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{\mathcal{T}\mathcal{T}} & \mathbf{A}_{\mathcal{T}\mathcal{F}} \\ \mathbf{A}_{\mathcal{T}\mathcal{F}}^{\top} & \mathbf{A}_{\mathcal{F}\mathcal{F}} \end{bmatrix}, \quad \mathbf{B} = \begin{bmatrix} \mathbf{B}_{\mathcal{T}} \\ \mathbf{B}_{\mathcal{F}} \end{bmatrix}.$$

For every mesh element  $T \in \mathcal{T}_h$ , the element-based displacement DOFs are only coupled with those face-based displacement DOFs that lie on the boundary of  $T$  and with the (element-based) pressure DOFs in  $T$ . This translates into the fact that the submatrix  $\mathbf{A}_{\mathcal{T}\mathcal{T}}$  is block-diagonal, i.e.,

$$\mathbf{A}_{\mathcal{T}\mathcal{T}} = \operatorname{diag}(\mathbf{A}_{TT})_{T \in \mathcal{T}_h},$$

with each elementary block  $\mathbf{A}_{TT}$  of size  $\dim(\mathcal{B}_T^k)^2$ . Additionally, it can be proved that the blocks  $\mathbf{A}_{TT}$ ,  $T \in \mathcal{T}_h$ , are invertible, so that the inverse of  $\mathbf{A}_{\mathcal{T}\mathcal{T}}$  can be efficiently computed setting

$$\mathbf{A}_{\mathcal{T}\mathcal{T}}^{-1} = \operatorname{diag}(\mathbf{A}_{TT}^{-1})_{T \in \mathcal{T}_h}. \quad (2.62)$$

The above remark can be exploited in practice to efficiently eliminate the element-based displacement DOFs from the global system. This process, usually referred to as ‘‘static condensation’’, is detailed in what follows.

For a given time step  $1 \leq n \leq N$ , the linear system corresponding to the discrete problem (2.26) is of the form

$$\begin{bmatrix} \mathbf{A}_{\mathcal{T}\mathcal{T}} & \mathbf{A}_{\mathcal{T}\mathcal{F}} & & \mathbf{B}_{\mathcal{T}} \\ \mathbf{A}_{\mathcal{T}\mathcal{F}}^{\top} & \mathbf{A}_{\mathcal{F}\mathcal{F}} & & \mathbf{B}_{\mathcal{F}} \\ & & & \\ -\mathbf{B}_{\mathcal{T}}^{\top} & -\mathbf{B}_{\mathcal{F}}^{\top} & \frac{\tau}{\theta} \mathbf{C} + c_0 \mathbf{M} & \end{bmatrix} \begin{bmatrix} \mathbf{U}_{\mathcal{T}}^n \\ \mathbf{U}_{\mathcal{F}}^n \\ \mathbf{P}^n \end{bmatrix} = \begin{bmatrix} \mathbf{F}_{\mathcal{T}}^n \\ \mathbf{0}_{\mathcal{F}} \\ \widetilde{\mathbf{G}}^n \end{bmatrix}, \quad (2.63)$$

where  $\mathbf{C}$  denotes the matrix that represents the bilinear form  $c_h$  in the selected basis,  $\mathbf{M}$  is the (block diagonal) pressure mass matrix,  $\mathbf{F}_{\mathcal{T}}^n$  is the vector corresponding to the discretization of the volumetric load  $\mathbf{f}^n$ , while  $\mathbf{0}_{\mathcal{F}}$  is the zero vector of length  $\dim(\mathbf{U}_{\mathcal{F}}^k)$ . Denoting by  $\mathbf{G}^n$  the vector corresponding to the discretization of the fluid source  $g^n$ , when the backward Euler method is used to march in time, we let  $\theta = 1$  and set

$$\widetilde{\mathbf{G}}^n := \tau \mathbf{G}^n - \mathbf{B}^{\top} \mathbf{U}^{n-1} + c_0 \mathbf{M} \mathbf{P}^{n-1}.$$

For the BDF2 method (and  $n \geq 2$ ), we let  $\theta = 3/2$  and set

$$\tilde{\mathbf{G}}^n := \frac{2}{3}\tau\mathbf{G}^n - \frac{4}{3}\mathbf{B}^T\mathbf{U}^{n-1} + \frac{1}{3}\mathbf{B}^T\mathbf{U}^{n-2} + \frac{4}{3}c_0\mathbf{M}\mathbf{P}^{n-1} - \frac{1}{3}c_0\mathbf{M}\mathbf{P}^{n-2}.$$

Recalling (2.62), instead of assembling the full system, we can effectively compute the Schur complement of  $\mathbf{A}_{\mathcal{T}\mathcal{T}}$  and code, instead, the following reduced version, where the element-based displacement DOFs collected in the subvector  $\mathbf{U}_{\mathcal{T}}^n$  no longer appear:

$$\left[ \begin{array}{c|c} \mathbf{A}_{\mathcal{F}\mathcal{F}} - \mathbf{A}_{\mathcal{T}\mathcal{F}}^T \mathbf{A}_{\mathcal{T}\mathcal{T}}^{-1} \mathbf{A}_{\mathcal{T}\mathcal{F}} & \mathbf{B}_{\mathcal{F}} - \mathbf{A}_{\mathcal{T}\mathcal{F}}^T \mathbf{A}_{\mathcal{T}\mathcal{T}}^{-1} \mathbf{B}_{\mathcal{T}} \\ \hline -\mathbf{B}_{\mathcal{F}}^T + \mathbf{B}_{\mathcal{T}}^T \mathbf{A}_{\mathcal{T}\mathcal{T}}^{-1} \mathbf{A}_{\mathcal{T}\mathcal{F}} & \frac{\tau}{\theta}\mathbf{C} + c_0\mathbf{M} + \mathbf{B}_{\mathcal{T}}^T \mathbf{A}_{\mathcal{T}\mathcal{T}}^{-1} \mathbf{B}_{\mathcal{T}} \end{array} \right] \begin{bmatrix} \mathbf{U}_{\mathcal{F}}^n \\ \mathbf{P}^n \end{bmatrix} = \begin{bmatrix} -\mathbf{A}_{\mathcal{T}\mathcal{F}}^T \mathbf{A}_{\mathcal{T}\mathcal{T}}^{-1} \mathbf{F}_{\mathcal{T}}^n \\ \tilde{\mathbf{G}}^n + \mathbf{B}_{\mathcal{T}}^T \mathbf{A}_{\mathcal{T}\mathcal{T}}^{-1} \mathbf{F}_{\mathcal{T}}^n \end{bmatrix}. \quad (2.64)$$

All matrix products appearing in (2.64) are directly assembled from their local counterparts (i.e., the factors need not be constructed separately). Specifically, introducing, for all  $T \in \mathcal{T}_h$ , the following local matrices  $\mathbf{A}(T)$  and  $\mathbf{B}(T)$  representing the local bilinear forms  $a_T$  (cf. (2.9)) and  $b_T$  (cf. (2.21)), respectively:

$$\mathbf{A}(T) = \begin{bmatrix} \mathbf{A}_{TT} & \mathbf{A}_{T\mathcal{F}_T} \\ \mathbf{A}_{T\mathcal{F}_T}^T & \mathbf{A}_{\mathcal{F}_T\mathcal{F}_T} \end{bmatrix}, \quad \mathbf{B}(T) = \begin{bmatrix} \mathbf{B}_T \\ \mathbf{B}_{\mathcal{F}_T} \end{bmatrix},$$

one has for the left-hand side matrix, denoting by  $\longleftarrow_{T \in \mathcal{T}_h}$  the usual assembly procedure based on a global DOF map,

$$\begin{aligned} \mathbf{A}_{\mathcal{T}\mathcal{F}}^T \mathbf{A}_{\mathcal{T}\mathcal{T}}^{-1} \mathbf{A}_{\mathcal{T}\mathcal{F}} &\longleftarrow_{T \in \mathcal{T}_h} \mathbf{A}_{T\mathcal{F}_T}^T \mathbf{A}_{TT}^{-1} \mathbf{A}_{T\mathcal{F}_T}, & \mathbf{A}_{\mathcal{T}\mathcal{F}}^T \mathbf{A}_{\mathcal{T}\mathcal{T}}^{-1} \mathbf{B}_{\mathcal{T}} &\longleftarrow_{T \in \mathcal{T}_h} \mathbf{A}_{T\mathcal{F}_T}^T \mathbf{A}_{TT}^{-1} \mathbf{B}_T, \\ \mathbf{B}_{\mathcal{F}}^T \mathbf{A}_{\mathcal{T}\mathcal{T}}^{-1} \mathbf{A}_{\mathcal{T}\mathcal{F}} &\longleftarrow_{T \in \mathcal{T}_h} \mathbf{B}_T^T \mathbf{A}_{TT}^{-1} \mathbf{A}_{T\mathcal{F}_T}, & \mathbf{B}_{\mathcal{F}}^T \mathbf{A}_{\mathcal{T}\mathcal{T}}^{-1} \mathbf{B}_{\mathcal{T}} &\longleftarrow_{T \in \mathcal{T}_h} \mathbf{B}_T^T \mathbf{A}_{TT}^{-1} \mathbf{B}_T, \end{aligned}$$

and, similarly, for the right-hand side vector

$$\mathbf{A}_{\mathcal{T}\mathcal{F}}^T \mathbf{A}_{\mathcal{T}\mathcal{T}}^{-1} \mathbf{F}_{\mathcal{T}}^n \longleftarrow_{T \in \mathcal{T}_h} \mathbf{A}_{T\mathcal{F}_T}^T \mathbf{A}_{TT}^{-1} \mathbf{F}_T^n, \quad \mathbf{B}_{\mathcal{F}}^T \mathbf{A}_{\mathcal{T}\mathcal{T}}^{-1} \mathbf{F}_{\mathcal{T}}^n \longleftarrow_{T \in \mathcal{T}_h} \mathbf{B}_T^T \mathbf{A}_{TT}^{-1} \mathbf{F}_T^n.$$

The advantage of implementing (2.64) over (2.63) is that the number of DOFs appearing in the linear system reduces to (compare with (2.61))

$$d \operatorname{card}(\mathcal{F}_h^i) N_{d-1}^k + \operatorname{card}(\mathcal{T}_h) N_d^k.$$

Additionally, since the reduced left-hand side matrix in (2.64) does not depend on the time step  $n$ , it can be assembled (and, possibly, factored) once and for all in a preliminary stage, thus leading to a further reduction in the computational cost. Finally, for all  $T \in \mathcal{T}_h$ , the local vector  $\mathbf{U}_T^n$  of element-based displacement DOFs can be recovered from the local right-hand side vector  $\mathbf{F}_T^n$  and the local vector of face-based displacement DOFs and (element-based) pressure DOFs  $(\mathbf{U}_{\mathcal{F}_T}^n, \mathbf{P}_T^n)$  by the following element-by-element post-processing:

$$\mathbf{U}_T^n = \mathbf{A}_{TT}^{-1} \left( \mathbf{F}_T^n - \mathbf{A}_{T\mathcal{F}_T} \mathbf{U}_{\mathcal{F}_T}^n - \mathbf{B}_T \mathbf{P}_T^n \right).$$

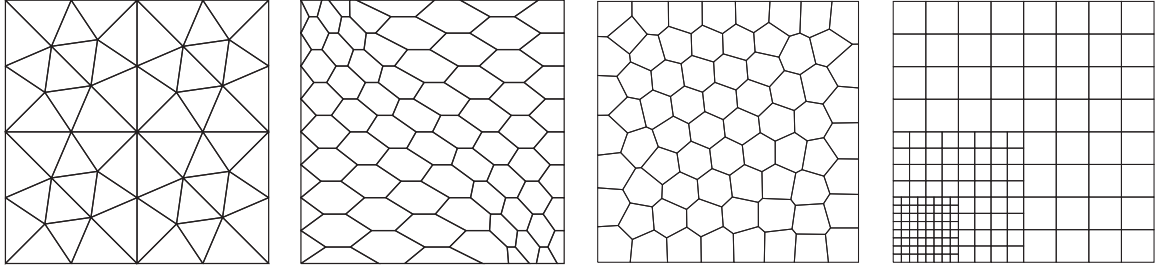


Figure 2.1: Triangular, hexagonal-dominant, Voronoi, and nonmatching quadrangular meshes for the numerical tests. The triangular and nonmatching quadrangular meshes were originally proposed for the FVCA5 benchmark [120], whereas the hexagonal-dominant mesh is the same used in [80, Section 4.2.3].

## 2.6 Numerical tests

In this section we present a comprehensive set of numerical tests to assess the properties of our method.

### 2.6.1 Convergence

We first consider a manufactured regular exact solution to confirm the convergence rates predicted in (2.52). Specifically, we solve the two-dimensional incompressible Biot problem ( $c_0 = 0$ ) in the unit square domain  $\Omega = (0, 1)^2$  with  $t_F = 1$  and physical parameters  $\mu = 1$ ,  $\lambda = 1$ , and  $\kappa = 1$ . The exact displacement  $\mathbf{u}$  and exact pressure  $p$  are given by, respectively

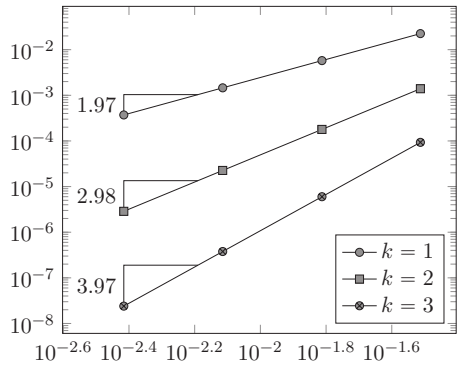
$$\begin{aligned} \mathbf{u}(\mathbf{x}, t) &= (-\sin(\pi t) \cos(\pi x_1) \cos(\pi x_2), \sin(\pi t) \sin(\pi x_1) \sin(\pi x_2)), \\ p(\mathbf{x}, t) &= -\cos(\pi t) \sin(\pi x_1) \cos(\pi x_2). \end{aligned}$$

The volumetric load is given by

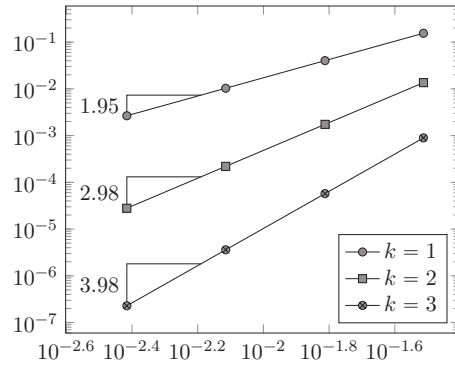
$$\mathbf{f}(\mathbf{x}, t) = 6\pi^2(\sin(\pi t) + \pi \cos(\pi t)) \times (-\cos(\pi x_1) \cos(\pi x_2), \sin(\pi x_1) \sin(\pi x_2)),$$

while  $g(\mathbf{x}, t) \equiv 0$ . Dirichlet boundary conditions for the displacement and Neumann boundary conditions for the pressure are inferred from exact solutions to  $\partial\Omega$ .

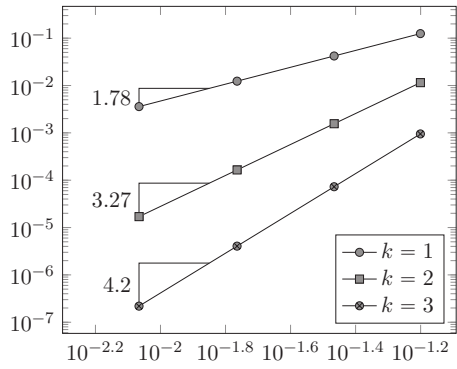
We consider the triangular, (predominantly) hexagonal, Voronoi, and nonmatching quadrangular mesh families depicted in Figure 2.1. The Voronoi mesh family was obtained using the PolyMesher algorithm of [188]. The nonmatching mesh is simply meant to show that the method supports nonconforming interfaces: refining in the corner has no particular meaning for the selected solution. The time discretization is based on the second order Backward Differentiation Formula (BDF2); cf. Remark 2.15. The time step  $\tau$  on the coarsest mesh is taken to be  $0.1/2^{\frac{(k+1)}{2}}$  for every choice of the spatial degree  $k$ , and it decreases with the mesh size  $h$  according to the theoretical convergence rates, thus, if  $h_2 = h_1/2$ , then  $\tau_2 = \tau_1/2^{\frac{(k+1)}{2}}$ . Figure 2.2 displays convergence results for the various mesh families and polynomial degrees up to 3. The error measures are  $\|p_h^N - \pi_h^k p^N\|$  for the pressure and  $\|\underline{\mathbf{u}}_h^N - \underline{I}_h^k \mathbf{u}^N\|_{a,h}$  for the displacement. Using the triangle inequality together with (2.52) and the approximation properties (2.2) of  $\pi_h^k$  and (2.6) of  $(\mathbf{r}_h^{k+1} \circ \underline{I}_h^k)$ , it is a simple matter to prove that these quantities



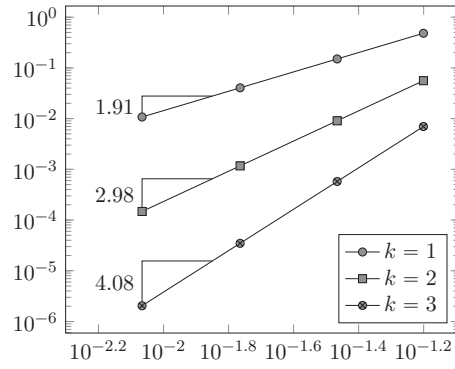
(a)  $\|p_h^N - \pi_h^k p^N\|$ , triangular



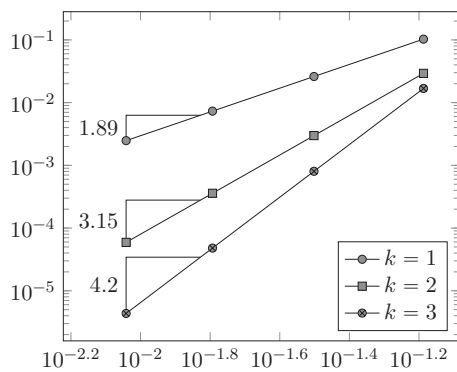
(b)  $\|\underline{u}_h^N - \underline{I}_h^k \underline{u}^N\|_{a,h}$ , triangular



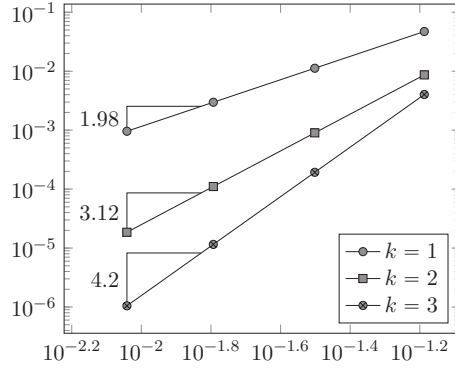
(c)  $\|p_h^N - \pi_h^k p^N\|$ , hexagonal



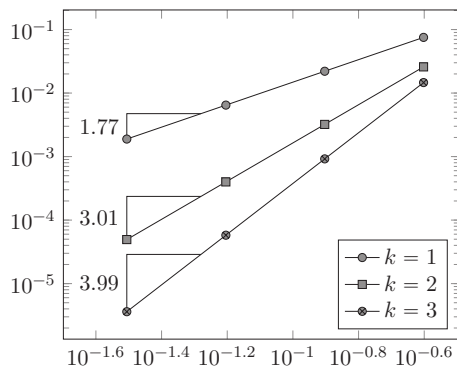
(d)  $\|\underline{u}_h^N - \underline{I}_h^k \underline{u}^N\|_{a,h}$ , hexagonal



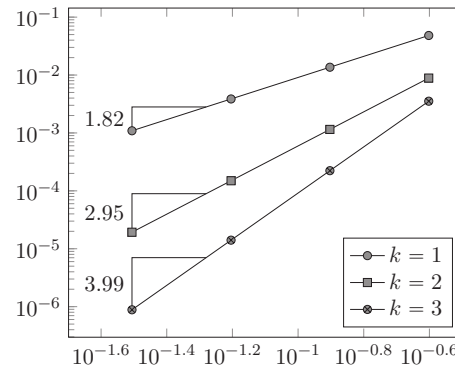
(e)  $\|p_h^N - \pi_h^k p^N\|$ , Voronoi



(f)  $\|\underline{u}_h^N - \underline{I}_h^k \underline{u}^N\|_{a,h}$ , Voronoi



(g)  $\|p_h^N - \pi_h^k p^N\|$ , nonmatching



(h)  $\|\underline{u}_h^N - \underline{I}_h^k \underline{u}^N\|_{a,h}$ , nonmatching

Figure 2.2: Errors vs.  $h$

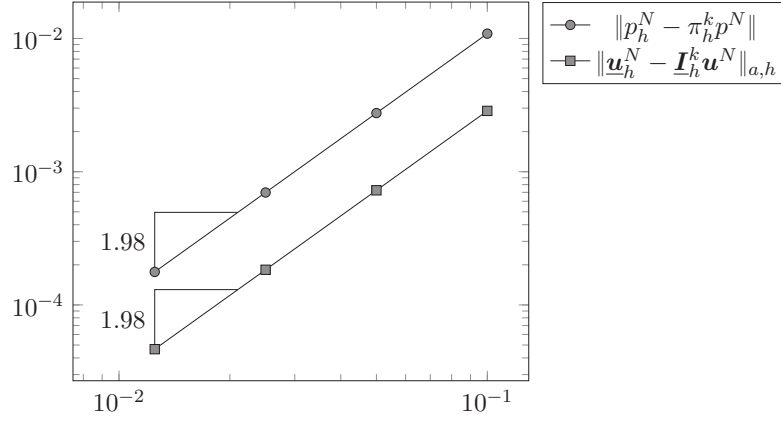


Figure 2.3: Time convergence rate with BDF2, hexagonal mesh

have the same convergence behaviour as the terms in the left-hand side of (2.52). In all the cases, the numerical results show asymptotic convergence rates that are in agreement with theoretical predictions. This test was also used to numerically check that the mechanical equilibrium and mass conservation relations of Lemma 2.18 hold up to machine precision.

The convergence in time was also separately checked considering the method with spatial degree  $k = 3$  on the hexagonal mesh with mesh size  $h = 0.0172$  and time step decreasing from  $\tau = 0.1$  to  $\tau = 0.0125$ . With this choice, the time-component of the error is dominant, and Figure 2.3 confirms the second order convergence of the BDF2 scheme.

## 2.6.2 Barry and Mercer's test case

A test case more representative of actual physical configurations is that of Barry and Mercer [14], for which an exact solution is available (we refer to the cited paper and also to [164, Section 4.2.1] for its expression). We let  $\Omega = (0, 1)^2$  and consider the following time-independent boundary conditions on  $\partial\Omega$

$$\mathbf{u} \cdot \boldsymbol{\tau} = 0, \quad \mathbf{n}^T \nabla \mathbf{u} \mathbf{n} = 0, \quad p = 0,$$

where  $\boldsymbol{\tau}$  denotes the tangent vector on  $\partial\Omega$ . The evolution of the displacement and pressure fields is driven by a periodic pointwise source (mimicking a well) located at  $\mathbf{x}_0 = (0.25, 0.25)$ :

$$g = \delta(\mathbf{x} - \mathbf{x}_0) \sin(\hat{t}),$$

with normalized time  $\hat{t} := \beta t$  for  $\beta := (\lambda + 2\mu)\kappa$ . As in [166, 176], we use the following values for the physical parameters:

$$c_0 = 0, \quad E = 1 \cdot 10^5, \quad \nu = 0.1, \quad \kappa = 1 \cdot 10^{-2},$$

where  $E$  and  $\nu$  denote Young's modulus and Poisson ratio, respectively, and

$$\lambda = \frac{E\nu}{(1+\nu)(1-2\nu)}, \quad \mu = \frac{E}{2(1+\nu)}.$$

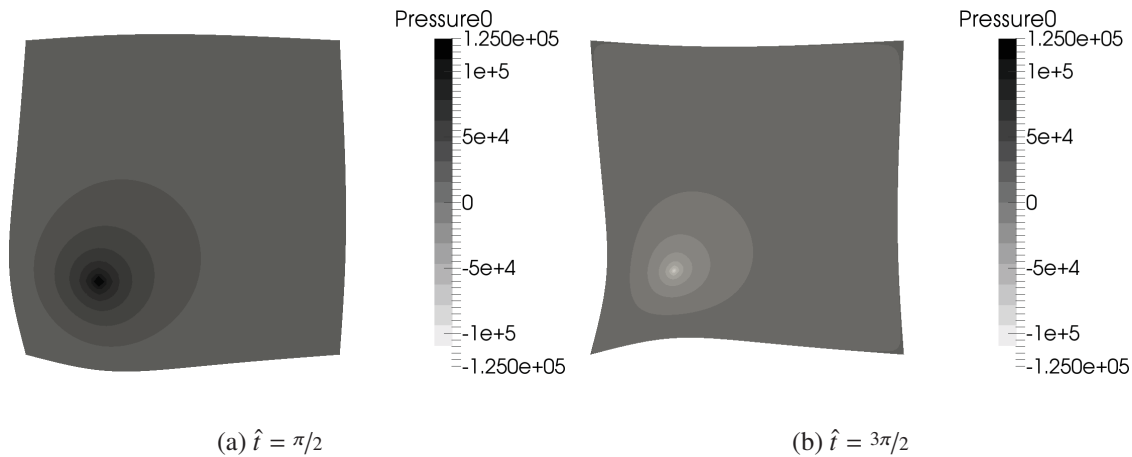


Figure 2.4: Pressure field on the deformed domain at different times for the finest Cartesian mesh containing 4,192 elements

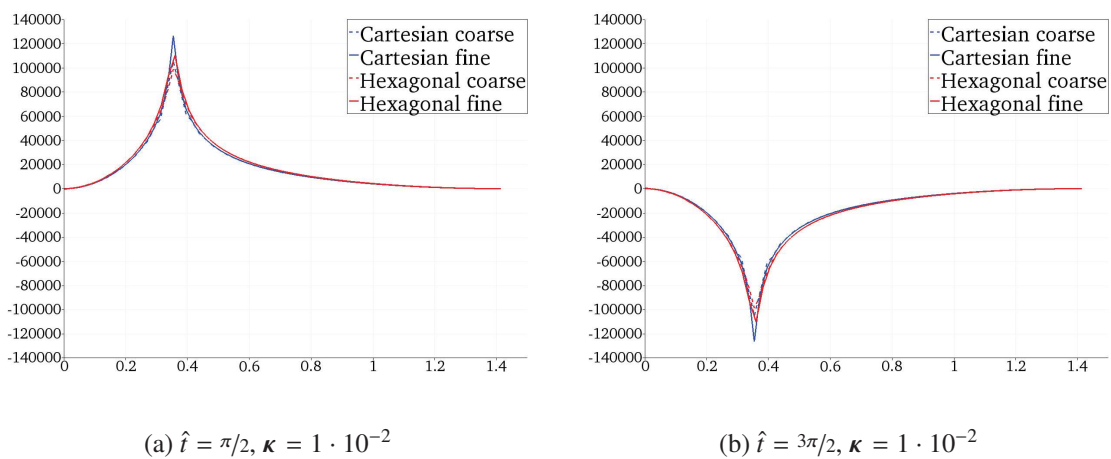


Figure 2.5: Pressure profiles along the diagonal  $(0,0)-(1,1)$  of the domain for different normalized times  $\hat{t}$  and meshes ( $k=1$ ). The time step is here  $\tau = (2\pi/\beta) \cdot 10^{-2}$ .

In the injection phase  $\hat{t} \in (0, \pi)$ , we observe an inflation of the domain, which reaches its maximum at  $\hat{t} = \pi/2$ ; cf. Figure 2.4a. In the extraction phase  $\hat{t} \in (\pi, 2\pi)$ , on the other hand, we have a contraction of the domain which reaches its maximum at  $\hat{t} = 3\pi/2$ ; cf. Figure 2.4b.

The following results have been obtained with the lowest-order version of the method corresponding to  $k=1$  (taking advantage of higher orders would require local mesh refinement, which is out of the scope of the present work). In Figure 2.5 we plot the pressure profile at normalized times  $\hat{t} = \pi/2$  and  $\hat{t} = 3\pi/2$  along the diagonal  $(0,0)-(1,1)$  of the domain. We consider two Cartesian meshes containing 1,024 and 4,096 elements, respectively, as well as two (predominantly) hexagonal meshes containing 1,072 and 4,192 elements, respectively. In all the cases, a time step  $\tau = (2\pi/\beta) \cdot 10^{-2}$  is used. We note that the behaviour of the pressure is well-captured even on the coarsest meshes. For the finest hexagonal mesh, the relative error on the pressure in the  $L^2$ -norm at times  $\hat{t} = \pi/2$  and  $\hat{t} = 3\pi/2$  is 2.85%.

To check the robustness of the method with respect to pressure oscillations for small

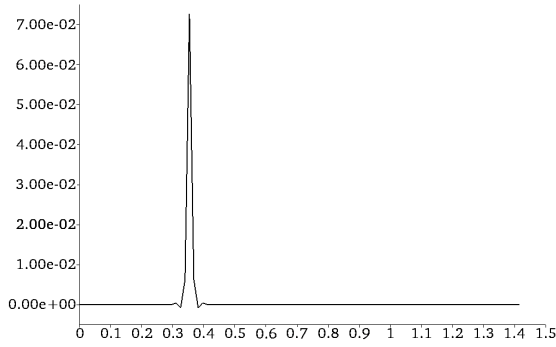
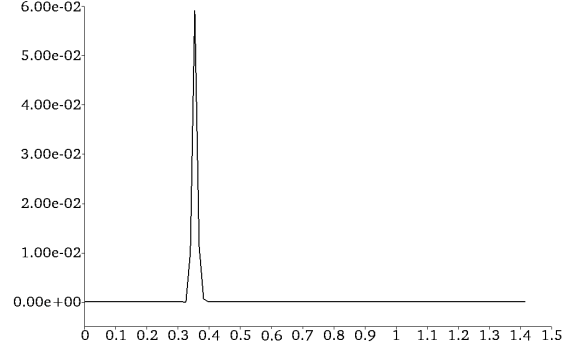
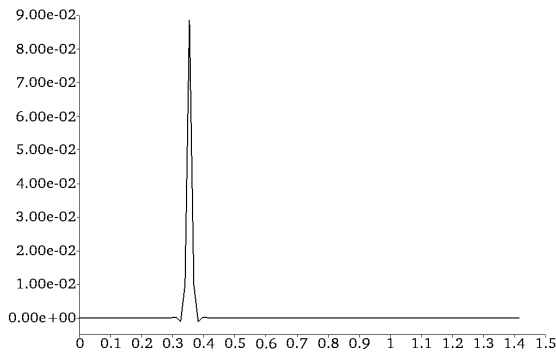
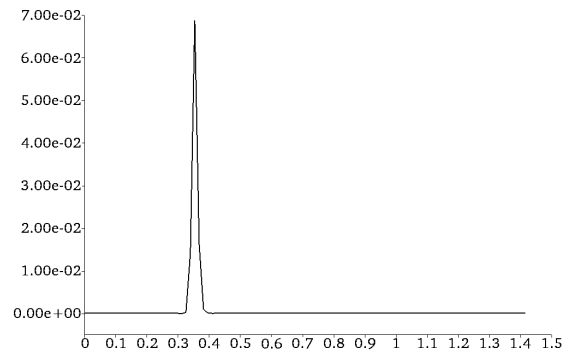
(a) Cartesian mesh ( $\text{card}(\mathcal{T}_h) = 4,028$ ), first step(b) Hexagonal mesh ( $\text{card}(\mathcal{T}_h) = 4,192$ ), first step(c) Cartesian mesh ( $\text{card}(\mathcal{T}_h) = 4,028$ ), second step(d) Hexagonal mesh ( $\text{card}(\mathcal{T}_h)=4,192$ ), second step

Figure 2.6: Pressure profiles along the diagonal  $(0,0)-(1,1)$  of the domain for  $\kappa = 1 \cdot 10^{-6}$  and time step  $\tau = 1 \cdot 10^{-4}$ . Small oscillations are present on the Cartesian mesh (left), whereas no sign of oscillations is present on the hexagonal mesh (right).

permeabilities combined with small time steps, we also show in Figure 2.6 the pressure profile after one and two step with  $\kappa = 1 \cdot 10^{-6}$  and  $\tau = 1 \cdot 10^{-4}$  on the Cartesian and hexagonal meshes with 4,096 and 4,192 elements, respectively. We remark that the first time step is performed using the backward Euler scheme, while the second with the second order BDF2 scheme. This situation corresponds to the one considered in [176, Figure 5.10] to highlight the onset of spurious oscillations. In our case, small oscillations can be observed for the Cartesian mesh (cf. Figure 2.6a and Figure 2.6c), whereas no sign of oscillations is present for the hexagonal mesh (cf. Figure 2.6b and Figure 2.6d). One possible conjecture is that increasing the number of element faces contributes to the monotonicity of the scheme.

## 2.7 Appendix: Flux formulation

In this section, we reformulate the discrete problem (2.26) to unveil the local conservation properties of the method. Before doing so, we need to introduce a few operators and some

notation to treat the boundary terms.

We start from the mechanical equilibrium. Let an element  $T \in \mathcal{T}_h$  be fixed and denote by  $\mathbf{U}_{\partial T} := \mathbb{P}_{d-1}^k(\mathcal{F}_T)^d$  the broken polynomial space of degree  $\leq k$  on the boundary  $\partial T$  of  $T$ . We define the boundary operator  $\mathbf{L}_T^k : \mathbf{U}_{\partial T} \rightarrow \mathbf{U}_{\partial T}$  such that, for all  $\boldsymbol{\varphi} \in \mathbf{U}_{\partial T}$ ,

$$\mathbf{L}_T^k \boldsymbol{\varphi}|_F := \pi_F^k \left( \boldsymbol{\varphi}|_F - \mathbf{r}_T^{k+1}(\mathbf{0}, (\boldsymbol{\varphi}|_F)_{F \in \mathcal{F}_T}) + \pi_T^k \mathbf{r}_T^{k+1}(\mathbf{0}, (\boldsymbol{\varphi}|_F)_{F \in \mathcal{F}_T}) \right) \quad \forall F \in \mathcal{F}_T.$$

We also need the adjoint  $\mathbf{L}_T^{k,*}$  of  $\mathbf{L}_T^k$  such that

$$\forall \boldsymbol{\varphi} \in \mathbf{U}_{\partial T}, \quad (\mathbf{L}_T^k \boldsymbol{\varphi}, \boldsymbol{\psi})_{\partial T} = (\boldsymbol{\varphi}, \mathbf{L}_T^{k,*} \boldsymbol{\psi})_{\partial T} \quad \forall \boldsymbol{\psi} \in \mathbf{U}_{\partial T}.$$

For a collection of DOFs  $\underline{\mathbf{v}}_T \in \underline{\mathbf{U}}_T^k$ , we denote in what follows by  $\mathbf{v}_{\partial T} \in \mathbf{U}_{\partial T}$  the function in  $\mathbf{U}_{\partial T}$  such that  $\mathbf{v}_{\partial T}|_F = \mathbf{v}_F$  for all  $F \in \mathcal{F}_T$ . Finally, it is convenient to define the discrete stress operator  $\mathbf{S}_T^k : \underline{\mathbf{U}}_T^k \rightarrow \mathbb{P}_d^k(T)^{d \times d}$  such that, for all  $\underline{\mathbf{v}}_T \in \underline{\mathbf{U}}_T^k$ ,

$$\mathbf{S}_T^k \underline{\mathbf{v}}_T := 2\mu \nabla_s \mathbf{r}_T^{k+1} \underline{\mathbf{v}}_T + \lambda \mathbf{I}_d D_T^k \underline{\mathbf{v}}_T. \quad (2.65)$$

To reformulate the mass conservation equation, we need to introduce the classical lifting operator  $R_{\kappa,h}^k : P_h^k \rightarrow \mathbb{P}_d^{k-1}(\mathcal{T}_h)^d$  such that, for all  $q_h \in P_h^k$ , it holds

$$(R_{\kappa,h}^k q_h, \boldsymbol{\xi}_h) = \sum_{F \in \mathcal{F}_h^i} ([q_h]_F, \{\boldsymbol{\kappa} \boldsymbol{\xi}_h\}_F \cdot \mathbf{n}_F)_F \quad \forall \boldsymbol{\xi}_h \in \mathbb{P}_d^{k-1}(\mathcal{T}_h)^d.$$

**Lemma 2.18** (Flux formulation of problem (2.26)). *Problem (2.26) can be reformulated as follows: Find  $(\underline{\mathbf{u}}_h^n, p_h^n) \in \underline{\mathbf{U}}_{h,D}^k \times P_h^k$  such that it holds, for all  $(\underline{\mathbf{v}}_h, q_h) \in \underline{\mathbf{U}}_{h,D}^k \times \mathbb{P}_d^k(\mathcal{T}_h)$  and all  $T \in \mathcal{T}_h$ ,*

$$(\mathbf{S}_T^k \underline{\mathbf{u}}_T^n - p_h^n \mathbf{I}_d, \nabla_s \mathbf{v}_T)_T + \sum_{F \in \mathcal{F}_T} (\boldsymbol{\Phi}_{TF}^k(\underline{\mathbf{u}}_T^n, p_h^n|_T), \mathbf{v}_F - \mathbf{v}_T)_F = (\mathbf{f}^n, \mathbf{v}_T)_T, \quad (2.66a)$$

$$(c_0 \delta_t p_h^n, q_h)_T - (\delta_t \mathbf{u}_T^n - \boldsymbol{\kappa}(\nabla_h p_h^n - R_{\kappa,h}^k p_h^n), \nabla_h q_h)_T - \sum_{F \in \mathcal{F}_T} (\phi_{TF}^k(\delta_t \mathbf{u}_F^n, p_h^n), q_h|_T)_F = (g^n, q_h), \quad (2.66b)$$

where, for all  $T \in \mathcal{T}_h$  and all  $F \in \mathcal{F}_T$ , the numerical traction  $\boldsymbol{\Phi}_{TF}^k : \underline{\mathbf{U}}_T^k \times \mathbb{P}_d^k(T) \rightarrow \mathbb{P}_{d-1}^k(F)^d$  and mass flux  $\phi_{TF}^k : \mathbb{P}_{d-1}^k(F)^d \times \mathbb{P}_d^k(\mathcal{T}_h) \rightarrow \mathbb{P}_{d-1}^k(F)$  are such that

$$\begin{aligned} \boldsymbol{\Phi}_{TF}^k(\underline{\mathbf{v}}_T, q) &:= (\mathbf{S}_T^k \underline{\mathbf{v}}_T - q \mathbf{I}_d) \mathbf{n}_{TF} + (2\mu) \mathbf{L}_T^{k,*} (\mathfrak{h}_{\partial T}^{-1} \mathbf{L}_T^k (\mathbf{v}_{\partial T} - \mathbf{v}_T)), \\ \phi_{TF}^k(\mathbf{v}_F, q_h) &:= \begin{cases} (-\mathbf{v}_F^n + \{\boldsymbol{\kappa} \nabla_h q_h\}_F) \cdot \mathbf{n}_{TF} - \frac{\varsigma \lambda_{\kappa,F}}{h_F} [q_h]_F (\mathbf{n}_{TF} \cdot \mathbf{n}_F) & \text{if } F \in \mathcal{F}_h^i, \\ 0 & \text{otherwise,} \end{cases} \end{aligned} \quad (2.67)$$

with  $\mathfrak{h}_{\partial T} \in \mathbb{P}_d^0(\mathcal{F}_T)$  such that  $\mathfrak{h}_{\partial T}|_F = h_F$  for all  $F \in \mathcal{F}_T$ , and it holds, for all  $F \in \mathcal{F}_h^i$  such that  $F \in \mathcal{F}_{T_1} \cap \mathcal{F}_{T_2}$ ,

$$\boldsymbol{\Phi}_{T_1 F}^k(\underline{\mathbf{u}}_{T_1}^n, p_h^n|_{T_1}) + \boldsymbol{\Phi}_{T_2 F}^k(\underline{\mathbf{u}}_{T_2}^n, p_h^n|_{T_2}) = \mathbf{0} \quad (2.68a)$$

$$\phi_{T_1 F}^k(\delta_t \mathbf{u}_F^n, p_h^n) + \phi_{T_2 F}^k(\delta_t \mathbf{u}_F^n, p_h^n) = 0. \quad (2.68b)$$



*Proof.* (1) *Proof of (2.66a).* Proceeding as in [57, Section 3.1], the stabilization bilinear form  $s_T$  defined by (2.10) can be rewritten as

$$s_T(\underline{\mathbf{w}}_T, \underline{\mathbf{v}}_T) = \sum_{F \in \mathcal{F}_T} (\mathbf{L}_T^{k,*} (\mathfrak{h}_{\partial T}^{-1} \mathbf{L}_T^k (\mathbf{w}_{\partial T} - \mathbf{w}_T)), \mathbf{v}_F - \mathbf{v}_T)_F.$$

Therefore, using the definitions (2.5) of  $\mathbf{r}_T^{k+1} \underline{\mathbf{v}}_T$  with  $\mathbf{w} = \mathbf{r}_T^{k+1} \underline{\mathbf{u}}_T^n$  and (2.7a) of  $D_T^k \underline{\mathbf{v}}_T$  with  $q = p_h^n|_T$ , and recalling the definition (2.65) of  $\mathbf{S}_T^k$ , one has

$$a_T(\underline{\mathbf{u}}_T^n, \underline{\mathbf{v}}_T) = (\mathbf{S}_T^k \underline{\mathbf{u}}_T^n, \nabla_s \mathbf{v}_T)_T + \sum_{F \in \mathcal{F}_T} (\mathbf{S}_T^k \underline{\mathbf{u}}_T^n \mathbf{n}_{TF} + (2\mu) \mathbf{L}_T^{k,*} (\mathfrak{h}_{\partial T}^{-1} \mathbf{L}_T^k (\underline{\mathbf{u}}_{\partial T}^n - \underline{\mathbf{u}}_T^n)), \mathbf{v}_F - \mathbf{v}_T)_F. \quad (2.69)$$

On the other hand, using again the definition (2.7a) of  $D_T^k \underline{\mathbf{v}}_T$  with  $q = p_h^n|_T$ , one has

$$b_T(\underline{\mathbf{v}}_T, p_h^n|_T) = -(p_h^n \mathbf{I}_d, \nabla_s \mathbf{v}_T)_T - \sum_{F \in \mathcal{F}_T} (p_h^n|_T \mathbf{n}_{TF}, \mathbf{v}_F - \mathbf{v}_T)_F. \quad (2.70)$$

Equation (2.66a) follows summing (2.69) and (2.70).

(2) *Proof of (2.66b).* Using the definition (2.7b) of  $D_T^k$  with  $\underline{\mathbf{v}}_T = \delta_t \underline{\mathbf{u}}_T^n$  and  $q = q_h|_T$ , it is inferred that

$$b_T(\delta_t \underline{\mathbf{u}}_T^n, q_h) = -(\delta_t \underline{\mathbf{u}}_T^n, \nabla_h q_h)_T + \sum_{F \in \mathcal{F}_T} (\delta_t \underline{\mathbf{u}}_F^n \cdot \mathbf{n}_{TF}, q_h|_T)_F. \quad (2.71)$$

On the other hand, adapting the results [76, Section 4.5.5] to the homogeneous Neumann boundary condition (2.1d), it is inferred

$$c_h(p_h^n, q_h) = \sum_{T \in \mathcal{T}_h} \left\{ (\boldsymbol{\kappa} (\nabla_h p_h^n - R_{\boldsymbol{\kappa}, h}^k p_h^n) \cdot \nabla_h q_h)_T - \sum_{F \in \mathcal{F}_T \cap \mathcal{F}_h^i} (\{\boldsymbol{\kappa} \nabla_h p_h^n\}_F \cdot \mathbf{n}_{TF} - \frac{\mathcal{S} \lambda_{\boldsymbol{\kappa}, F}}{h_F} [p_h^n]_F (\mathbf{n}_{TF} \cdot \mathbf{n}_F), q_h|_T)_F \right\}. \quad (2.72)$$

Equation (2.66b) follows summing (2.71) and (2.72).

(3) *Proof of (2.68).* To prove (2.68a), let an internal face  $F \in \mathcal{F}_h^i$  be fixed, and make  $\underline{\mathbf{v}}_h$  in (2.68a) such that  $\mathbf{v}_T \equiv \mathbf{0}$  for all  $T \in \mathcal{T}_h$ ,  $\mathbf{v}_{F'} \equiv \mathbf{0}$  for all  $F' \in \mathcal{F}_h \setminus \{F\}$ , let  $\mathbf{v}_F$  span  $\mathbb{P}_{d-1}^k(F)$  and rearrange the sums. The mass flux conservation (2.68b) follows immediately from the expression of  $\phi_{TF}^k$  observing that, for all  $(\underline{\mathbf{v}}_h, q_h) \in \underline{\mathbf{U}}_h^k \times P_h^k$  and all  $F \in \mathcal{F}_h^i$ , the quantity

$$(-\mathbf{v}_F + \{\boldsymbol{\kappa} \nabla_h q_h\}_F) \cdot \mathbf{n}_F - \frac{\mathcal{S} \lambda_{\boldsymbol{\kappa}, F}}{h_F} [q_h]_F$$

is single-valued on  $F$ . □

Let now an element  $T \in \mathcal{T}_h$  be fixed. Choosing as test functions in (2.66a)  $\underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_{h,D}^k$  such that  $\mathbf{v}_F \equiv \mathbf{0}$  for all  $F \in \mathcal{F}_h$ ,  $\mathbf{v}_{T'} \equiv \mathbf{0}$  for all  $T' \in \mathcal{T}_h \setminus \{T\}$ , and  $\mathbf{v}_T$  spans  $\mathbb{P}_d^k(T)^d$ , we infer the following local mechanical equilibrium relation: For all  $\mathbf{v}_T \in \mathbb{P}_d^k(T)^d$ ,

$$(\mathbf{S}_T^k \underline{\mathbf{u}}_T^n - p_h^n \mathbf{I}_d, \nabla_s \mathbf{v}_T)_T - \sum_{F \in \mathcal{F}_T} (\boldsymbol{\Phi}_{TF}^k (\underline{\mathbf{u}}_T^n, p_h^n|_T), \mathbf{v}_T)_F = (\mathbf{f}^n, \mathbf{v}_T)_T.$$

Similarly, selecting  $q_h$  in (2.66b) such that  $q_{h|T'} \equiv 0$  for all  $T' \in \mathcal{T}_h \setminus \{T\}$  and  $q_T := q_{h|T}$  spans  $\mathbb{P}_d^k(T)$ , we infer the following local mass conservation relation: For all  $q_T \in \mathbb{P}_d^k(T)$ ,

$$(c_0 \delta_t p_h^n, q_T)_T - (\delta_t \mathbf{u}_T^n - \kappa(\nabla_h p_h^n - R_{\kappa,h}^k p_h^n), \nabla q_T)_T - \sum_{F \in \mathcal{F}_T} (\phi_{TF}^k(\delta_t \mathbf{u}_F^n, p_h^n), q_T)_F = (g^n, q_T).$$

To actually compute the numerical fluxes defined by (2.67), besides the operator  $\mathbf{S}_T^k$  defined by (2.65) (which is readily available once  $\mathbf{r}_T^{k+1}$  and  $D_T^k$  have been computed; cf. (2.5) and (2.7a), respectively), one also needs to compute the operators  $L_T^k$  and  $L_T^{k,*}$ . The latter operation can be performed at marginal cost, since it only requires to invert the face mass matrices of  $\mathbb{P}_{d-1}^k(F)$  for all  $F \in \mathcal{F}_T$ .



# Chapter 3

---

## Nonlinear elasticity

---

This chapter has been published in the following peer-reviewed journal (see [33]):

**SIAM Journal on Numerical Analysis,**  
Volume 55, Issue 6, 2017, Pages 2687–2717.

### Contents

---

<b>3.1</b>	<b>Introduction</b> . . . . .	<b>50</b>
<b>3.2</b>	<b>Setting and examples</b> . . . . .	<b>52</b>
<b>3.3</b>	<b>Notation and basic results</b> . . . . .	<b>54</b>
<b>3.4</b>	<b>The Hybrid High-Order method</b> . . . . .	<b>55</b>
3.4.1	Degrees of freedom . . . . .	55
3.4.2	Local reconstructions . . . . .	56
3.4.3	Discrete problem . . . . .	57
<b>3.5</b>	<b>Analysis</b> . . . . .	<b>58</b>
3.5.1	Existence and uniqueness . . . . .	58
3.5.2	Convergence . . . . .	60
3.5.3	Error estimate . . . . .	65
<b>3.6</b>	<b>Local principle of virtual work and law of action and reaction</b> . . . . .	<b>68</b>
<b>3.7</b>	<b>Numerical results</b> . . . . .	<b>70</b>
3.7.1	Convergence for the Hencky–Mises model . . . . .	70
3.7.2	Tensile and shear test cases . . . . .	71
<b>3.8</b>	<b>Appendix: Technical results</b> . . . . .	<b>74</b>
3.8.1	Discrete Korn inequality . . . . .	74
3.8.2	Discrete compactness . . . . .	76
3.8.3	Consistency of the discrete symmetric gradient operator . . . . .	77

---

### 3.1 Introduction

In this chapter we develop and analyze a novel Hybrid High-Order (HHO) method for a class of (linear and) nonlinear elasticity problems in the small deformation regime.

Let  $\Omega \subset \mathbb{R}^d$ ,  $d \in \{2, 3\}$ , denote a bounded connected open polyhedral domain with Lipschitz boundary  $\Gamma := \partial\Omega$  and outward normal  $\mathbf{n}$ . We consider a body that occupies the region  $\Omega$  and is subjected to a volumetric force field  $\mathbf{f} \in L^2(\Omega; \mathbb{R}^d)$ . For the sake of simplicity, we assume the body fixed on  $\Gamma$  (extensions to other standard boundary conditions are possible). The nonlinear elasticity problem consists in finding a vector-valued displacement field  $\mathbf{u} : \Omega \rightarrow \mathbb{R}^d$  solution of

$$-\nabla \cdot \boldsymbol{\sigma}(\cdot, \nabla_s \mathbf{u}) = \mathbf{f} \quad \text{in } \Omega, \quad (3.1a)$$

$$\mathbf{u} = \mathbf{0} \quad \text{on } \Gamma, \quad (3.1b)$$

where  $\nabla_s$  denotes the symmetric gradient. The stress-strain law  $\boldsymbol{\sigma} : \Omega \times \mathbb{R}_{\text{sym}}^{d \times d} \rightarrow \mathbb{R}_{\text{sym}}^{d \times d}$  is assumed to satisfy regularity requirements closely inspired by [87], including conditions on its growth, coercivity, and monotonicity; cf. Assumption 3.1 below for a precise statement. Problem (3.1) is relevant, e.g., in modeling the mechanical behavior of soft materials [191] and metal alloys [169]. Examples of stress-strain laws of common use in the engineering practice are collected in Section 3.2.

The HHO discretization studied in this chapter is inspired by recent works on linear elasticity [74] (where HHO methods were originally introduced) and Leray–Lions operators [69, 70]. It hinges on degrees of freedom (DOFs) that are discontinuous polynomials of degree  $k \geq 1$  on the mesh and on the mesh skeleton. Based on these DOFs, we reconstruct discrete counterparts of the symmetric gradient and of the displacement by solving local linear problems inside each mesh element. These reconstruction operators are used to formulate a local contribution composed of two terms: a consistency term inspired by the weak formulation of problem (3.1) with  $\nabla_s$  replaced by its discrete counterpart, and a stabilization term penalizing cleverly designed face-based residuals. The resulting method has several advantageous features: (i) it is valid in arbitrary space dimension; (ii) it supports arbitrary polynomial orders  $\geq 1$  on fairly general meshes including, e.g., polyhedral elements and nonmatching interfaces; (iii) it satisfies inside each mesh element a local principle of virtual work with numerical tractions that obey the law of action and reaction; (iv) it can be efficiently implemented thanks to the possibility of statically condensing a large subset of the unknowns for linearized versions of the problem encountered, e.g., when solving the corresponding system of nonlinear algebraic equations by the Newton method (a numerical comparison between HHO methods and standard conforming finite element methods in the context of scalar diffusion problems can be found in [79]). Additionally, as shown by the numerical tests of Section 3.7, the method is robust with respect to strong nonlinearities.

In the context of structural mechanics, discretization methods supporting polyhedral meshes and nonconforming interfaces can be useful for several reasons including, e.g., the use of hanging nodes for contact [26, 202] and interface [106] elasticity problems, the simplicity in mesh refinement [186] and coarsening [16] for adaptivity, and the greater robustness to mesh distortion [55] and fracture [139]. The use of high-order methods, on the other hand,

can classically accelerate the convergence in the presence of regular exact solutions or when combined with local mesh refinement. Over the last few years, several discretization schemes supporting polyhedral meshes and/or high-order have been proposed for the linear version of problem (3.1); a non-exhaustive list includes [20, 74, 80, 81, 185, 197, 199]. For the nonlinear version, the literature is more scarce. Conforming approximations on standard meshes have been considered in [104, 105], where the convergence analysis is carried out assuming regularity for the exact displacement field  $\mathbf{u}$  and the constraint tensor  $\boldsymbol{\sigma}(\cdot, \nabla_s \mathbf{u})$  beyond the minimal regularity required by the weak formulation. Discontinuous Galerkin methods on standard meshes have been considered in [163], where convergence is proved for  $d = 2$  assuming  $\mathbf{u} \in H^{m+1}(\Omega; \mathbb{R}^2)$  for some  $m > 2$ , and in [25], where convergence to minimal regularity solutions is proved for stress-strain functions similar to [21]. General meshes are considered, on the other hand, in [21] and [54], where the authors propose a low-order Virtual Element method, whose convergence analysis is carried out for nonlinear elasticity problems in the small deformation regime (more general problems are considered numerically). In [21], an energy-norm convergence estimate in  $h$  (with  $h$  denoting, as usual, the meshsize) is proved when  $\mathbf{u} \in H^2(\Omega; \mathbb{R}^d)$  under the assumption that the function  $\tau \mapsto \boldsymbol{\sigma}(\cdot, \tau)$  is piecewise  $C^1$  with positive definite and bounded differential inside each element. A closer look at the proof reveals that properties essentially analogous to the ones considered in Assumption 3.14 below are in fact sufficient for the analysis, while  $C^1$ -regularity is used for the evaluation of the stability constant. Convergence to solutions that exhibit only the minimal regularity required by the weak formulation and for stress-strain functions as in Assumption 3.1 is proved in [87] for Gradient Schemes [85]. In this case, convergence rates are only proved for the linear case. We note, in passing, that the HHO method studied here fails to enter the Gradient Scheme framework essentially because the stabilization term is not embedded in the discrete symmetric gradient operator; see [72].

We carry out a complete analysis for the proposed HHO discretization of problem (3.1). Existence of a discrete solution is proved in Theorem 3.7, where we also identify a strict monotonicity assumption on the stress-strain law which ensures uniqueness. Convergence to minimal regularity solutions  $\mathbf{u} \in H_0^1(\Omega; \mathbb{R}^d)$  is proved in Theorem 3.9 using a compactness argument inspired by [69, 87]. More precisely, we prove for monotone stress-strain laws that (i) the discrete displacement field strongly converges (up to a subsequence) to  $\mathbf{u}$  in  $L^q(\Omega; \mathbb{R}^d)$  with  $1 \leq q < +\infty$  if  $d = 2$  and  $1 \leq q < 6$  if  $d = 3$ ; (ii) the discrete strain tensor weakly converges (up to a subsequence) to  $\nabla_s \mathbf{u}$  in  $L^2(\Omega, \mathbb{R}^{d \times d})$ . Notice that our results are slightly stronger than [87, Theorem 3.5] (cf. also Remark 3.6 therein) because the HHO discretization is compact as proved in Lemma 3.12. If, additionally, strict monotonicity holds for  $\boldsymbol{\sigma}$ , the strain tensor strongly converges and convergence extends to the whole sequence. An optimal energy-norm error estimate in  $h^{k+1}$  is then proved in Theorem 3.16 under the additional conditions of Lipschitz continuity and strong monotonicity on the stress-strain law; cf. Assumption 3.14. The performance of the method is investigated in Section 3.7 on a complete panel of model problems using stress-strain laws corresponding to real materials.

The rest of the chapter is organized as follows. In Section 3.2 we formulate the assumptions on the stress-strain function  $\boldsymbol{\sigma}$ , provide several examples of models relevant in the engineering practice, and write the weak formulation of problem (3.1). In Section 3.3 we introduce the notation for the mesh and recall a few known results. In Section 3.4 we discuss the choice of

DOFs, formulate the local reconstructions and the discrete problem. The main results, stated and proved in Section 3.5, are collected in Theorems 3.7, 3.9, and 3.16 below. In Section 3.6 we show that the HHO method satisfies on each mesh element a discrete counterpart of the principle of virtual work, and that interface tractions obey the law of action and reaction. Section 3.7 contains numerical tests. Finally, Appendix 3.8 contains the proofs of intermediate technical results.

## 3.2 Setting and examples

For the stress-strain function, we make the following

**Assumption 3.1** (Stress-strain function I). The stress-strain function  $\sigma : \Omega \times \mathbb{R}_{\text{sym}}^{d \times d} \rightarrow \mathbb{R}_{\text{sym}}^{d \times d}$  is a Caratheodory function, namely

$$\sigma(\mathbf{x}, \cdot) \text{ is continuous on } \mathbb{R}_{\text{sym}}^{d \times d} \text{ for a.e. } \mathbf{x} \in \Omega, \quad (3.2a)$$

$$\sigma(\cdot, \boldsymbol{\tau}) \text{ is measurable on } \Omega \text{ for all } \boldsymbol{\tau} \in \mathbb{R}_{\text{sym}}^{d \times d}, \quad (3.2b)$$

and it holds  $\sigma(\cdot, \mathbf{0}) \in L^2(\Omega; \mathbb{R}^{d \times d})$ . Moreover, there exist real numbers  $\bar{\sigma}, \underline{\sigma} \in (0, +\infty)$  such that, for a.e.  $\mathbf{x} \in \Omega$ , and all  $\boldsymbol{\tau}, \boldsymbol{\eta} \in \mathbb{R}_{\text{sym}}^{d \times d}$ , the following conditions hold:

$$\|\sigma(\mathbf{x}, \boldsymbol{\tau}) - \sigma(\mathbf{x}, \mathbf{0})\|_{d \times d} \leq \bar{\sigma} \|\boldsymbol{\tau}\|_{d \times d}, \quad (\text{growth}) \quad (3.2c)$$

$$\sigma(\mathbf{x}, \boldsymbol{\tau}) : \boldsymbol{\tau} \geq \underline{\sigma} \|\boldsymbol{\tau}\|_{d \times d}^2, \quad (\text{coercivity}) \quad (3.2d)$$

$$(\sigma(\mathbf{x}, \boldsymbol{\tau}) - \sigma(\mathbf{x}, \boldsymbol{\eta})) : (\boldsymbol{\tau} - \boldsymbol{\eta}) \geq 0, \quad (\text{monotonicity}) \quad (3.2e)$$

where  $\boldsymbol{\tau} : \boldsymbol{\eta} := \sum_{i,j=1}^d \tau_{i,j} \eta_{i,j}$  and  $\|\boldsymbol{\tau}\|_{d \times d}^2 := \boldsymbol{\tau} : \boldsymbol{\tau}$ .

We next discuss a number of meaningful models that satisfy the above assumptions.

**Example 3.2** (Linear elasticity). The linear elasticity model corresponds to

$$\sigma(\cdot, \nabla_s \mathbf{u}) = \mathbf{C}(\cdot) \nabla_s \mathbf{u},$$

where  $\mathbf{C}$  is a fourth order tensor. Being linear, the previous stress-strain relation clearly satisfies Assumption 3.1 provided that  $\mathbf{C}$  is uniformly elliptic. A particular case of the previous stress-strain relation is the usual linear elasticity Cauchy stress tensor

$$\sigma(\nabla_s \mathbf{u}) = \lambda \operatorname{tr}(\nabla_s \mathbf{u}) \mathbf{I}_d + 2\mu \nabla_s \mathbf{u}, \quad (3.3)$$

where  $\operatorname{tr}(\boldsymbol{\tau}) := \boldsymbol{\tau} : \mathbf{I}_d$  and  $\lambda \geq 0, \mu > 0$  are the Lamé's parameters.

**Example 3.3** (Hencky–Mises model). The nonlinear Hencky–Mises model of [105, 157] corresponds to the stress-strain relation

$$\sigma(\nabla_s \mathbf{u}) = \tilde{\lambda}(\operatorname{dev}(\nabla_s \mathbf{u})) \operatorname{tr}(\nabla_s \mathbf{u}) \mathbf{I}_d + 2\tilde{\mu}(\operatorname{dev}(\nabla_s \mathbf{u})) \nabla_s \mathbf{u}, \quad (3.4)$$

where  $\tilde{\lambda}$  and  $\tilde{\mu}$  are the nonlinear Lamé's scalar functions and  $\operatorname{dev} : \mathbb{R}_{\text{sym}}^{d \times d} \rightarrow \mathbb{R}$  defined by  $\operatorname{dev}(\boldsymbol{\tau}) = \operatorname{tr}(\boldsymbol{\tau}^2) - \frac{1}{d} \operatorname{tr}(\boldsymbol{\tau})^2$  is the deviatoric operator. Conditions on  $\tilde{\lambda}$  and  $\tilde{\mu}$  such that  $\sigma$  satisfies Assumption 3.1 can be found in [13, 21].

**Example 3.4** (An isotropic damage model). The isotropic damage model of [50] corresponds to the stress-strain relation

$$\boldsymbol{\sigma}(\cdot, \nabla_s \mathbf{u}) = (1 - D(\nabla_s \mathbf{u})) \mathbb{C}(\cdot) \nabla_s \mathbf{u}, \quad (3.5)$$

where  $D : \mathbb{R}_{\text{sym}}^{d \times d} \rightarrow \mathbb{R}$  is the scalar damage function and  $\mathbb{C}$  is a fourth-order symmetric and uniformly elliptic tensor. If there exists a continuous and bounded function  $f : [0, +\infty) \rightarrow [a, b]$  for some  $0 < a \leq b$ , such that  $s \in [0, +\infty) \rightarrow sf(s)$  is non-decreasing and, for all  $\boldsymbol{\tau} \in \mathbb{R}_{\text{sym}}^{d \times d}$ ,  $D(\boldsymbol{\tau}) = 1 - f(|\boldsymbol{\tau}|)$ , the damage model constitutive relation satisfies Assumption 3.1.

In the numerical experiments of Section 3.7 we will also consider the following model, relevant in engineering applications, which however does not satisfy Assumption 3.1 in general.

**Example 3.5** (The second-order elasticity model). The nonlinear second-order isotropic elasticity model of [66, 124, 135] corresponds to the stress-strain relation

$$\begin{aligned} \boldsymbol{\sigma}(\nabla_s \mathbf{u}) &= \lambda \operatorname{tr}(\nabla_s \mathbf{u}) \mathbf{I}_d + 2\mu \nabla_s \mathbf{u} \\ &+ B \operatorname{tr}((\nabla_s \mathbf{u})^2) \mathbf{I}_d + 2B \operatorname{tr}(\nabla_s \mathbf{u}) \nabla_s \mathbf{u} + C \operatorname{tr}(\nabla_s \mathbf{u})^2 \mathbf{I}_d + A(\nabla_s \mathbf{u})^2, \end{aligned} \quad (3.6)$$

where  $\lambda$  and  $\mu$  are the standard Lamé's parameter, and  $A, B, C \in \mathbb{R}$  are the second-order moduli.

*Remark 3.6* (Energy density functions). Examples 3.2, 3.3, and 3.5, used in numerical tests of Section 3.7, can be interpreted in the framework of hyperelasticity. Hyperelasticity is a type of constitutive model for ideally elastic materials in which the stress-strain relation derives from a stored energy density function  $\Psi : \mathbb{R}_{\text{sym}}^{d \times d} \rightarrow \mathbb{R}$ , namely

$$\boldsymbol{\sigma}(\boldsymbol{\tau}) := \frac{\partial \Psi(\boldsymbol{\tau})}{\partial \boldsymbol{\tau}}.$$

The stored energy density function leading to the linear Cauchy stress tensor (3.3) is

$$\Psi_{\text{lin}}(\boldsymbol{\tau}) := \frac{\lambda}{2} \operatorname{tr}(\boldsymbol{\tau})^2 + \mu \operatorname{tr}(\boldsymbol{\tau}^2), \quad (3.7)$$

while, in the Hencky–Mises model (3.4), it is defined such that

$$\Psi_{\text{hm}}(\boldsymbol{\tau}) := \frac{\alpha}{2} \operatorname{tr}(\boldsymbol{\tau})^2 + \Phi(\operatorname{dev}(\boldsymbol{\tau})). \quad (3.8)$$

Here  $\alpha \in (0, +\infty)$ , while  $\Phi : [0, +\infty) \rightarrow \mathbb{R}$  is a function of class  $C^2$  satisfying, for some positive constants  $C_1, C_2$ , and  $C_3$ ,

$$C_1 \leq \Phi'(\rho) < \alpha, \quad |\rho \Phi''(\rho)| \leq C_2, \quad \text{and} \quad \Phi'(\rho) + 2\rho \Phi''(\rho) \geq C_3 \quad \forall \rho \in [0, +\infty). \quad (3.9)$$

Deriving the energy density function (3.8) yields the stress-strain relation (3.4) with nonlinear Lamé's functions  $\tilde{\mu}(\rho) := \Phi'(\rho)$  and  $\tilde{\lambda}(\rho) := \alpha - \Phi'(\rho)$ . Taking  $\alpha = \lambda + \mu$  and  $\Phi(\rho) = \mu\rho$  in (3.8) leads to the linear case. Finally, the second-order elasticity model (3.6) is obtained by adding third-order terms to the linear stored energy density function defined in (3.7):

$$\Psi_{\text{snd}}(\boldsymbol{\tau}) := \frac{\lambda}{2} \operatorname{tr}(\boldsymbol{\tau})^2 + \mu \operatorname{tr}(\boldsymbol{\tau}^2) + \frac{C}{3} \operatorname{tr}(\boldsymbol{\tau})^3 + B \operatorname{tr}(\boldsymbol{\tau}) \operatorname{tr}(\boldsymbol{\tau}^2) + \frac{A}{3} \operatorname{tr}(\boldsymbol{\tau}^3). \quad (3.10)$$



The weak formulation of problem (3.1) that will serve as a starting point for the development and analysis of the HHO method reads

$$\text{Find } \mathbf{u} \in H_0^1(\Omega; \mathbb{R}^d) \text{ such that, for all } \mathbf{v} \in H_0^1(\Omega; \mathbb{R}^d), a(\mathbf{u}, \mathbf{v}) = \int_{\Omega} \mathbf{f} \cdot \mathbf{v}, \quad (3.11)$$

where  $H_0^1(\Omega; \mathbb{R}^d)$  is the zero-trace subspace of  $H^1(\Omega; \mathbb{R}^d)$  and the function  $a : H_0^1(\Omega; \mathbb{R}^d) \times H_0^1(\Omega; \mathbb{R}^d) \rightarrow \mathbb{R}$  is such that

$$a(\mathbf{v}, \mathbf{w}) := \int_{\Omega} \boldsymbol{\sigma}(\mathbf{x}, \nabla_s \mathbf{v}(\mathbf{x})) : \nabla_s \mathbf{w}(\mathbf{x}) \, d\mathbf{x}.$$

Throughout the rest of the chapter, to alleviate the notation, we omit the dependence on the space variable  $\mathbf{x}$  and the differential  $d\mathbf{x}$  from integrals.

### 3.3 Notation and basic results

We consider refined sequences of general polytopal meshes as in [85, Definition 7.2] matching the regularity requirements detailed in [82, Definition 3]. The main points are summarized hereafter. Denote by  $\mathcal{H} \subset \mathbb{R}_*^+$  a countable set of meshsizes having 0 as its unique accumulation point, and let  $(\mathcal{T}_h)_{h \in \mathcal{H}}$  be a refined mesh sequence where each  $\mathcal{T}_h$  is a finite collection of nonempty disjoint open polyhedral elements  $T$  with boundary  $\partial T$  such that  $\overline{\Omega} = \bigcup_{T \in \mathcal{T}_h} \overline{T}$  and  $h = \max_{T \in \mathcal{T}_h} h_T$  with  $h_T$  diameter of  $T$ .

For each  $h \in \mathcal{H}$ , let  $\mathcal{F}_h$  be a set of faces with disjoint interiors which partitions the mesh skeleton, i.e.,  $\bigcup_{F \in \mathcal{F}_h} \overline{F} = \bigcup_{T \in \mathcal{T}_h} \partial T$ . A face  $F$  is defined here as a hyperplanar closed connected subset of  $\overline{\Omega}$  with positive  $(d-1)$ -dimensional Hausdorff measure such that (i) either there exist distinct  $T_1, T_2 \in \mathcal{T}_h$  such that  $F \subset \partial T_1 \cap \partial T_2$  and  $F$  is called an interface or (ii) there exists  $T \in \mathcal{T}_h$  such that  $F \subset \partial T \cap \Gamma$  and  $F$  is called a boundary face. Interfaces are collected in the set  $\mathcal{F}_h^i$  and boundary faces in  $\mathcal{F}_h^b$ , so that  $\mathcal{F}_h := \mathcal{F}_h^i \cup \mathcal{F}_h^b$ . For all  $T \in \mathcal{T}_h$ ,  $\mathcal{F}_T := \{F \in \mathcal{F}_h : F \subset \partial T\}$  denotes the set of faces contained in  $\partial T$  and, for all  $F \in \mathcal{F}_T$ ,  $\mathbf{n}_{TF}$  is the unit normal to  $F$  pointing out of  $T$ .

Mesh regularity holds in the sense that, for all  $h \in \mathcal{H}$ ,  $\mathcal{T}_h$  admits a matching simplicial submesh  $\mathfrak{T}_h$  and there exists a real number  $\varrho > 0$  such that, for all  $h \in \mathcal{H}$ , (i) for any simplex  $S \in \mathfrak{T}_h$  of diameter  $h_S$  and inradius  $r_S$ ,  $\varrho h_S \leq r_S$  and (ii) for any  $T \in \mathcal{T}_h$  and all  $S \in \mathfrak{T}_h$  such that  $S \subset T$ ,  $\varrho h_T \leq h_S$ .

Let  $X$  be a mesh element or face. For an integer  $l \geq 0$ , we denote by  $\mathbb{P}_d^l(X; \mathbb{R})$  the space spanned by the restriction to  $X$  of scalar-valued,  $d$ -variate polynomials of total degree  $l$ . The  $L^2$ -projector  $\pi_X^l : L^1(X; \mathbb{R}) \rightarrow \mathbb{P}_d^l(X; \mathbb{R})$  is defined such that, for all  $v \in L^1(X; \mathbb{R})$ ,

$$\int_X (\pi_X^l v - v) w = 0 \quad \forall w \in \mathbb{P}_d^l(X; \mathbb{R}). \quad (3.12)$$

When dealing with the vector-valued polynomial space  $\mathbb{P}_d^l(X; \mathbb{R}^d)$  or with the tensor-valued polynomial space  $\mathbb{P}_d^l(X; \mathbb{R}^{d \times d})$ , we use the boldface notation  $\boldsymbol{\pi}_X^l$  for the corresponding  $L^2$ -orthogonal projectors acting component-wise.

On regular mesh sequences, we have the following optimal approximation properties for  $\pi_T^l$  (for a proof, cf. [76, Lemmas 1.58 and 1.59] and, in a more general framework, [70, Lemmas 3.4 and 3.6]): There exists a real number  $C_{\text{app}} > 0$  such that, for all  $s \in \{0, \dots, l+1\}$ , all  $h \in \mathcal{H}$ , all  $T \in \mathcal{T}_h$ , and all  $v \in H^s(T; \mathbb{R})$ ,

$$|v - \pi_T^l v|_{H^m(T; \mathbb{R})} \leq C_{\text{app}} h_T^{s-m} |v|_{H^s(T; \mathbb{R})} \quad \forall m \in \{0, \dots, s\}, \quad (3.13a)$$

and, if  $s \geq 1$ ,

$$|v - \pi_T^l v|_{H^m(\mathcal{F}_T; \mathbb{R})} \leq C_{\text{app}} h_T^{s-m-\frac{1}{2}} |v|_{H^s(T; \mathbb{R})} \quad \forall m \in \{0, \dots, s-1\}, \quad (3.13b)$$

with  $H^m(\mathcal{F}_T; \mathbb{R})$  denoting the broken hilbert space on the faces of  $T$ . Other useful geometric and functional analytic results on regular mesh sequences can be found in [76, Chapter 1] and [69, 70].

At the global level, we define broken versions of polynomial and Sobolev spaces. In particular, for an integer  $l \geq 0$ , we denote by  $\mathbb{P}_d^l(\mathcal{T}_h; \mathbb{R})$ ,  $\mathbb{P}_d^l(\mathcal{T}_h; \mathbb{R}^d)$ , and  $\mathbb{P}_d^l(\mathcal{T}_h; \mathbb{R}^{d \times d})$ , respectively, the space of scalar-valued, vector-valued, and tensor-valued broken polynomial functions on  $\mathcal{T}_h$  of total degree  $l$ . The space of broken vector-valued polynomial functions of total degree  $l$  on the trace of the mesh on the domain boundary  $\Gamma$  is denoted by  $\mathbb{P}_d^l(\mathcal{F}_h^{\text{b}}; \mathbb{R}^d)$ . Similarly, for an integer  $s \geq 1$ ,  $H^s(\mathcal{T}_h; \mathbb{R})$ ,  $H^s(\mathcal{T}_h; \mathbb{R}^d)$ , and  $H^s(\mathcal{T}_h; \mathbb{R}^{d \times d})$  are the scalar-valued, vector-valued, and tensor-valued broken Sobolev spaces of index  $s$ .

Throughout the rest of the chapter, for  $X \subset \bar{\Omega}$ , we denote by  $\|\cdot\|_X$  the standard norm in  $L^2(X; \mathbb{R})$ , with the convention that the subscript is omitted whenever  $X = \Omega$ . The same notation is used for the vector- and tensor-valued spaces  $L^2(X; \mathbb{R}^d)$  and  $L^2(X; \mathbb{R}^{d \times d})$ .

## 3.4 The Hybrid High-Order method

In this section we define the space of DOFs and the local reconstructions, and we state the discrete problem.

### 3.4.1 Degrees of freedom

Let a polynomial degree  $k \geq 1$  be fixed. The DOFs for the displacement are collected in the space

$$\underline{U}_h^k := \left( \bigotimes_{T \in \mathcal{T}_h} \mathbb{P}_d^k(T; \mathbb{R}^d) \right) \times \left( \bigotimes_{F \in \mathcal{F}_h} \mathbb{P}_d^k(F; \mathbb{R}^d) \right),$$

see Figure 3.1. Observe that naming  $\underline{U}_h^k$  the space of DOFs involves a shortcut: the actual DOFs can be chosen in several equivalent ways (polynomial moments, point values, etc.), and the specific choice does not affect the following discussion. Details concerning the actual choice made in the implementation are given in Section 3.7 below.

For a generic collection of DOFs in  $\underline{U}_h^k$ , we use the classical HHO underlined notation  $\underline{\mathbf{v}}_h := ((\mathbf{v}_T)_{T \in \mathcal{T}_h}, (\mathbf{v}_F)_{F \in \mathcal{F}_h})$ . We also denote by  $\mathbf{v}_h \in \mathbb{P}_d^k(\mathcal{T}_h; \mathbb{R}^d)$  and  $\mathbf{v}_{\Gamma, h} \in \mathbb{P}_d^k(\mathcal{F}_h^{\text{b}}; \mathbb{R}^d)$  (not underlined) the broken polynomial functions such that

$$(\mathbf{v}_h)|_T = \mathbf{v}_T \quad \forall T \in \mathcal{T}_h \quad \text{and} \quad (\mathbf{v}_{\Gamma, h})|_F = \mathbf{v}_F \quad \forall F \in \mathcal{F}_h^{\text{b}}.$$

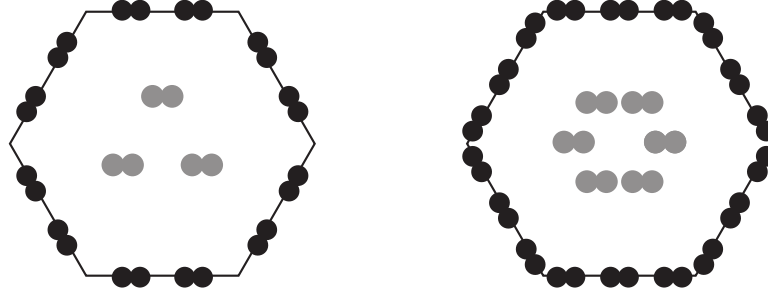


Figure 3.1: Local DOFs for  $k = 1$  (left) and  $k = 2$  (right). Shaded DOFs can be locally eliminated by static condensation when solving linearized versions of problem (3.21).

The restrictions of  $\underline{U}_h^k$  and  $\underline{v}_h$  to a mesh element  $T$  are denoted by  $\underline{U}_T^k$  and  $\underline{v}_T = (\mathbf{v}_T, (\mathbf{v}_F)_{F \in \mathcal{F}_T})$ , respectively. The space  $\underline{U}_h^k$  is equipped with the following discrete strain semi-norm:

$$\|\underline{v}_h\|_{\epsilon, h} := \left( \sum_{T \in \mathcal{T}_h} \|\underline{v}_T\|_{\epsilon, T}^2 \right)^{1/2}, \quad \|\underline{v}_h\|_{\epsilon, T}^2 := \|\nabla_s \mathbf{v}_T\|_T^2 + \sum_{F \in \mathcal{F}_T} h_F^{-1} \|\mathbf{v}_F - \mathbf{v}_T\|_F^2. \quad (3.14)$$

The DOFs corresponding to a given function  $\mathbf{v} \in H^1(\Omega; \mathbb{R}^d)$  are obtained by means of the reduction map  $\underline{I}_h^k : H^1(\Omega; \mathbb{R}^d) \rightarrow \underline{U}_h^k$  such that

$$\underline{I}_h^k \mathbf{v} := ((\pi_T^k \mathbf{v})_{T \in \mathcal{T}_h}, (\pi_F^k \mathbf{v})_{F \in \mathcal{F}_h}), \quad (3.15)$$

where we remind the reader that  $\pi_T^k$  and  $\pi_F^k$  denote the  $L^2$ -orthogonal projectors on  $\mathbb{P}_d^k(T; \mathbb{R}^d)$  and  $\mathbb{P}_d^k(F; \mathbb{R}^d)$ , respectively. For all mesh elements  $T \in \mathcal{T}_h$ , the local reduction map  $\underline{I}_T^k : H^1(T; \mathbb{R}^d) \rightarrow \underline{U}_T^k$  is obtained by a restriction of  $\underline{I}_h^k$ , and is therefore such that for all  $\mathbf{v} \in H^1(T; \mathbb{R}^d)$

$$\underline{I}_T^k \mathbf{v} = (\pi_T^k \mathbf{v}, (\pi_F^k \mathbf{v})_{F \in \mathcal{F}_T}). \quad (3.16)$$

### 3.4.2 Local reconstructions

We introduce symmetric gradient and displacement reconstruction operators devised at the element level that are instrumental to the formulation of the method.

Let a mesh element  $T \in \mathcal{T}_h$  be fixed. The local symmetric gradient reconstruction operator

$$\mathbf{G}_{s,T}^k : \underline{U}_T^k \rightarrow \mathbb{P}_d^k(T; \mathbb{R}_{\text{sym}}^{d \times d})$$

is obtained by solving the following pure traction problem: For a given local collection of DOFs  $\underline{v}_T = (\mathbf{v}_T, (\mathbf{v}_F)_{F \in \mathcal{F}_T}) \in \underline{U}_T^k$ , find  $\mathbf{G}_{s,T}^k \underline{v}_T \in \mathbb{P}_d^k(T; \mathbb{R}_{\text{sym}}^{d \times d})$  such that, for all  $\boldsymbol{\tau} \in \mathbb{P}_d^k(T; \mathbb{R}_{\text{sym}}^{d \times d})$ ,

$$\int_T \mathbf{G}_{s,T}^k \underline{v}_T : \boldsymbol{\tau} = - \int_T \mathbf{v}_T \cdot (\nabla \cdot \boldsymbol{\tau}) + \sum_{F \in \mathcal{F}_T} \int_F \mathbf{v}_F \cdot (\boldsymbol{\tau} \mathbf{n}_{TF}) \quad (3.17a)$$

$$= \int_T \nabla_s \mathbf{v}_T : \boldsymbol{\tau} + \sum_{F \in \mathcal{F}_T} \int_F (\mathbf{v}_F - \mathbf{v}_T) \cdot (\boldsymbol{\tau} \mathbf{n}_{TF}). \quad (3.17b)$$

The right-hand side of (3.17a) is designed to resemble an integration by parts formula where the role of the function represented by the DOFs in  $\underline{\mathbf{v}}_T$  is played by  $\mathbf{v}_T$  inside the volumetric integral and by  $(\mathbf{v}_F)_{F \in \mathcal{F}_T}$  inside boundary integrals. The reformulation (3.17b), obtained integrating by parts the first term in the right-hand side of (3.17a), highlights the fact that our method is nonconforming, as the second addend accounts for the difference between  $\mathbf{v}_F$  and  $\mathbf{v}_T$ .

The definition of the symmetric gradient reconstruction is justified observing that, using the definitions (3.16) of the local reduction map  $\underline{I}_T^k$  and (3.12) of the  $L^2$ -orthogonal projectors  $\pi_T^k$  and  $\pi_F^k$  in (3.17a), one can prove the following commuting property: For all  $T \in \mathcal{T}_h$  and all  $\mathbf{v} \in H^1(T; \mathbb{R}^d)$ ,

$$\mathbf{G}_{s,T}^k \underline{I}_T^k \mathbf{v} = \pi_T^k(\nabla_s \mathbf{v}). \quad (3.18)$$

As a result of (3.18) and (3.13),  $\mathbf{G}_{s,T}^k \underline{I}_T^k$  has optimal approximation properties in  $\mathbb{P}_d^k(T; \mathbb{R}_{\text{sym}}^{d \times d})$ .

From  $\mathbf{G}_{s,T}^k$ , one can define the local displacement reconstruction operator

$$\mathbf{r}_T^{k+1} : \underline{U}_T^k \rightarrow \mathbb{P}_d^{k+1}(T; \mathbb{R}^d)$$

such that, for all  $\underline{\mathbf{v}}_T \in \underline{U}_T^k$ ,  $\nabla_s \mathbf{r}_T^{k+1} \underline{\mathbf{v}}_T$  is the orthogonal projection of  $\mathbf{G}_{s,T}^k \underline{\mathbf{v}}_T$  on  $\nabla_s \mathbb{P}_d^{k+1}(T; \mathbb{R}^d) \subset \mathbb{P}_d^k(T; \mathbb{R}_{\text{sym}}^{d \times d})$  and rigid-body motions are prescribed according to [74, Eq. (15)]. More precisely, we let  $\mathbf{r}_T^{k+1} \underline{\mathbf{v}}_T$  be such that for all  $\mathbf{w} \in \mathbb{P}_d^{k+1}(T; \mathbb{R}^d)$  it holds

$$\int_T (\nabla_s \mathbf{r}_T^{k+1} \underline{\mathbf{v}}_T - \mathbf{G}_{s,T}^k \underline{\mathbf{v}}_T) : \nabla_s \mathbf{w} = 0$$

and, denoting by  $\nabla_{\text{ss}}$  the skew-symmetric part of the gradient operator, we have

$$\int_T \mathbf{r}_T^{k+1} \underline{\mathbf{v}}_T = \int_T \mathbf{v}_T, \quad \int_T \nabla_{\text{ss}} \mathbf{r}_T^{k+1} \underline{\mathbf{v}}_T = \sum_{F \in \mathcal{F}_T} \int_F \frac{1}{2} (\mathbf{v}_F \otimes \mathbf{n}_{TF} - \mathbf{n}_{TF} \otimes \mathbf{v}_F).$$

Notice that, for a given  $\underline{\mathbf{v}}_T \in \underline{U}_T^k$ , the displacement reconstruction  $\mathbf{r}_T^{k+1} \underline{\mathbf{v}}_T$  is a vector-valued polynomial function one degree higher than the element-based DOFs  $\mathbf{v}_T$ . It was proved in [74, Lemma 2] that  $\mathbf{r}_T^{k+1} \underline{I}_T^k$  has optimal approximation properties in  $\mathbb{P}_d^{k+1}(T; \mathbb{R}^d)$ .

In what follows, we will also need the global counterparts of the discrete gradient and displacement operators  $\mathbf{G}_{s,h}^k : \underline{U}_h^k \rightarrow \mathbb{P}_d^k(\mathcal{T}_h; \mathbb{R}_{\text{sym}}^{d \times d})$  and  $\mathbf{r}_h^{k+1} : \underline{U}_h^k \rightarrow \mathbb{P}_d^{k+1}(\mathcal{T}_h; \mathbb{R}^d)$  defined setting, for all  $\underline{\mathbf{v}}_h \in \underline{U}_h^k$  and all  $T \in \mathcal{T}_h$ ,

$$(\mathbf{G}_{s,h}^k \underline{\mathbf{v}}_h)|_T = \mathbf{G}_{s,T}^k \underline{\mathbf{v}}_T, \quad (\mathbf{r}_h^{k+1} \underline{\mathbf{v}}_h)|_T = \mathbf{r}_T^{k+1} \underline{\mathbf{v}}_T. \quad (3.19)$$

### 3.4.3 Discrete problem

We define the following subspace of  $\underline{U}_h^k$  strongly accounting for the homogeneous Dirichlet boundary condition (3.1b):

$$\underline{U}_{h,D}^k := \left\{ \underline{\mathbf{v}}_h \in \underline{U}_h^k : \mathbf{v}_F = \mathbf{0} \quad \forall F \in \mathcal{F}_h^b \right\}, \quad (3.20)$$

and we notice that the map  $\|\cdot\|_{\epsilon,h}$  defined by (3.14) is a norm on  $\underline{U}_{h,D}^k$ . The HHO approximation of problem (3.11) reads:

$$\text{Find } \underline{\mathbf{u}}_h \in \underline{U}_{h,D}^k \text{ such that, for all } \underline{\mathbf{v}}_h \in \underline{U}_{h,D}^k, \\ a_h(\underline{\mathbf{u}}_h, \underline{\mathbf{v}}_h) := A_h(\underline{\mathbf{u}}_h, \underline{\mathbf{v}}_h) + s_h(\underline{\mathbf{u}}_h, \underline{\mathbf{v}}_h) = \int_{\Omega} \mathbf{f} \cdot \mathbf{v}_h, \quad (3.21)$$

where the consistency contribution  $A_h : \underline{U}_h^k \times \underline{U}_h^k \rightarrow \mathbb{R}$  and the stability contribution  $s_h : \underline{U}_h^k \times \underline{U}_h^k \rightarrow \mathbb{R}$  are respectively defined setting

$$A_h(\underline{\mathbf{u}}_h, \underline{\mathbf{v}}_h) := \int_{\Omega} \boldsymbol{\sigma}(\cdot, \mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h) : \mathbf{G}_{s,h}^k \underline{\mathbf{v}}_h, \\ s_h(\underline{\mathbf{u}}_h, \underline{\mathbf{v}}_h) := \sum_{T \in \mathcal{T}_h} s_T(\underline{\mathbf{u}}_T, \underline{\mathbf{v}}_T), \quad s_T(\underline{\mathbf{u}}_T, \underline{\mathbf{v}}_T) := \sum_{F \in \mathcal{F}_T} \frac{\gamma}{h_F} \int_F \Delta_{TF}^k \underline{\mathbf{u}}_T \cdot \Delta_{TF}^k \underline{\mathbf{v}}_T. \quad (3.22)$$

The scaling parameter  $\gamma > 0$  in (3.22) can depend on  $\bar{\sigma}$  and  $\underline{\sigma}$  but is independent of the meshsize  $h$ . In the numerical tests of Section 3.7 we take  $\gamma = 2\mu$  for the linear (3.3) and second-order (3.6) models and  $\gamma = 2\tilde{\mu}(\mathbf{0})$  for the Hencky–Mises model (3.4). In  $s_T$ , we penalize in a least-square sense the face-based residual  $\Delta_{TF}^k : \underline{U}_T^k \rightarrow \mathbb{P}_d^k(F; \mathbb{R}^d)$  such that, for all  $T \in \mathcal{T}_h$ , all  $\underline{\mathbf{v}}_T \in \underline{U}_T^k$ , and all  $F \in \mathcal{F}_T$ ,

$$\Delta_{TF}^k \underline{\mathbf{v}}_T := \boldsymbol{\pi}_F^k(\mathbf{r}_T^{k+1} \underline{\mathbf{v}}_T - \mathbf{v}_F) - \boldsymbol{\pi}_T^k(\mathbf{r}_T^{k+1} \underline{\mathbf{v}}_T - \mathbf{v}_T). \quad (3.23)$$

This particular choice ensures that  $\Delta_{TF}^k$  vanishes whenever its argument is of the form  $\mathbf{I}_T^k \mathbf{w}$  with  $\mathbf{w} \in \mathbb{P}_d^{k+1}(T; \mathbb{R}^d)$ , a crucial property to obtain an energy-norm error estimate in  $h^{k+1}$ ; cf. Theorem 3.16 below. Additionally,  $s_h$  is stabilizing in the sense that the following uniform norm equivalence holds (the proof is a straightforward modification of [74, Lemma 4]; cf. also Corollary 6 therein): There exists a real number  $\eta > 0$  independent of  $h$  such that, for all  $\underline{\mathbf{v}}_h \in \underline{U}_{h,D}^k$ ,

$$\eta^{-1} \|\underline{\mathbf{v}}_h\|_{\epsilon,h}^2 \leq \|\mathbf{G}_{s,h}^k \underline{\mathbf{v}}_h\|^2 + s_h(\underline{\mathbf{v}}_h, \underline{\mathbf{v}}_h) \leq \eta \|\underline{\mathbf{v}}_h\|_{\epsilon,h}^2. \quad (3.24)$$

By (3.2d), this implies the coercivity of  $a_h$ .

## 3.5 Analysis

In this section we carry out a complete analysis of the method. To alleviate the notation, inside the proofs we abridge into  $a \lesssim b$  the inequality  $a \leq Cb$  with real number  $C > 0$  independent of  $h$ .

### 3.5.1 Existence and uniqueness

We start by discussing existence and uniqueness of the discrete solution.

**Theorem 3.7** (Existence and uniqueness of a discrete solution). *Let Assumption 3.1 hold and let  $(\mathcal{T}_h)_{h \in \mathcal{H}}$  be a regular mesh sequence. Then, for all  $h \in \mathcal{H}$ , there exists at least one solution  $\underline{\mathbf{u}}_h \in \underline{\mathbf{U}}_{h,D}^k$  to problem (3.21). Additionally, if the stress-strain function  $\sigma$  is strictly monotone (i.e., if the inequality in (3.2e) is strict for  $\tau \neq \eta$ ), the solution is unique.*

*Proof.* 1) *Existence.* We follow the argument of [64, Theorem 3.3]. If  $(E, (\cdot, \cdot)_E, \|\cdot\|_E)$  is a Euclidean space and  $\Phi : E \rightarrow E$  is a continuous map such that  $\frac{(\Phi(x), x)_E}{\|x\|_E} \rightarrow +\infty$ , as  $\|x\|_E \rightarrow +\infty$ , then  $\Phi$  is surjective. We take  $E = \underline{\mathbf{U}}_{h,D}^k$ , endowed with the inner product

$$(\underline{\mathbf{v}}_h, \underline{\mathbf{w}}_h)_{\epsilon, h} := \sum_{T \in \mathcal{T}_h} \left( \int_T \nabla_s \mathbf{v}_T : \nabla_s \mathbf{w}_T + \sum_{F \in \mathcal{F}_T} \frac{1}{h_F} \int_F (\mathbf{v}_F - \mathbf{v}_T) \cdot (\mathbf{w}_F - \mathbf{w}_T) \right),$$

and we define  $\Phi : \underline{\mathbf{U}}_{h,D}^k \rightarrow \underline{\mathbf{U}}_{h,D}^k$  such that, for all  $\underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_{h,D}^k$ ,  $(\Phi(\underline{\mathbf{v}}_h), \underline{\mathbf{w}}_h)_{\epsilon, h} = a_h(\underline{\mathbf{v}}_h, \underline{\mathbf{w}}_h)$  for all  $\underline{\mathbf{w}}_h \in \underline{\mathbf{U}}_{h,D}^k$ . The coercivity (3.2d) of  $\sigma$  together with the norm equivalence (3.24) yields  $(\Phi(\underline{\mathbf{v}}_h), \underline{\mathbf{v}}_h)_{\epsilon, h} \geq \min\{1, \underline{\sigma}\} \eta^{-1} \|\underline{\mathbf{v}}_h\|_{\epsilon, h}^2$  for all  $\underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_{h,D}^k$ , so that  $\Phi$  is surjective. Let now  $\underline{\mathbf{y}}_h \in \underline{\mathbf{U}}_{h,D}^k$  be such that  $(\underline{\mathbf{y}}_h, \underline{\mathbf{w}}_h)_{\epsilon, h} = \int_{\Omega} \mathbf{f} \cdot \mathbf{w}_h$  for all  $\underline{\mathbf{w}}_h \in \underline{\mathbf{U}}_{h,D}^k$ . By the surjectivity of  $\Phi$ , there exists  $\underline{\mathbf{u}}_h \in \underline{\mathbf{U}}_{h,D}^k$  such that  $\Phi(\underline{\mathbf{u}}_h) = \underline{\mathbf{y}}_h$ . By definition of  $\Phi$  and  $\underline{\mathbf{y}}_h$ ,  $\underline{\mathbf{u}}_h$  is a solution to the problem (3.21).

2) *Uniqueness.* Let  $\underline{\mathbf{u}}_{h,1}, \underline{\mathbf{u}}_{h,2} \in \underline{\mathbf{U}}_{h,D}^k$  solve (3.21). We assume  $\underline{\mathbf{u}}_{h,1} \neq \underline{\mathbf{u}}_{h,2}$  and proceed by contradiction. Subtracting (3.21) for  $\underline{\mathbf{u}}_{h,2}$  from (3.21) for  $\underline{\mathbf{u}}_{h,1}$ , it is inferred that  $a_h(\underline{\mathbf{u}}_{h,1}, \underline{\mathbf{v}}_h) - a_h(\underline{\mathbf{u}}_{h,2}, \underline{\mathbf{v}}_h) = 0$  for all  $\underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_{h,D}^k$ . Hence in particular, taking  $\underline{\mathbf{v}}_h = \underline{\mathbf{u}}_{h,1} - \underline{\mathbf{u}}_{h,2}$  we obtain that

$$a_h(\underline{\mathbf{u}}_{h,1}, \underline{\mathbf{u}}_{h,1} - \underline{\mathbf{u}}_{h,2}) - a_h(\underline{\mathbf{u}}_{h,2}, \underline{\mathbf{u}}_{h,1} - \underline{\mathbf{u}}_{h,2}) = 0.$$

On the other hand, owing to the strict monotonicity of  $\sigma$  and to the fact that the bilinear form  $s_h$  is positive semidefinite, we have that

$$\begin{aligned} & a_h(\underline{\mathbf{u}}_{h,1}, \underline{\mathbf{u}}_{h,1} - \underline{\mathbf{u}}_{h,2}) - a_h(\underline{\mathbf{u}}_{h,2}, \underline{\mathbf{u}}_{h,1} - \underline{\mathbf{u}}_{h,2}) \\ &= \int_{\Omega} (\sigma(\cdot, \mathbf{G}_{s,h}^k \underline{\mathbf{u}}_{h,1}) - \sigma(\cdot, \mathbf{G}_{s,h}^k \underline{\mathbf{u}}_{h,2})) : \mathbf{G}_{s,h}^k (\underline{\mathbf{u}}_{h,1} - \underline{\mathbf{u}}_{h,2}) \\ & \quad + s_h(\underline{\mathbf{u}}_{h,1} - \underline{\mathbf{u}}_{h,2}, \underline{\mathbf{u}}_{h,1} - \underline{\mathbf{u}}_{h,2}) > 0. \end{aligned}$$

Hence,  $\underline{\mathbf{u}}_{h,1} = \underline{\mathbf{u}}_{h,2}$  and the conclusion follows.  $\square$

*Remark 3.8* (Strict monotonicity of the stress-strain function). The strict monotonicity assumption is fulfilled, e.g., by the Hencky–Mises model (3.4) and by the damage model (3.5) when  $D(\tau) = 1 - f(|\tau|)$ , with  $f$  continuous, bounded, and such that  $[0, +\infty) \ni s \mapsto sf(s)$  is strictly increasing. We observe, in passing, that the strict monotonicity is weaker than the strong monotonicity (3.40b) used in Theorem 3.16 below to prove error estimates.

### 3.5.2 Convergence

We now consider the convergence to solutions that only exhibit the minimal regularity required by the variational formulation (3.11). More specifically, in this section we prove the following **Theorem 3.9** (Convergence). *Let Assumption 3.1 hold, let  $k \geq 1$ , and let  $(\mathcal{T}_h)_{h \in \mathcal{H}}$  be a regular mesh sequence. Further assume the existence of a real number  $C_K > 0$  independent of  $h$  but possibly depending on  $\Omega$ ,  $\varrho$ , and on  $k$  such that, for all  $\underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_{h,D}^k$ ,*

$$\|\mathbf{v}_h\| + \|\nabla_h \mathbf{v}_h\| \leq C_K \|\underline{\mathbf{v}}_h\|_{\epsilon,h}, \quad (3.25)$$

where  $\nabla_h$  denotes the broken gradient on  $H^1(\mathcal{T}_h; \mathbb{R}^d)$ . For all  $h \in \mathcal{H}$ , let  $\underline{\mathbf{u}}_h \in \underline{\mathbf{U}}_{h,D}^k$  be a solution to the discrete problem (3.21) on  $\mathcal{T}_h$ . Then, for all  $q$  such that  $1 \leq q < +\infty$  if  $d = 2$  or  $1 \leq q < 6$  if  $d = 3$ , as  $h \rightarrow 0$  it holds, up to a subsequence,

- $\mathbf{u}_h \rightarrow \mathbf{u}$  strongly in  $L^q(\Omega; \mathbb{R}^d)$ ,
- $\mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h \rightarrow \nabla_s \mathbf{u}$  weakly in  $L^2(\Omega; \mathbb{R}^{d \times d})$ ,

where  $\mathbf{u} \in H_0^1(\Omega; \mathbb{R}^d)$  solves the weak formulation (3.11). Moreover, if we assume strict monotonicity for  $\sigma$  (i.e., the inequality in (3.2e) is strict for  $\tau \neq \eta$ ), it holds that

- $\mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h \rightarrow \nabla_s \mathbf{u}$  strongly in  $L^2(\Omega; \mathbb{R}^{d \times d})$ .

Finally, if the solution to (3.11) is unique, convergence extends to the whole sequence.

Some remarks are of order before proceeding with the proof of Theorem 3.9.

*Remark 3.10* (Existence of a solution to the continuous problem). A side result of the existence of discrete solutions proved in Theorem 3.7 together with the convergence results of Theorem 3.9 is the existence of a solution to the weak formulation (3.11).

*Remark 3.11* (Discrete Korn inequality). In Proposition 3.20 below we give a proof of the discrete Korn inequality (3.25) based on the results of [39], which require further assumptions on the mesh. While we have the feeling that these assumptions could probably be relaxed, we postpone this topic to a future work. Notice that inequality (3.25) is not required to prove the error estimate of Theorem 3.16 below.

The proof of Theorem 3.9 hinges on two technical results stated hereafter: a discrete Rellich–Kondrachov Lemma (cf. [40, Theorem 9.16]) and a proposition showing the approximation properties of the discrete symmetric gradient  $\mathbf{G}_{s,h}^k$  defined by (3.19). Their proofs are given in Appendix 3.8.

**Lemma 3.12** (Discrete compactness). *Let the assumptions of Theorem 3.9 hold. Let  $(\underline{\mathbf{v}}_h)_{h \in \mathcal{H}} \in (\underline{\mathbf{U}}_{h,D}^k)_{h \in \mathcal{H}}$ , and assume that there is a real number  $C \geq 0$  such that*

$$\|\underline{\mathbf{v}}_h\|_{\epsilon,h} \leq C \quad \forall h \in \mathcal{H}. \quad (3.26)$$

Then, for all  $q$  such that  $1 \leq q < +\infty$  if  $d = 2$  or  $1 \leq q < 6$  if  $d = 3$ , the sequence  $(\mathbf{v}_h)_{h \in \mathcal{H}} \in (\mathbb{P}_d^k(\mathcal{T}_h; \mathbb{R}^d))_{h \in \mathcal{H}}$  is relatively compact in  $L^q(\Omega; \mathbb{R}^d)$ . As a consequence, there is a function  $\mathbf{v} \in L^q(\Omega; \mathbb{R}^d)$  such that as  $h \rightarrow 0$ , up to a subsequence,  $\mathbf{v}_h \rightarrow \mathbf{v}$  strongly in  $L^q(\Omega; \mathbb{R}^d)$ .

*Proof.* See Appendix 3.8.2.  $\square$

**Proposition 3.13** (Consistency of the discrete symmetric gradient operator). *Let  $(\mathcal{T}_h)_{h \in \mathcal{H}}$  be a regular mesh sequence, and let  $\mathbf{G}_{s,h}^k$  be as in (3.19) with  $\mathbf{G}_{s,T}^k$  defined by (3.17) for all  $T \in \mathcal{T}_h$ .*

1) Strong consistency. *For all  $\mathbf{v} \in H^1(\Omega; \mathbb{R}^d)$  with  $\underline{I}_h^k$  defined by (3.15), it holds as  $h \rightarrow 0$*

$$\mathbf{G}_{s,h}^k \underline{I}_h^k \mathbf{v} \rightarrow \nabla_s \mathbf{v} \text{ strongly in } L^2(\Omega; \mathbb{R}^{d \times d}). \quad (3.27)$$

2) Sequential consistency. *For all  $h \in \mathcal{H}$  and all  $\boldsymbol{\tau} \in H^1(\Omega; \mathbb{R}_{\text{sym}}^{d \times d})$ , denoting by  $\boldsymbol{\gamma}_n(\boldsymbol{\tau})$  the normal trace of  $\boldsymbol{\tau}$  on  $\Gamma$ , it holds*

$$\lim_{h \rightarrow 0} \left( \max_{\underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_h^k, \|\underline{\mathbf{v}}_h\|_{\epsilon,h}=1} \left| \int_{\Omega} \mathbf{G}_{s,h}^k \underline{\mathbf{v}}_h : \boldsymbol{\tau} + \underline{\mathbf{v}}_h \cdot (\nabla \cdot \boldsymbol{\tau}) - \int_{\Gamma} \underline{\mathbf{v}}_{\Gamma,h} \cdot \boldsymbol{\gamma}_n(\boldsymbol{\tau}) \right| \right) = 0, \quad (3.28)$$

*Proof.* See Appendix 3.8.3.  $\square$

We are now ready to prove convergence.

*Proof of Theorem 3.9.* The proof is subdivided into four steps: in **Step 1** we prove a uniform a priori bound on the solutions of the discrete problem (3.21); in **Step 2** we infer the existence of a limit for the sequence of discrete solutions and investigate its regularity; in **Step 3** we show that this limit solves the continuous problem (3.11); finally, in **Step 4** we prove strong convergence.

**Step 1: A priori bound.** We start by showing the following uniform a priori bound on the sequence of discrete solutions:

$$\|\underline{\mathbf{u}}_h\|_{\epsilon,h} \leq C \|\mathbf{f}\|, \quad (3.29)$$

where the real number  $C > 0$  only depends on  $\Omega, \underline{\sigma}, \gamma, \varrho$ , and  $k$ . Making  $\underline{\mathbf{v}}_h = \underline{\mathbf{u}}_h$  in (3.21) and using the coercivity property (3.2d) of  $\boldsymbol{\sigma}$  in the left-hand side together with the Cauchy–Schwarz inequality in the right-hand side yields

$$\sum_{T \in \mathcal{T}_h} \left( \underline{\sigma} \|\mathbf{G}_{s,T}^k \underline{\mathbf{u}}_T\|_T^2 + \sum_{F \in \mathcal{F}_h} \frac{\gamma}{h_F} \|\Delta_{TF}^k \underline{\mathbf{u}}_T\|_F^2 \right) \leq \|\mathbf{f}\| \|\underline{\mathbf{u}}_h\|.$$

Owing to the norm equivalence (3.24), and using the discrete Korn inequality (3.25) to estimate the right-hand side of the previous inequality, it is inferred that

$$\eta^{-1} \min(1, \underline{\sigma}) \|\underline{\mathbf{u}}_h\|_{\epsilon,h}^2 \leq \|\mathbf{f}\| \|\underline{\mathbf{u}}_h\| \leq C_K \|\mathbf{f}\| \|\underline{\mathbf{u}}_h\|_{\epsilon,h}.$$

Dividing by  $\|\underline{\mathbf{u}}_h\|_{\epsilon,h}$  yields (3.29) with  $C = \eta \min(1, \underline{\sigma})^{-1} C_K$ .

**Step 2: Existence of a limit and regularity.** Let  $1 \leq q < +\infty$  if  $d = 2$  or  $1 \leq q < 6$  if  $d = 3$ . Owing to the a priori bound (3.29) and the norm equivalence (3.24), the sequences  $(\|\underline{\mathbf{u}}_h\|_{\epsilon,h})_{h \in \mathcal{H}}$  and  $(\|\mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h\|)_{h \in \mathcal{H}}$  are uniformly bounded. Therefore, Lemma 3.12



and the Kakutani theorem [40, Theorem 3.17] yield the existence of  $\mathbf{u} \in L^q(\Omega; \mathbb{R}^d)$  and  $\mathcal{G} \in L^2(\Omega; \mathbb{R}^{d \times d})$  such that as  $h \rightarrow 0$ , up to a subsequence,

$$\mathbf{u}_h \rightarrow \mathbf{u} \text{ strongly in } L^q(\Omega; \mathbb{R}^d) \text{ and } \mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h \rightarrow \mathcal{G} \text{ weakly in } L^2(\Omega; \mathbb{R}^{d \times d}).$$

This together with the fact that  $\mathbf{u}_{h,\Gamma} = \mathbf{0}$  on  $\Gamma$ , shows that, for any  $\boldsymbol{\tau} \in H^1(\Omega; \mathbb{R}^{d \times d}_{\text{sym}})$ ,

$$\begin{aligned} & \left| \int_{\Omega} \mathcal{G} : \boldsymbol{\tau} + \mathbf{u} \cdot (\nabla \cdot \boldsymbol{\tau}) \right| \\ &= \lim_{h \rightarrow 0} \left| \int_{\Omega} \mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h : \boldsymbol{\tau} + \mathbf{u}_h \cdot (\nabla \cdot \boldsymbol{\tau}) - \int_{\Gamma} \mathbf{u}_{h,\Gamma} \cdot \boldsymbol{\gamma}_n(\boldsymbol{\tau}) \right| \\ &\leq \lim_{h \rightarrow 0} \left( \|\underline{\mathbf{u}}_h\|_{\epsilon,h} \max_{\mathbf{v}_h \in \underline{\mathbf{U}}_h^k, \|\mathbf{v}_h\|_{\epsilon,h}=1} \left| \int_{\Omega} \mathbf{G}_{s,h}^k \underline{\mathbf{v}}_h : \boldsymbol{\tau} + \mathbf{v}_h \cdot (\nabla \cdot \boldsymbol{\tau}) - \int_{\Gamma} \mathbf{v}_{h,\Gamma} \cdot \boldsymbol{\gamma}_n(\boldsymbol{\tau}) \right| \right) \\ &= 0. \end{aligned} \tag{3.30}$$

To infer the previous limit we have used the uniform bound (3.29) on  $\|\underline{\mathbf{u}}_h\|_{\epsilon,h}$  and the sequential consistency (3.28) of  $\mathbf{G}_{s,h}^k$ . Applying (3.30) with  $\boldsymbol{\tau} \in C_c^\infty(\Omega; \mathbb{R}^{d \times d}_{\text{sym}})$  leads to

$$\int_{\Omega} \mathcal{G} : \boldsymbol{\tau} + \mathbf{u} \cdot (\nabla \cdot \boldsymbol{\tau}) = 0,$$

thus  $\mathcal{G} = \nabla_s \mathbf{u}$  in the sense of distributions on  $\Omega$ . As a result, owing to the isomorphism of Hilbert spaces between  $H^1(\Omega; \mathbb{R}^d)$  and  $\{\mathbf{v} \in L^2(\Omega; \mathbb{R}^d) \mid \nabla_s \mathbf{v} \in L^2(\Omega; \mathbb{R}^{d \times d}_{\text{sym}})\}$  proved in [89, Theorem 3.1], we have  $\mathbf{u} \in H^1(\Omega; \mathbb{R}^d)$ . Using again (3.30) with  $\boldsymbol{\tau} \in H^1(\Omega; \mathbb{R}^{d \times d}_{\text{sym}})$  and integrating by parts, we obtain

$$\int_{\Gamma} \boldsymbol{\gamma}(\mathbf{u}) \cdot \boldsymbol{\gamma}_n(\boldsymbol{\tau}) = 0$$

with  $\boldsymbol{\gamma}(\mathbf{u})$  denoting the trace of  $\mathbf{u}$ . As the set  $\{\boldsymbol{\gamma}_n(\boldsymbol{\tau}) : \boldsymbol{\tau} \in H^1(\Omega; \mathbb{R}^{d \times d}_{\text{sym}})\}$  is dense in  $L^2(\Gamma; \mathbb{R}^d)$ , we deduce that  $\boldsymbol{\gamma}(\mathbf{u}) = \mathbf{0}$  on  $\Gamma$ . In conclusion, with convergences up to a subsequence,

$$\begin{aligned} \mathbf{u} &\in H_0^1(\Omega; \mathbb{R}^d), \mathbf{u}_h \rightarrow \mathbf{u} \text{ strongly in } L^q(\Omega; \mathbb{R}^d) \\ \mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h &\rightarrow \nabla_s \mathbf{u} \text{ weakly in } L^2(\Omega; \mathbb{R}^{d \times d}). \end{aligned}$$

**Step 3: Identification of the limit.** Let us now prove that  $\mathbf{u}$  is a solution to (3.11). The growth property (3.2c) on  $\boldsymbol{\sigma}$  and the bound on  $(\|\mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h\|)_{h \in \mathcal{H}}$  ensure that the sequence  $(\boldsymbol{\sigma}(\cdot, \mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h))_{h \in \mathcal{H}}$  is bounded in  $L^2(\Omega; \mathbb{R}^{d \times d}_{\text{sym}})$ . Hence, there exists  $\boldsymbol{\eta} \in L^2(\Omega; \mathbb{R}^{d \times d}_{\text{sym}})$  such that, up to a subsequence as  $h \rightarrow 0$ ,

$$\boldsymbol{\sigma}(\cdot, \mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h) \rightarrow \boldsymbol{\eta} \quad \text{weakly in } L^2(\Omega; \mathbb{R}^{d \times d}). \tag{3.31}$$

Plugging into (3.21)  $\underline{\mathbf{v}}_h = \underline{I}_h^k \boldsymbol{\phi}$ , with  $\boldsymbol{\phi} \in C_c^\infty(\Omega; \mathbb{R}^d)$ , gives

$$\int_{\Omega} \boldsymbol{\sigma}(\cdot, \mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h) : \mathbf{G}_{s,h}^k \underline{I}_h^k \boldsymbol{\phi} = \int_{\Omega} \mathbf{f} \cdot \boldsymbol{\pi}_h^k \boldsymbol{\phi} - s_h(\underline{\mathbf{u}}_h, \underline{I}_h^k \boldsymbol{\phi}), \tag{3.32}$$

with  $\pi_h^k$  denoting the  $L^2$ -projector on the broken polynomial spaces  $\mathbb{P}_d^k(\mathcal{T}_h; \mathbb{R}^d)$  and  $s_h$  defined by (3.22). Using the Cauchy–Schwarz inequality followed by the norm equivalence (3.24) to bound the first factor, we infer

$$|s_h(\underline{\mathbf{u}}_h, \underline{I}_h^k \boldsymbol{\phi})| \leq s_h(\underline{\mathbf{u}}_h, \underline{\mathbf{u}}_h)^{1/2} s_h(\underline{I}_h^k \boldsymbol{\phi}, \underline{I}_h^k \boldsymbol{\phi})^{1/2} \leq \|\underline{\mathbf{u}}_h\|_{\epsilon, h} s_h(\underline{I}_h^k \boldsymbol{\phi}, \underline{I}_h^k \boldsymbol{\phi})^{1/2}. \quad (3.33)$$

It was proved in [74, Eq. (35)] using the optimal approximation properties of  $\mathbf{r}_T^{k+1} \underline{I}_T^k$  that it holds for all  $h \in \mathcal{H}$ , all  $T \in \mathcal{T}_h$ , all  $\mathbf{v} \in H^{k+2}(T; \mathbb{R}^d)$ , and all  $F \in \mathcal{F}_T$  that

$$h_F^{-1/2} \|\Delta_{TF}^k \underline{I}_T^k \mathbf{v}\|_F \lesssim h_T^{k+1} \|\mathbf{v}\|_{H^{k+2}(T; \mathbb{R}^d)}, \quad (3.34)$$

with  $\Delta_{TF}^k$  defined by (3.23). As a consequence, recalling the definition (3.22) of  $s_h$ , we have the following convergence result:

$$\forall \mathbf{v} \in H^1(\Omega; \mathbb{R}^d) \cap H^2(\mathcal{T}_h; \mathbb{R}^d), \quad \lim_{h \rightarrow 0} s_h(\underline{I}_h^k \mathbf{v}, \underline{I}_h^k \mathbf{v}) = 0. \quad (3.35)$$

Recalling the a priori bound (3.29) on the discrete solution and the convergence property (3.35), it follows from (3.33) that  $|s_h(\underline{\mathbf{u}}_h, \underline{I}_h^k \boldsymbol{\phi})| \rightarrow 0$  as  $h \rightarrow 0$ . Additionally, by the approximation property (3.13a) of the  $L^2$ -projector, one has  $\pi_h^k \boldsymbol{\phi} \rightarrow \boldsymbol{\phi}$  strongly in  $L^2(\Omega; \mathbb{R}^d)$  and, by virtue of Proposition 3.13, that  $\mathbf{G}_{s,h}^k \underline{I}_h^k \boldsymbol{\phi} \rightarrow \nabla_s \boldsymbol{\phi}$  strongly in  $L^2(\Omega; \mathbb{R}^{d \times d})$ . Thus, we can pass to the limit  $h \rightarrow 0$  in (3.32) and obtain

$$\int_{\Omega} \boldsymbol{\eta} : \nabla_s \boldsymbol{\phi} = \int_{\Omega} \mathbf{f} \cdot \boldsymbol{\phi}. \quad (3.36)$$

By density of  $C_c^\infty(\Omega; \mathbb{R}^d)$  in  $H_0^1(\Omega; \mathbb{R}^d)$ , this relation still holds if  $\boldsymbol{\phi} \in H_0^1(\Omega; \mathbb{R}^d)$ . On the other hand, plugging  $\underline{\mathbf{v}}_h = \underline{\mathbf{u}}_h$  into (3.21) and using the fact that  $s_h(\underline{\mathbf{u}}_h, \underline{\mathbf{u}}_h) \geq 0$ , we obtain

$$\mathfrak{I}_h := \int_{\Omega} \boldsymbol{\sigma}(\cdot, \mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h) : \mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h \leq \int_{\Omega} \mathbf{f} \cdot \underline{\mathbf{u}}_h.$$

Thus, using the previous bound, the strong convergence  $\underline{\mathbf{u}}_h \rightarrow \mathbf{u}$ , and (3.36), it is inferred that

$$\lim_{h \rightarrow 0} \mathfrak{I}_h \leq \int_{\Omega} \mathbf{f} \cdot \mathbf{u} = \int_{\Omega} \boldsymbol{\eta} : \nabla_s \mathbf{u}. \quad (3.37)$$

We now use the monotonicity assumption on  $\boldsymbol{\sigma}$  and the Minty trick [152] to prove that  $\boldsymbol{\eta} = \boldsymbol{\sigma}(\cdot, \nabla_s \mathbf{u})$ . Let  $\boldsymbol{\Lambda} \in L^2(\Omega; \mathbb{R}^{d \times d})$  and write, using the monotonicity (3.2e) of  $\boldsymbol{\sigma}$ , the convergence (3.31) of  $\boldsymbol{\sigma}(\cdot, \mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h)$ , and the bound (3.37),

$$0 \leq \lim_{h \rightarrow 0} \left( \int_{\Omega} (\boldsymbol{\sigma}(\cdot, \mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h) - \boldsymbol{\sigma}(\cdot, \boldsymbol{\Lambda})) : (\mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h - \boldsymbol{\Lambda}) \right) \leq \int_{\Omega} (\boldsymbol{\eta} - \boldsymbol{\sigma}(\cdot, \boldsymbol{\Lambda})) : (\nabla_s \mathbf{u} - \boldsymbol{\Lambda}). \quad (3.38)$$

Applying the previous relation with  $\boldsymbol{\Lambda} = \nabla_s \mathbf{u} \pm t \nabla_s \mathbf{v}$ , for  $t > 0$  and  $\mathbf{v} \in H_0^1(\Omega; \mathbb{R}^d)$ , and dividing by  $t$ , leads to

$$0 \leq \pm \int_{\Omega} (\boldsymbol{\eta} - \boldsymbol{\sigma}(\cdot, \nabla_s \mathbf{u} \mp t \nabla_s \mathbf{v})) : \nabla_s \mathbf{v}.$$

Owing to the growth property (3.2c) and the Caratheodory property (3.2a) of  $\sigma$ , we can let  $t \rightarrow 0$  and pass the limit inside the integral and then inside the argument of  $\sigma$ . In conclusion, for all  $\mathbf{v} \in H_0^1(\Omega; \mathbb{R}^d)$ , we infer

$$\int_{\Omega} \sigma(\cdot, \nabla_s \mathbf{u}) : \nabla_s \mathbf{v} = \int_{\Omega} \boldsymbol{\eta} : \nabla_s \mathbf{v} = \int_{\Omega} \mathbf{f} \cdot \mathbf{v},$$

where we have used (3.36) with  $\boldsymbol{\phi} = \mathbf{v}$  in order to obtain the second equality. The above equation shows that  $\boldsymbol{\eta} = \sigma(\cdot, \nabla_s \mathbf{u})$  and that  $\mathbf{u}$  solves (3.11).

**Step 4: Strong convergence.** We prove that if  $\sigma$  is strictly monotone then  $\mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h \rightarrow \nabla_s \mathbf{u}$  strongly in  $L^2(\Omega; \mathbb{R}^{d \times d})$ . We define the function  $\mathcal{D}_h : \Omega \rightarrow \mathbb{R}$  such that

$$\mathcal{D}_h := (\sigma(\cdot, \mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h) - \sigma(\cdot, \nabla_s \mathbf{u})) : (\mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h - \nabla_s \mathbf{u}).$$

For all  $h \in \mathcal{H}$ , the function  $\mathcal{D}_h$  is non-negative as a result of the monotonicity property (3.2e) and, by (3.38) with  $\boldsymbol{\Lambda} = \nabla_s \mathbf{u}$ , it is inferred that  $\lim_{h \rightarrow 0} \int_{\Omega} \mathcal{D}_h = 0$ . Hence,  $(\mathcal{D}_h)_{h \in \mathcal{H}}$  converges to 0 in  $L^1(\Omega)$  and, therefore, also almost everywhere on  $\Omega$  up to a subsequence. Let us take  $\bar{\mathbf{x}} \in \Omega$  such that the above mentioned convergence hold at  $\bar{\mathbf{x}}$ . Developing the products in  $\mathcal{D}_h$  and using the coercivity and growth properties (3.2d) and (3.2c) of  $\sigma$ , one has

$$\mathcal{D}_h(\bar{\mathbf{x}}) \geq \underline{\sigma} \|\mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h(\bar{\mathbf{x}})\|_{d \times d}^2 - 2\bar{\sigma} \|\mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h(\bar{\mathbf{x}})\|_{d \times d} \|\nabla_s \mathbf{u}(\bar{\mathbf{x}})\|_{d \times d} + \underline{\sigma} \|\nabla_s \mathbf{u}(\bar{\mathbf{x}})\|_{d \times d}^2.$$

Since the right hand side is quadratic in  $\|\mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h(\bar{\mathbf{x}})\|_{d \times d}$  and  $(\mathcal{D}_h(\bar{\mathbf{x}}))_{h \in \mathcal{H}}$  is bounded, we deduce that also  $(\mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h(\bar{\mathbf{x}}))_{h \in \mathcal{H}}$  is bounded. Passing to the limit in the definition of  $\mathcal{D}_h(\bar{\mathbf{x}})$  yields

$$(\sigma(\bar{\mathbf{x}}, \mathbf{L}_{\bar{\mathbf{x}}}) - \sigma(\bar{\mathbf{x}}, \nabla_s \mathbf{u}(\bar{\mathbf{x}}))) : (\mathbf{L}_{\bar{\mathbf{x}}} - \nabla_s \mathbf{u}(\bar{\mathbf{x}})) = 0,$$

where  $\mathbf{L}_{\bar{\mathbf{x}}}$  is an adherence value of  $(\mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h(\bar{\mathbf{x}}))_{h \in \mathcal{H}}$ . The strict monotonicity assumption forces  $\mathbf{L}_{\bar{\mathbf{x}}} = \nabla_s \mathbf{u}(\bar{\mathbf{x}})$  to be the unique adherence value of  $(\mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h(\bar{\mathbf{x}}))_{h \in \mathcal{H}}$ , and therefore the sequence converges to this value. As a result,

$$\mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h \rightarrow \nabla_s \mathbf{u} \text{ a.e. on } \Omega. \quad (3.39)$$

Using (3.37) together with Fatou's Lemma, we see that

$$\lim_{h \rightarrow 0} \int_{\Omega} \sigma(\cdot, \mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h) : \mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h = \int_{\Omega} \sigma(\cdot, \nabla_s \mathbf{u}) : \nabla_s \mathbf{u}.$$

Moreover, owing to (3.39),  $(\sigma(\cdot, \mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h) : \mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h)_{h \in \mathcal{H}}$  is a non-negative sequence converging almost everywhere on  $\Omega$ . Using [84, Lemma 8.4] we see that this sequence also converges in  $L^1(\Omega)$  and, therefore, it is equi-integrable in  $L^1(\Omega)$ . Thus, the coercivity (3.2d) of  $\sigma$  ensures that  $(\mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h)_{h \in \mathcal{H}}$  is equi-integrable in  $L^2(\Omega; \mathbb{R}^{d \times d})$  and Vitali's theorem shows that

$$\mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h \rightarrow \nabla_s \mathbf{u} \text{ strongly in } L^2(\Omega; \mathbb{R}^{d \times d}).$$

This concludes the proof.  $\square$

### 3.5.3 Error estimate

We next estimate the order of convergence of the method for solutions that display further regularity beyond that required by the variational formulation. We stipulate the following additional assumptions on the stress-strain function  $\sigma$ .

**Assumption 3.14** (Stress-strain relation II). There exist real numbers  $\sigma^*, \sigma_* \in (0, +\infty)$  such that, for a.e.  $\mathbf{x} \in \Omega$ , and all  $\boldsymbol{\tau}, \boldsymbol{\eta} \in \mathbb{R}_{\text{sym}}^{d \times d}$ ,

$$\|\sigma(\mathbf{x}, \boldsymbol{\tau}) - \sigma(\mathbf{x}, \boldsymbol{\eta})\|_{d \times d} \leq \sigma^* \|\boldsymbol{\tau} - \boldsymbol{\eta}\|_{d \times d}, \quad (\text{Lipschitz continuity}) \quad (3.40a)$$

$$(\sigma(\mathbf{x}, \boldsymbol{\tau}) - \sigma(\mathbf{x}, \boldsymbol{\eta})) : (\boldsymbol{\tau} - \boldsymbol{\eta}) \geq \sigma_* \|\boldsymbol{\tau} - \boldsymbol{\eta}\|_{d \times d}^2. \quad (\text{strong monotonicity}) \quad (3.40b)$$

*Remark 3.15* (Lipschitz continuity and strong monotonicity). It has been proved in [13, Lemma 4.1] that, under the assumptions (3.9), the stress-strain tensor function for the Hencky–Mises model is strongly monotone and Lipschitz-continuous, namely Assumption 3.14 holds. Also the isotropic damage model satisfies Assumption 3.14 if the damage function in (3.5) is, for instance, such that  $D(|\boldsymbol{\tau}|) = 1 - (1 + |\boldsymbol{\tau}|)^{-\frac{1}{2}}$ .

**Theorem 3.16** (Error estimate). *Let Assumptions 3.1 and 3.14 hold, and let  $(\mathcal{T}_h)_{h \in \mathcal{H}}$  be a regular mesh sequence. Let  $\mathbf{u}$  be the unique solution to (3.1). Let a polynomial degree  $k \geq 1$  be fixed, and, for all  $h \in \mathcal{H}$ , let  $\underline{\mathbf{u}}_h$  be the unique solution to (3.21) on the mesh  $\mathcal{T}_h$ . Then, under the additional regularity  $\mathbf{u} \in H^{k+2}(\mathcal{T}_h; \mathbb{R}^d)$  and  $\sigma(\cdot, \nabla_s \mathbf{u}) \in H^{k+1}(\mathcal{T}_h; \mathbb{R}^{d \times d})$ , it holds*

$$\|\nabla_s \mathbf{u} - \mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h\| + s_h(\underline{\mathbf{u}}_h, \underline{\mathbf{u}}_h)^{\frac{1}{2}} \leq Ch^{k+1} \left( \|\mathbf{u}\|_{H^{k+2}(\mathcal{T}_h; \mathbb{R}^d)} + \|\sigma(\cdot, \nabla_s \mathbf{u})\|_{H^{k+1}(\mathcal{T}_h; \mathbb{R}^{d \times d})} \right), \quad (3.41)$$

where  $C$  is a positive constant depending only on  $\Omega$ ,  $k$ , the mesh regularity parameter  $\varrho$ , the real numbers  $\bar{\sigma}$ ,  $\underline{\sigma}$ ,  $\sigma^*$ ,  $\sigma_*$  appearing in (3.2) and in (3.40), and an upper bound of  $\|f\|$ .

*Remark 3.17* (Locking-free error estimate). The proposed scheme, although different from the one of [74], is robust in the quasi-incompressible limit. The reason is that, as a result of the commuting property (3.18), we have  $\pi_T^k(\nabla \cdot \mathbf{v}) = \text{tr}(\mathbf{G}_{s,T}^k \mathbf{I}_T^k \mathbf{v})$ . Thus, considering, e.g., the linear elasticity stress-strain relation (3.3), we can proceed as in [74, Theorem 8] in order to prove that, when  $\mathbf{u} \in H^{k+2}(\mathcal{T}_h; \mathbb{R}^d)$  and  $\nabla \cdot \mathbf{u} \in H^{k+1}(\mathcal{T}_h; \mathbb{R})$ , and choosing  $\gamma = 2\mu$ , it holds

$$(2\mu)^{1/2} \|\nabla_s \mathbf{u} - \mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h\| \leq Ch^{k+1} \left( 2\mu \|\mathbf{u}\|_{H^{k+2}(\mathcal{T}_h; \mathbb{R}^d)} + \lambda \|\nabla \cdot \mathbf{u}\|_{H^{k+1}(\mathcal{T}_h; \mathbb{R})} \right),$$

with real number  $C > 0$  independent of  $h$ ,  $\mu$  and  $\lambda$ . The previous bound leads to a locking-free estimate; see [74, Remark 9]. Note that the locking-free nature of polyhedral element methods has also been observed in [197] for the Weak Galerkin method and in [20] for the Virtual Element method.

*Proof of Theorem 3.16.* For the sake of conciseness, throughout the proof we let  $\widehat{\underline{\mathbf{u}}}_h := \mathbf{I}_h^k \mathbf{u}$  and use the following abridged notations for the constraint field and its approximations:

$$\boldsymbol{\varsigma} := \sigma(\cdot, \nabla_s \mathbf{u}) \text{ and, for all } T \in \mathcal{T}_h, \boldsymbol{\varsigma}_T := \sigma(\cdot, \mathbf{G}_{s,T}^k \underline{\mathbf{u}}_T) \text{ and } \widehat{\boldsymbol{\varsigma}}_T := \sigma(\cdot, \mathbf{G}_{s,T}^k \widehat{\underline{\mathbf{u}}}_T).$$

First we want to show that (3.41) holds assuming that

$$\|\underline{\mathbf{u}}_h - \widehat{\underline{\mathbf{u}}}_h\|_{\epsilon, h} \lesssim h^{k+1} \left( \|\mathbf{u}\|_{H^{k+2}(\mathcal{T}_h; \mathbb{R}^d)} + \|\boldsymbol{\varsigma}\|_{H^{k+1}(\mathcal{T}_h; \mathbb{R}^{d \times d})} \right). \quad (3.42)$$

Using the triangle inequality, we obtain

$$\begin{aligned} \|\mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h - \nabla_s \mathbf{u}\| + s_h(\underline{\mathbf{u}}_h, \underline{\mathbf{u}}_h)^{1/2} &\leq \|\mathbf{G}_{s,h}^k(\underline{\mathbf{u}}_h - \widehat{\underline{\mathbf{u}}}_h)\| + s_h(\underline{\mathbf{u}}_h - \widehat{\underline{\mathbf{u}}}_h, \underline{\mathbf{u}}_h - \widehat{\underline{\mathbf{u}}}_h)^{1/2} \\ &\quad + \|\mathbf{G}_{s,h}^k \widehat{\underline{\mathbf{u}}}_h - \nabla_s \mathbf{u}\| + s_h(\widehat{\underline{\mathbf{u}}}_h, \widehat{\underline{\mathbf{u}}}_h)^{1/2}. \end{aligned} \quad (3.43)$$

Using the norm equivalence (3.24) followed by (3.42) we obtain for the terms in the first line of (3.43)

$$\|\mathbf{G}_{s,h}^k(\underline{\mathbf{u}}_h - \widehat{\underline{\mathbf{u}}}_h)\| + s_h(\underline{\mathbf{u}}_h - \widehat{\underline{\mathbf{u}}}_h, \underline{\mathbf{u}}_h - \widehat{\underline{\mathbf{u}}}_h)^{1/2} \lesssim h^{k+1} \left( \|\mathbf{u}\|_{H^{k+2}(\mathcal{T}_h; \mathbb{R}^d)} + \|\mathbf{s}\|_{H^{k+1}(\mathcal{T}_h; \mathbb{R}^{d \times d})} \right).$$

For the terms in the second line, using the approximation properties of  $\mathbf{G}_{s,h}^k$  resulting from (3.18) together with (3.13a) for the first addend and (3.34) for the second, we get

$$\|\mathbf{G}_{s,h}^k \widehat{\underline{\mathbf{u}}}_h - \nabla_s \mathbf{u}\| + s_h(\widehat{\underline{\mathbf{u}}}_h, \widehat{\underline{\mathbf{u}}}_h)^{1/2} \lesssim h^{k+1} \|\mathbf{u}\|_{H^{k+2}(\mathcal{T}_h; \mathbb{R}^d)}.$$

It only remains to prove (3.42), which we do in two steps: in **Step 1** we prove a basic estimate in terms of a conformity error, which is then bounded in **Step 2**.

**Step 1: Basic error estimate.** Using for all  $T \in \mathcal{T}_h$  the strong monotonicity (3.40b) with  $\boldsymbol{\tau} = \mathbf{G}_{s,T}^k \widehat{\underline{\mathbf{u}}}_T$  and  $\boldsymbol{\eta} = \mathbf{G}_{s,T}^k \underline{\mathbf{u}}_T$ , we infer

$$\|\mathbf{G}_{s,h}^k(\widehat{\underline{\mathbf{u}}}_h - \underline{\mathbf{u}}_h)\|^2 \lesssim \sum_{T \in \mathcal{T}_h} \int_T (\widehat{\boldsymbol{\mathcal{S}}}_T - \boldsymbol{\mathcal{S}}_T) : \mathbf{G}_{s,T}^k(\widehat{\underline{\mathbf{u}}}_T - \underline{\mathbf{u}}_T).$$

Owing to the norm equivalence (3.24) and the previous bound, we get

$$\begin{aligned} \|\widehat{\underline{\mathbf{u}}}_h - \underline{\mathbf{u}}_h\|_{\epsilon,h}^2 &\lesssim \sum_{T \in \mathcal{T}_h} \int_T (\widehat{\boldsymbol{\mathcal{S}}}_T - \boldsymbol{\mathcal{S}}_T) : \mathbf{G}_{s,T}^k(\widehat{\underline{\mathbf{u}}}_T - \underline{\mathbf{u}}_T) + s_h(\widehat{\underline{\mathbf{u}}}_h - \underline{\mathbf{u}}_h, \widehat{\underline{\mathbf{u}}}_h - \underline{\mathbf{u}}_h) \\ &= a_h(\widehat{\underline{\mathbf{u}}}_h, \widehat{\underline{\mathbf{u}}}_h - \underline{\mathbf{u}}_h) - \int_{\Omega} \mathbf{f} \cdot (\widehat{\underline{\mathbf{u}}}_h - \underline{\mathbf{u}}_h). \end{aligned}$$

where we have used the discrete problem (3.21) to conclude. Hence, dividing by  $\|\widehat{\underline{\mathbf{u}}}_h - \underline{\mathbf{u}}_h\|_{\epsilon,h}$  and passing to the supremum in the right-hand side, we arrive at the following error estimate:

$$\|\widehat{\underline{\mathbf{u}}}_h - \underline{\mathbf{u}}_h\|_{\epsilon,h} \lesssim \sup_{\underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_{h,D}^k, \|\underline{\mathbf{v}}_h\|_{\epsilon,h}=1} \mathcal{E}_h(\underline{\mathbf{v}}_h), \quad (3.44)$$

with conformity error such that, for all  $\underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_{h,D}^k$ ,

$$\mathcal{E}_h(\underline{\mathbf{v}}_h) := \sum_{T \in \mathcal{T}_h} \int_T \widehat{\boldsymbol{\mathcal{S}}}_T : \mathbf{G}_{s,T}^k \underline{\mathbf{v}}_T - \int_{\Omega} \mathbf{f} \cdot \underline{\mathbf{v}}_h + s_h(\widehat{\underline{\mathbf{u}}}_h, \underline{\mathbf{v}}_h). \quad (3.45)$$

**Step 2: Bound of the conformity error.** We bound the quantity  $\mathcal{E}_h(\underline{\mathbf{v}}_h)$  defined above for a generic  $\underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_{h,D}^k$ . Denote by  $\mathfrak{I}_1$ ,  $\mathfrak{I}_2$ , and  $\mathfrak{I}_3$  the three addends in the right-hand side of (3.45).

Using for all  $T \in \mathcal{T}_h$  the definition (3.17) of  $\mathbf{G}_{s,T}^k$  with  $\boldsymbol{\tau} = \boldsymbol{\pi}_T^k \widehat{\boldsymbol{\mathcal{S}}}_T$ , we have that

$$\mathfrak{I}_1 = \sum_{T \in \mathcal{T}_h} \left( \int_T \widehat{\boldsymbol{\mathcal{S}}}_T : \nabla_s \mathbf{v}_T + \sum_{F \in \mathcal{F}_T} \int_F \boldsymbol{\pi}_T^k \widehat{\boldsymbol{\mathcal{S}}}_T \mathbf{n}_{TF} \cdot (\mathbf{v}_F - \mathbf{v}_T) \right), \quad (3.46)$$

where we have used the fact that  $\nabla_s \mathbf{v}_T \in \mathbb{P}_d^{k-1}(T; \mathbb{R}^{d \times d})$  together with the definition (3.12) of the orthogonal projector to cancel  $\boldsymbol{\pi}_T^k$  in the first term.

On the other hand, using the fact that  $\mathbf{f} = -\nabla \cdot \boldsymbol{\mathcal{S}}$  a.e. in  $\Omega$  and integrating by parts element by element, we get that

$$\mathfrak{I}_2 = - \sum_{T \in \mathcal{T}_h} \left( \int_T \boldsymbol{\mathcal{S}} : \nabla_s \mathbf{v}_T + \sum_{F \in \mathcal{F}_T} \int_F \boldsymbol{\mathcal{S}} \mathbf{n}_{TF} \cdot (\mathbf{v}_F - \mathbf{v}_T) \right), \quad (3.47)$$

where we have additionally used that  $\boldsymbol{\mathcal{S}}_{|T_1} \mathbf{n}_{T_1 F} + \boldsymbol{\mathcal{S}}_{|T_2} \mathbf{n}_{T_2 F} = \mathbf{0}$  for all interfaces  $F \subset \partial T_1 \cap \partial T_2$  and that  $\mathbf{v}_F$  vanishes on  $\Gamma$  (cf. (3.20)) to insert  $\mathbf{v}_F$  into the second term.

Summing (3.46) and (3.47), taking absolute values, and using the Cauchy–Schwarz inequality to bound the right-hand side, we infer that

$$|\mathfrak{I}_1 + \mathfrak{I}_2| \leq \left( \sum_{T \in \mathcal{T}_h} \left( \|\boldsymbol{\mathcal{S}} - \widehat{\boldsymbol{\mathcal{S}}}_T\|_T^2 + h_T \|\boldsymbol{\mathcal{S}} - \boldsymbol{\pi}_T^k \widehat{\boldsymbol{\mathcal{S}}}_T\|_{\partial T}^2 \right) \right)^{1/2} \|\underline{\mathbf{v}}_h\|_{\epsilon, h}. \quad (3.48)$$

It only remains to bound the first factor. Let a mesh element  $T \in \mathcal{T}_h$  be fixed. Using the Lipschitz continuity (3.40a) with  $\boldsymbol{\tau} = \mathbf{G}_{s,T}^k \widehat{\underline{\mathbf{u}}}_T$  and  $\boldsymbol{\eta} = \nabla_s \mathbf{u}$  and the optimal approximation properties of  $\mathbf{G}_{s,T}^k I_T^k$  resulting from (3.18) together with (3.13a) with  $m = 1$  and  $s = k + 2$ , leads to

$$\|\boldsymbol{\mathcal{S}} - \widehat{\boldsymbol{\mathcal{S}}}_T\|_T \lesssim \|\nabla_s \mathbf{u} - \mathbf{G}_{s,T}^k \widehat{\underline{\mathbf{u}}}_T\|_T \lesssim h^{k+1} \|\mathbf{u}\|_{H^{k+2}(T; \mathbb{R}^d)}, \quad (3.49)$$

which provides an estimate for the first term inside the summation in the right-hand side of (3.48). To estimate the second term, we use the triangle inequality, the discrete trace inequality of [76, Lemma 1.46], and the boundedness of  $\boldsymbol{\pi}_T^k$  to write

$$h_T^{1/2} \|\boldsymbol{\mathcal{S}} - \boldsymbol{\pi}_T^k \widehat{\boldsymbol{\mathcal{S}}}_T\|_{\partial T} \lesssim \|\boldsymbol{\pi}_T^k (\boldsymbol{\mathcal{S}} - \widehat{\boldsymbol{\mathcal{S}}}_T)\|_T + h_T^{1/2} \|\boldsymbol{\mathcal{S}} - \boldsymbol{\pi}_T^k \boldsymbol{\mathcal{S}}\|_{\partial T} \leq \|\boldsymbol{\mathcal{S}} - \widehat{\boldsymbol{\mathcal{S}}}_T\|_T + h_T^{1/2} \|\boldsymbol{\mathcal{S}} - \boldsymbol{\pi}_T^k \boldsymbol{\mathcal{S}}\|_{\partial T}.$$

The first term in the right-hand side is bounded by (3.49). For the second, using the approximation properties (3.13b) of  $\boldsymbol{\pi}_T^k$  with  $m = 0$  and  $s = k + 1$ , we get  $h_T^{1/2} \|\boldsymbol{\mathcal{S}} - \boldsymbol{\pi}_T^k \boldsymbol{\mathcal{S}}\|_{\partial T} \lesssim h^{k+1} \|\boldsymbol{\mathcal{S}}\|_{H^{k+1}(T; \mathbb{R}^{d \times d})}$  so that, in conclusion,

$$h_T^{1/2} \|\boldsymbol{\mathcal{S}} - \boldsymbol{\pi}_T^k \widehat{\boldsymbol{\mathcal{S}}}_T\|_{\partial T} \lesssim h^{k+1} \left( \|\mathbf{u}\|_{H^{k+2}(T; \mathbb{R}^d)} + \|\boldsymbol{\mathcal{S}}\|_{H^{k+1}(T; \mathbb{R}^{d \times d})} \right). \quad (3.50)$$

Plugging the estimates (3.49) and (3.50) into (3.48) finally yields

$$|\mathfrak{I}_1 + \mathfrak{I}_2| \lesssim h^{k+1} \left( \|\mathbf{u}\|_{H^{k+2}(\mathcal{T}_h; \mathbb{R}^d)} + \|\boldsymbol{\mathcal{S}}\|_{H^{k+1}(\mathcal{T}_h; \mathbb{R}^{d \times d})} \right) \|\underline{\mathbf{v}}_h\|_{\epsilon, h}. \quad (3.51)$$

It only remains to bound  $\mathfrak{I}_3 = s_h(\widehat{\underline{\mathbf{u}}}_h, \underline{\mathbf{v}}_h)$ . Using the Cauchy–Schwarz inequality, the definition (3.22) of  $s_h$ , the approximation property (3.34) of  $\Delta_{TF}^k$ , and the norm equivalence (3.24), we infer

$$|\mathfrak{I}_3| \lesssim \left( \sum_{T \in \mathcal{T}_h} \sum_{F \in \mathcal{F}_T} \frac{\gamma}{h_F} \|\Delta_{TF}^k \widehat{\underline{\mathbf{u}}}_T\|_F^2 \right)^{1/2} s_h(\underline{\mathbf{v}}_h, \underline{\mathbf{v}}_h)^{1/2} \lesssim h^{k+1} \|\mathbf{u}\|_{H^{k+2}(\mathcal{T}_h; \mathbb{R}^d)} \|\underline{\mathbf{v}}_h\|_{\epsilon, h}. \quad (3.52)$$

Using (3.51) and (3.52), we finally get that, for all  $\underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_{h, \mathbf{D}}^k$ ,

$$\mathcal{E}_h(\underline{\mathbf{v}}_h) \lesssim h^{k+1} \left( \|\mathbf{u}\|_{H^{k+2}(\mathcal{T}_h; \mathbb{R}^d)} + \|\mathcal{S}\|_{H^{k+1}(\mathcal{T}_h; \mathbb{R}^{d \times d})} \right) \|\underline{\mathbf{v}}_h\|_{\epsilon, h}. \quad (3.53)$$

Thus, using (3.53) to bound the right-hand side of (3.44), (3.42) follows.  $\square$

## 3.6 Local principle of virtual work and law of action and reaction

We show in this section that the solution of the discrete problem (3.21) satisfies inside each element a local principle of virtual work with numerical tractions that obey the law of action and reaction. This property is important from both the mathematical and engineering points of view, and it can simplify the derivation of a posteriori error estimators based on equilibrated tractions; see, e.g., [3, 158]. It is worth emphasizing that local equilibrium properties on the primal mesh are a distinguishing feature of hybrid (face-based) methods: the derivation of similar properties for vertex-based methods usually requires to perform reconstructions on a dual mesh.

Define, for all  $T \in \mathcal{T}_h$ , the space

$$\underline{\mathbf{D}}_{\partial T}^k := \times_{F \in \mathcal{F}_T} \mathbb{P}_d^k(F; \mathbb{R}^d),$$

as well as the boundary difference operator  $\underline{\delta}_{\partial T}^k : \underline{\mathbf{U}}_T^k \rightarrow \underline{\mathbf{D}}_{\partial T}^k$  such that, for all  $\underline{\mathbf{v}}_T \in \underline{\mathbf{U}}_T^k$ ,

$$\underline{\delta}_{\partial T}^k \underline{\mathbf{v}}_T = (\delta_F^k \underline{\mathbf{v}}_T)_{F \in \mathcal{F}_T} := (\mathbf{v}_F - \mathbf{v}_T)_{F \in \mathcal{F}_T}.$$

The following proposition shows that the stabilization can be reformulated in terms of boundary differences.

**Proposition 3.18** (Reformulation of the local stabilization bilinear form). *For all mesh elements  $T \in \mathcal{T}_h$ , the local stabilization bilinear form  $s_T$  defined by (3.22) satisfies, for all  $\underline{\mathbf{u}}_T, \underline{\mathbf{v}}_T \in \underline{\mathbf{U}}_T^k$ ,*

$$s_T(\underline{\mathbf{u}}_T, \underline{\mathbf{v}}_T) = s_T((\mathbf{0}, \underline{\delta}_{\partial T}^k \underline{\mathbf{u}}_T), (\mathbf{0}, \underline{\delta}_{\partial T}^k \underline{\mathbf{v}}_T)). \quad (3.54)$$

*Proof.* Let a mesh element  $T \in \mathcal{T}_h$  be fixed. Using the fact that  $\mathbf{r}_T^{k+1} \underline{\mathbf{I}}_T^k \underline{\mathbf{v}}_T = \mathbf{v}_T$  for all  $\underline{\mathbf{v}}_T \in \mathbb{P}_d^k(T)^d$  (this because  $\mathbf{r}_T^{k+1} \underline{\mathbf{I}}_T^k$  is a projector on  $\mathbb{P}_d^{k+1}(T; \mathbb{R}^d)$ , cf. [74, Eq. (20)]) together

with the linearity of  $\mathbf{r}_T^{k+1}$ , it is inferred that, for all  $F \in \mathcal{F}_T$ , the face-based residual defined by (3.23) satisfies

$$\Delta_{TF}^k \mathbf{v}_T = \boldsymbol{\pi}_F^k(\mathbf{r}_T^{k+1}(\mathbf{0}, \underline{\boldsymbol{\delta}}_{\partial T}^k \mathbf{v}_T) - \boldsymbol{\delta}_F^k \mathbf{v}_T) - \boldsymbol{\pi}_T^k \mathbf{r}_T^{k+1}(\mathbf{0}, \underline{\boldsymbol{\delta}}_{\partial T}^k \mathbf{v}_T) = \Delta_{TF}^k(\mathbf{0}, \underline{\boldsymbol{\delta}}_{\partial T}^k \mathbf{v}_T)$$

for all  $\mathbf{v}_T \in \underline{\mathbf{U}}_T^k$ . Plugging this expression into (3.22) yields the assertion.  $\square$

Define now the boundary residual operator  $\mathbf{R}_{\partial T}^k : \underline{\mathbf{U}}_T^k \rightarrow \underline{\mathbf{D}}_{\partial T}^k$  such that, for all  $\mathbf{v}_T \in \underline{\mathbf{U}}_T^k$ ,  $\mathbf{R}_{\partial T}^k \mathbf{v}_T := (\mathbf{R}_{TF}^k \mathbf{v}_T)_{F \in \mathcal{F}_T}$  satisfies

$$- \sum_{F \in \mathcal{F}_T} \int_F \mathbf{R}_{TF}^k \mathbf{v}_T \cdot \boldsymbol{\alpha}_F = s_T((\mathbf{0}, \underline{\boldsymbol{\delta}}_{\partial T}^k \mathbf{v}_T), (\mathbf{0}, \underline{\boldsymbol{\alpha}}_{\partial T})) \quad \forall \underline{\boldsymbol{\alpha}}_{\partial T} \in \underline{\mathbf{D}}_{\partial T}^k. \quad (3.55)$$

Problem (3.55) is well-posed, and computing  $\mathbf{R}_{TF}^k \mathbf{v}_T$  requires to invert the boundary mass matrix.

**Lemma 3.19** (Local principle of virtual work and law of action and reaction). *Denote by  $\underline{\mathbf{u}}_h \in \underline{\mathbf{U}}_{h,D}^k$  a solution of problem (3.21) and, for all  $T \in \mathcal{T}_h$  and all  $F \in \mathcal{F}_T$ , define the numerical traction*

$$\mathbf{T}_{TF}(\underline{\mathbf{u}}_T) := -\boldsymbol{\pi}_T^k \boldsymbol{\sigma}(\cdot, \mathbf{G}_{s,T}^k \underline{\mathbf{u}}_T) \mathbf{n}_{TF} + \mathbf{R}_{TF}^k \underline{\mathbf{u}}_T.$$

Then, for all  $T \in \mathcal{T}_h$  we have the following discrete principle of virtual work: For all  $\mathbf{v}_T \in \mathbb{P}_d^k(T; \mathbb{R}^d)$ ,

$$\int_T \boldsymbol{\sigma}(\cdot, \mathbf{G}_{s,T}^k \underline{\mathbf{u}}_T) : \nabla_s \mathbf{v}_T + \sum_{F \in \mathcal{F}_T} \int_F \mathbf{T}_{TF}(\underline{\mathbf{u}}_T) \cdot \mathbf{v}_T = \int_T \mathbf{f} \cdot \mathbf{v}_T, \quad (3.56)$$

and, for any interface  $F \in \mathcal{F}_{T_1} \cap \mathcal{F}_{T_2}$ , the numerical tractions satisfy the law of action and reaction:

$$\mathbf{T}_{T_1 F}(\underline{\mathbf{u}}_{T_1}) + \mathbf{T}_{T_2 F}(\underline{\mathbf{u}}_{T_2}) = \mathbf{0}. \quad (3.57)$$

*Proof.* For all  $T \in \mathcal{T}_h$ , use the definition (3.17) of  $\mathbf{G}_{s,T}^k \mathbf{v}_T$  with  $\boldsymbol{\tau} = \boldsymbol{\pi}_T^k \boldsymbol{\sigma}(\cdot, \mathbf{G}_{s,T}^k \underline{\mathbf{u}}_T)$  in  $A_h$  and the rewriting (3.54) of  $s_T$  together with the definition (3.55) of  $\mathbf{R}_{TF}^k$  to infer that it holds, for all  $\mathbf{v}_h \in \underline{\mathbf{U}}_h^k$ ,

$$\begin{aligned} \int_{\Omega} \mathbf{f} \cdot \mathbf{v}_h &= A_h(\underline{\mathbf{u}}_h, \mathbf{v}_h) + s_h(\underline{\mathbf{u}}_h, \mathbf{v}_h) \\ &= \sum_{T \in \mathcal{T}_h} \int_T \boldsymbol{\sigma}(\cdot, \mathbf{G}_{s,T}^k \underline{\mathbf{u}}_T) : \nabla_s \mathbf{v}_T \\ &\quad + \sum_{T \in \mathcal{T}_h} \sum_{F \in \mathcal{F}_T} \int_F (\boldsymbol{\pi}_T^k \boldsymbol{\sigma}(\cdot, \mathbf{G}_{s,T}^k \underline{\mathbf{u}}_T) \mathbf{n}_{TF} - \mathbf{R}_{TF}^k \underline{\mathbf{u}}_T) \cdot (\mathbf{v}_F - \mathbf{v}_T), \end{aligned}$$

where to cancel  $\boldsymbol{\pi}_T^k$  inside the first integral in the second line we have used the fact that  $\nabla_s \mathbf{v}_T \in \mathbb{P}_d^{k-1}(T; \mathbb{R}^{d \times d})$  for all  $T \in \mathcal{T}_h$ . Selecting  $\mathbf{v}_h$  such that  $\mathbf{v}_T$  spans  $\mathbb{P}_d^k(T; \mathbb{R}^d)$  for a selected mesh element  $T \in \mathcal{T}_h$  while  $\mathbf{v}_{T'} \equiv \mathbf{0}$  for all  $T' \in \mathcal{T}_h \setminus \{T\}$  and  $\mathbf{v}_F \equiv \mathbf{0}$  for all  $F \in \mathcal{F}_h$ , we obtain (3.56). On the other hand, selecting  $\mathbf{v}_h$  such that  $\mathbf{v}_T \equiv \mathbf{0}$  for all  $T \in \mathcal{T}_h$ ,  $\mathbf{v}_F$  spans  $\mathbb{P}_d^k(F; \mathbb{R}^d)$  for a selected interface  $F \in \mathcal{F}_{T_1} \cap \mathcal{F}_{T_2}$ , and  $\mathbf{v}_{F'} \equiv \mathbf{0}$  for all  $F' \in \mathcal{F}_h \setminus \{F\}$  yields (3.57).  $\square$



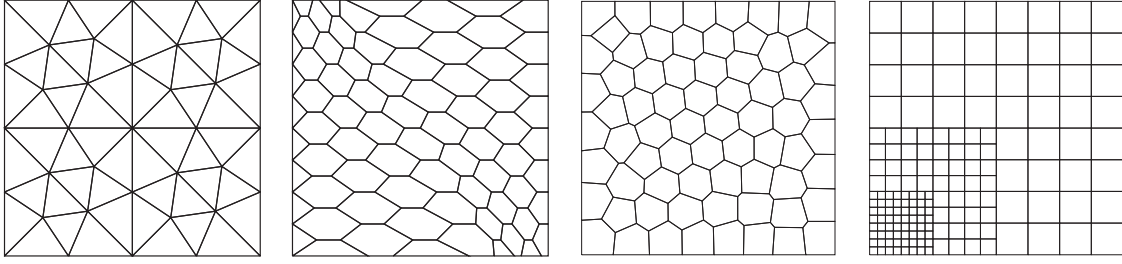


Figure 3.2: Triangular, hexagonal-dominant, Voronoi, and nonmatching quadrangular meshes for the numerical tests. The triangular and nonmatching quadrangular meshes were originally proposed for the FVCA5 benchmark [120]. The (predominantly) hexagonal was used in [80]. The Voronoi mesh family was obtained using the PolyMesher algorithm of [188].

## 3.7 Numerical results

In this section we present a comprehensive set of numerical tests to assess the properties of our method using the models of Examples 3.2, 3.3, and 3.5 (cf. also Remark 3.6). Note that an important step in the implementation of HHO methods consists in selecting a basis for each of the polynomial spaces that appear in the construction. In the numerical tests of the present section, for all  $T \in \mathcal{T}_h$ , we take as a basis for  $\mathbb{P}_d^k(T; \mathbb{R}^d)$  the Cartesian product of the monomials in the translated and scaled coordinates  $(h_T^{-1}(x_i - x_{T,i}))_{1 \leq i \leq d}$ , where  $\mathbf{x}_T$  is the barycenter of  $T$ . Similarly, for all  $F \in \mathcal{F}_h$  we define a basis for  $\mathbb{P}_d^k(F; \mathbb{R}^d)$  by taking the monomials with respect to a local frame scaled using the face diameter  $h_F$  and the middle point of  $F$ . Further details on implementation aspects are given in [74, Section 6.1].

### 3.7.1 Convergence for the Hencky–Mises model

In order to check the error estimates stated in Theorem 3.16, we first solve a manufactured two-dimensional hyperelasticity problem. We consider the Hencky–Mises model with  $\Phi(\rho) = \mu(e^{-\rho} + 2\rho)$  and  $\alpha = \lambda + \mu$  in (3.8), so that the conditions (3.9) are satisfied. This choice leads to the following stress-strain relation:

$$\boldsymbol{\sigma}(\nabla_s \mathbf{u}) = ((\lambda - \mu) + \mu e^{-\text{dev}(\nabla_s \mathbf{u})}) \text{tr}(\nabla_s \mathbf{u}) \mathbf{I}_d + \mu(2 - e^{-\text{dev}(\nabla_s \mathbf{u})}) \nabla_s \mathbf{u}.$$

We consider the unit square domain  $\Omega = (0, 1)^2$  and take  $\mu = 2$ ,  $\lambda = 1$ , and an exact displacement  $\mathbf{u}$  given by

$$\mathbf{u}(\mathbf{x}) = (\sin(\pi x_1) \sin(\pi x_2), \sin(\pi x_1) \sin(\pi x_2)).$$

The volumetric load  $\mathbf{f} = -\nabla \cdot \boldsymbol{\sigma}(\nabla_s \mathbf{u})$  is inferred from the exact solution  $\mathbf{u}$ . In this case, since the selected exact displacement vanishes on  $\Gamma$ , we simply consider homogeneous Dirichlet conditions. We consider the triangular, hexagonal, Voronoi, and nonmatching quadrangular mesh families depicted in Figure 3.2 and polynomial degrees  $k$  ranging from 1 to 4. The nonmatching mesh is simply meant to show that the method supports nonconforming interfaces: refining in the corner has no particular meaning for the selected solution. The

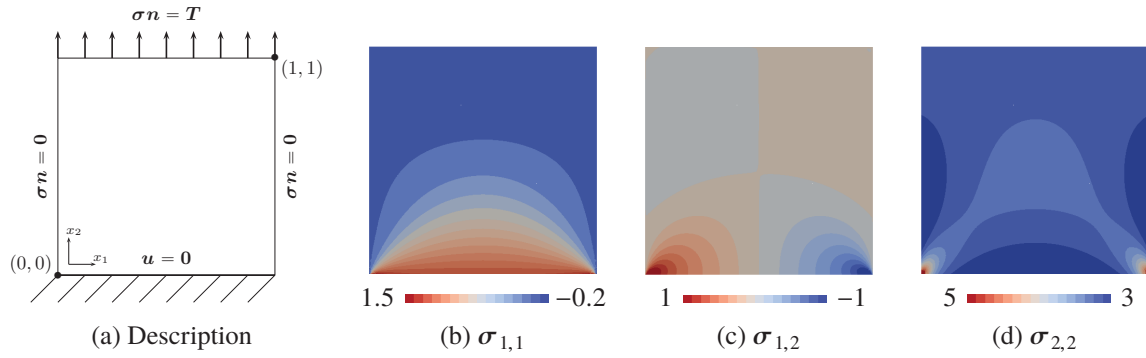


Figure 3.3: Tensile test description and resulting stress components for the linear case. Values in  $10^5$ Pa

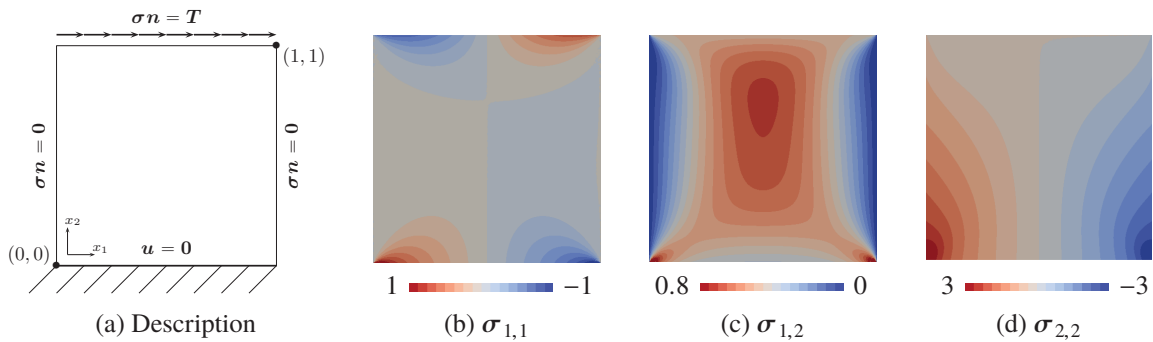


Figure 3.4: Shear test description and resulting stress components for the linear case. Values in  $10^5$ Pa

initialization of our iterative linearization procedure (Newton scheme) is obtained solving the linear elasticity model. This initial guess leads to a 40% reduction of the number of iterations with respect to a null initial guess. The energy-norm orders of convergence (OCV) displayed in the third column of Tables 3.1–3.4 at the end of the chapter are in agreement with the theoretical predictions. In particular, it is observed that the optimal convergence in  $h^{k+1}$  is reached for the triangular, nonmatching Cartesian, and hexagonal meshes for  $1 \leq k \leq 3$ , whereas for  $k = 4$  the asymptotic convergence order does not appear to have been reached in the last mesh refinement. It can also be observed in Table 3.2 that the convergence rate exceeds the estimated one on the locally refined Cartesian mesh for  $k = 1$  and  $k = 2$ . For the sake of completeness, we also display in the fourth column of Tables 3.1–3.4 the  $L^2$ -norm of the error defined as the difference between the  $L^2$ -projection  $\pi_h^k \mathbf{u}$  of the exact solution on  $\mathbb{P}_d^k(\mathcal{T}_h; \mathbb{R}^d)$  and the broken polynomial function  $\mathbf{u}_h$  obtained from element-based DOFs, while in the fifth column we display the corresponding observed convergence rates. In this case, orders of convergence up to  $h^{k+2}$  are observed.

### 3.7.2 Tensile and shear test cases

We next consider the two test cases schematically depicted in Figures 3.3 and 3.4. On the unit square domain  $\Omega$ , we solve problem (3.1a) considering three different models of hyperelasticity (see Remark 3.6):

- (i) *Linear*. The linear model corresponding to the stored energy density function (3.7) with Lamé's parameters

$$\lambda = 11 \times 10^5 \text{Pa}, \quad \mu = 82 \times 10^4 \text{Pa}. \quad (3.58)$$

- (ii) *Hencky–Mises*. The Hencky–Mises model (3.4) obtained by taking  $\Phi(\rho) = \mu(\frac{\rho}{2} + (1 + \rho)^{1/2})$  and  $\alpha = \lambda + \mu$  in (3.8), with  $\lambda, \mu$  as in (3.58) (also in this case the conditions (3.9) hold). This choice leads to

$$\begin{aligned} \boldsymbol{\sigma}(\nabla_s \mathbf{u}) &= ((\lambda + \frac{\mu}{2}) - \frac{\mu}{2}(1 + \text{dev}(\nabla_s \mathbf{u}))^{-1/2}) \text{tr}(\nabla_s \mathbf{u}) \mathbf{I}_d \\ &\quad + \mu(1 + (1 + \text{dev}(\nabla_s \mathbf{u}))^{-1/2}) \nabla_s \mathbf{u}. \end{aligned} \quad (3.59)$$

The Lamé's functions of the previous relation are inspired from those proposed in [21, Section 5.1]. In particular, the function  $\tilde{\mu}(\rho) = \mu(1 + (1 + \text{dev}(\nabla_s \mathbf{u}))^{-1/2})$  corresponds to the Carreau law for viscoplastic materials.

- (iii) *Second-order*. The second-order model (3.6) with Lamé's parameter as in (3.58) and second-order moduli

$$A = 11 \times 10^6 \text{Pa}, \quad B = -48 \times 10^5 \text{Pa}, \quad C = 13.2 \times 10^5 \text{Pa}.$$

These values correspond to the estimates provided in [124] for the Armco Iron. We recall that the second-order elasticity stress-strain relation does not satisfy in general the assumptions under which we are able to prove the convergence and error estimates. In particular, we observe that the stored energy density function defined in (3.10) is not convex.

The bottom part of the boundary of the domain is assumed to be fixed, the normal stress is equal to zero on the two lateral parts, and a traction is imposed at the top of the boundary. So, mixed boundary conditions are imposed as follows

$$\begin{aligned} \mathbf{u} &= \mathbf{0} && \text{on } \{\mathbf{x} \in \Gamma, x_2 = 0\}, \\ \boldsymbol{\sigma} \mathbf{n} &= \mathbf{T} && \text{on } \{\mathbf{x} \in \Gamma, x_2 = 1\}, \\ \boldsymbol{\sigma} \mathbf{n} &= \mathbf{0} && \text{on } \{\mathbf{x} \in \Gamma, x_1 = 0\}, \\ \boldsymbol{\sigma} \mathbf{n} &= \mathbf{0} && \text{on } \{\mathbf{x} \in \Gamma, x_1 = 1\}. \end{aligned}$$

For the tensile case, we impose a vertical traction at the top of the boundary equal to  $\mathbf{T} = (0, 3.2 \times 10^5 \text{Pa})$ . This type of boundary conditions produces large normal stresses (i.e., the diagonal components of  $\boldsymbol{\sigma}$ ) and minor shear stresses (i.e., the off-diagonal components of  $\boldsymbol{\sigma}$ ). It can be observed in Figure 3.3, where the components of the stress tensor are depicted for the linear case. In Figure 3.5 we plot the stress norm on the deformed domain obtained for the three hyperelasticity models. The results of Fig. 3.3, 3.4, 3.5, and 3.6 are obtained on a mesh with 3584 triangles (corresponding to a typical mesh-size of  $3.84 \times 10^{-3}$ ) and with polynomial degree  $k = 2$ . Obviously, the symmetry of the results is visible, and we observe that the three displacement fields are very close. This is motivated by the fact that, with our choice of the parameters in (3.58) and in (3.59), the linear model exactly corresponds to the linear approximation at the origin of the nonlinear ones. The maximum value of the stress

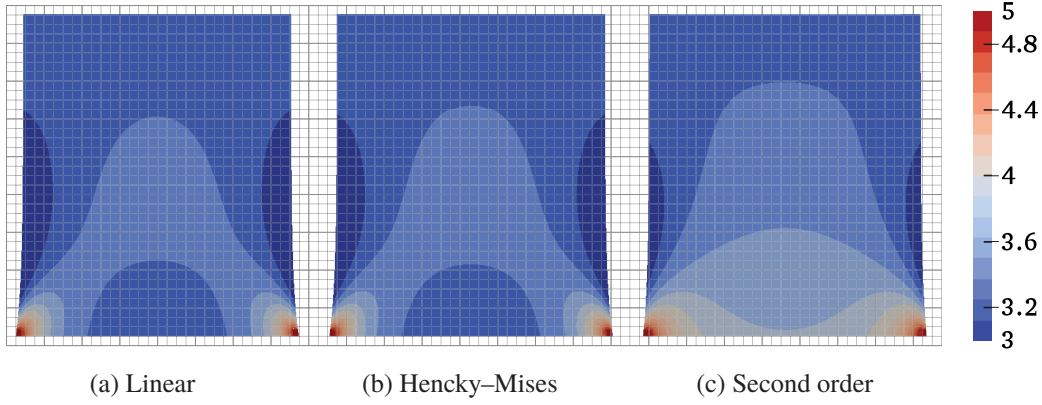


Figure 3.5: Tensile test case: Stress norm on the deformed domain. Values in  $10^5$ Pa

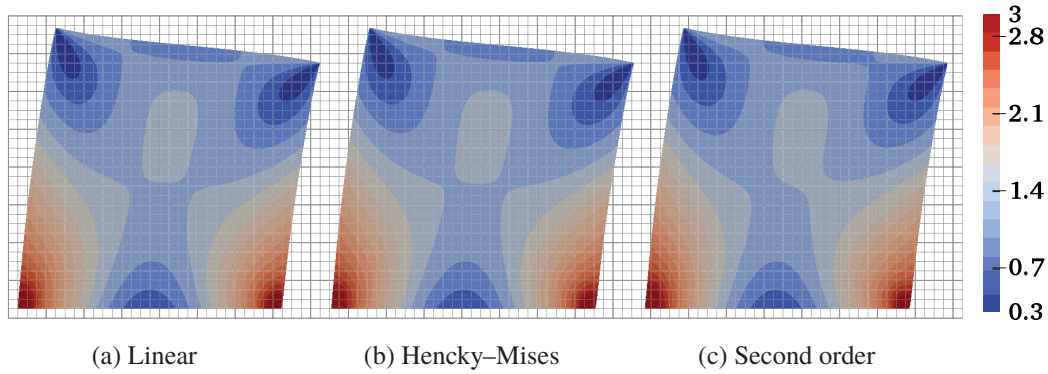


Figure 3.6: Shear test case: Stress norm on the deformed domain. Values in  $10^4$ Pa

concentrates on the two bottom corners due to the homogeneous Dirichlet condition that totally locks the displacement when  $x_1 = 0$ . The repartition of the stress on the domain with the second-order model is visibly different from those obtained with the linear and Hencky–Mises models. At the energy level, we also have a higher difference between the second-order model and the linear one since  $|\mathbb{E}_{\text{lin}} - \mathbb{E}_{\text{hm}}|/\mathbb{E}_{\text{lin}} = 0.44\%$  while  $|\mathbb{E}_{\text{lin}} - \mathbb{E}_{\text{snd}}|/\mathbb{E}_{\text{lin}} = 4.45\%$ , where  $\mathbb{E}_{\bullet}$  is the total elastic energy obtained by integrating over the domain the strain energy density functions defined by (3.7), (3.8), and (3.10):

$$\mathbb{E}_{\bullet} := \int_{\Omega} \Psi_{\bullet}, \quad \text{with } \bullet \in \{\text{lin}, \text{hm}, \text{snd}\}.$$

The reference values for the total energy, used in Figure 3.7 in order to assess convergence, are obtained on a fine Cartesian mesh having a mesh-size of  $1.95 \times 10^{-3}$  and  $k = 3$ . For the shear case, we consider an horizontal traction equal to  $\mathbf{T} = (4.5 \times 10^4 \text{Pa}, 0)$  which induces the stress pattern illustrated in Figure 3.4. The computed stress norm on the deformed domain is depicted in Figure 3.6, and we can see that the displacement fields associated with the three models are very close as for the tensile test case. Here, the maximum values of the stress are localized in the lower part of the domain near the lateral parts. Unlike the tensile test, the

difference between the three models is tiny as confirmed by the elastic energy equal to 3180 J, 3184 J, and 3190 J respectively. The decreasing of the energy difference in comparison with the previous test can be explained by the fact that the value of the Neumann boundary data on the top is divided by a factor 7 in order to obtain maximum displacements roughly equal to 15%.

## 3.8 Appendix: Technical results

This appendix contains the proof of technical results.

### 3.8.1 Discrete Korn inequality

**Proposition 3.20** (Discrete Korn inequality). *Assume that the mesh further verifies the assumption of [39, Theorem 4.2] if  $d = 2$  and [39, Theorem 5.2] if  $d = 3$ . Then, the discrete Korn inequality (3.25) holds.*

*Proof.* Using the broken Korn inequality [39, Eq. (1.22)] on  $H^1(\mathcal{T}_h; \mathbb{R}^d)$  followed by the Cauchy–Schwarz inequality, one has

$$\begin{aligned} \|\mathbf{v}_h\|^2 + \|\nabla_h \mathbf{v}_h\|^2 &\lesssim \|\nabla_{s,h} \mathbf{v}_h\|^2 + \sum_{F \in \mathcal{F}_h^i} h_F^{-1} \|[\mathbf{v}_h]_F\|_F^2 \\ &\quad + \sup_{\mathbf{m} \in \mathbb{P}_d^1(\mathcal{T}_h; \mathbb{R}^d), \|\boldsymbol{\gamma}_n(\mathbf{m})\|_\Gamma=1} \left( \int_\Gamma \boldsymbol{\gamma}(\mathbf{v}_h) \cdot \boldsymbol{\gamma}_n(\mathbf{m}) \right)^2 \\ &\lesssim \|\nabla_{s,h} \mathbf{v}_h\|^2 + \sum_{F \in \mathcal{F}_h^i} h_F^{-1} \|[\mathbf{v}_h]_F\|_F^2 + \sum_{F \in \mathcal{F}_h^b} \|\mathbf{v}_h|_F\|_F^2. \end{aligned} \quad (3.60)$$

For an interface  $F \in \mathcal{F}_{T_1} \cap \mathcal{F}_{T_2}$ , we have introduced the jump  $[\mathbf{v}_h]_F := \mathbf{v}_{T_1} - \mathbf{v}_{T_2}$ . Thus, using the triangle inequality, we get  $\|[\mathbf{v}_h]_F\|_F \leq \|\mathbf{v}_F - \mathbf{v}_{T_1}\|_F + \|\mathbf{v}_F - \mathbf{v}_{T_2}\|_F$ . For a boundary face  $F \in \mathcal{F}_h^b$  such that  $F \in \mathcal{F}_T \cap \mathcal{F}_h^b$  for some  $T \in \mathcal{T}_h$  we have, on the other hand,  $\|\mathbf{v}_h|_F\|_F = \|\mathbf{v}_F - \mathbf{v}_T\|_F$  since  $\mathbf{v}_F \equiv \mathbf{0}$  (cf. (3.20)). Using these relations in the right-hand side of (3.60) and rearranging the sums leads to

$$\begin{aligned} \|\mathbf{v}_h\|^2 + \|\nabla_h \mathbf{v}_h\|^2 &\lesssim \sum_{T \in \mathcal{T}_h} \left( \|\nabla_s \mathbf{v}_T\|_T^2 + \sum_{F \in \mathcal{F}_T \cap \mathcal{F}_h^i} h_F^{-1} \|\mathbf{v}_F - \mathbf{v}_T\|_F^2 \right) + h \sum_{F \in \mathcal{F}_h^b} h_F^{-1} \|\mathbf{v}_F - \mathbf{v}_h|_F\|_F^2 \\ &\lesssim \max\{1, d_\Omega\} \sum_{T \in \mathcal{T}_h} \left( \|\nabla_s \mathbf{v}_T\|_T^2 + \sum_{F \in \mathcal{F}_T} h_F^{-1} \|\mathbf{v}_F - \mathbf{v}_T\|_F^2 \right), \end{aligned}$$

where  $d_\Omega$  denotes the diameter of  $\Omega$ . Owing to the definition (3.14) of the discrete strain seminorm, the latter yields the assertion.  $\square$

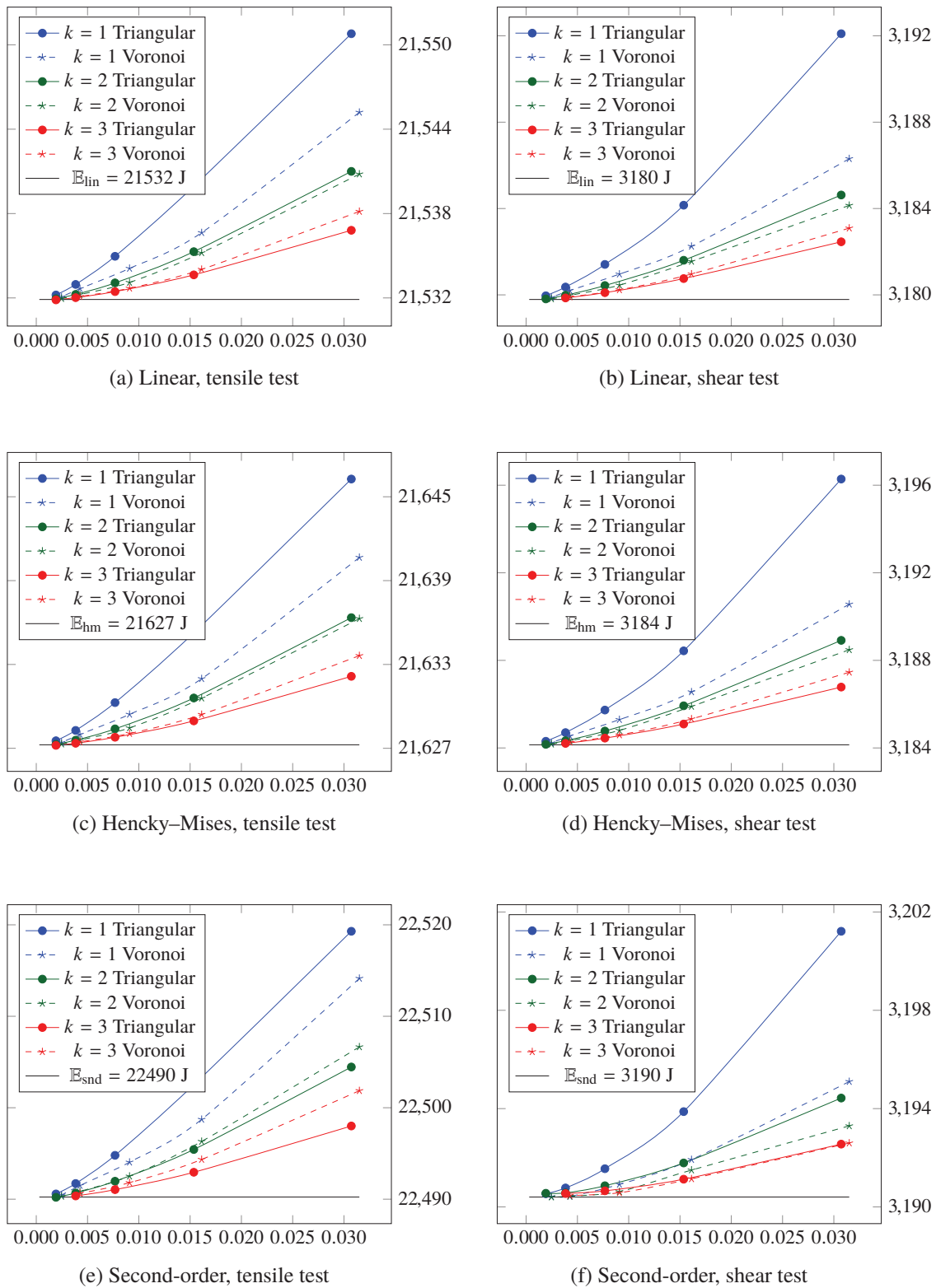


Figure 3.7: Energy vs  $h$ , tensile and shear test cases

### 3.8.2 Discrete compactness

*Proof of Lemma 3.12.* In the proof we use the same notation for functions in  $L^2(\Omega; \mathbb{R}^d) \subset L^1(\Omega; \mathbb{R}^d)$  and for their extension by zero outside  $\Omega$ . Let  $(\underline{\mathbf{v}}_h)_{h \in \mathcal{H}} \in (\underline{\mathbf{U}}_{h,D}^k)_{h \in \mathcal{H}}$  be such that (3.26) holds. Define the space of integrable functions with bounded variation  $\text{BV}(\mathbb{R}^d) := \{\mathbf{v} \in L^1(\mathbb{R}^d; \mathbb{R}^d) \mid \|\mathbf{v}\|_{\text{BV}} < +\infty\}$ , where

$$\|\mathbf{v}\|_{\text{BV}} := \sum_{i=1}^d \sup \left\{ \int_{\mathbb{R}^d} \mathbf{v} \cdot \partial_i \boldsymbol{\phi} \mid \boldsymbol{\phi} \in C_c^\infty(\mathbb{R}^d; \mathbb{R}^d), \|\boldsymbol{\phi}\|_{L^\infty(\mathbb{R}^d; \mathbb{R}^d)} \leq 1 \right\}.$$

Here,  $\partial_i \boldsymbol{\phi}$  denotes the  $i$ -th column of  $\nabla \boldsymbol{\phi}$ . Let  $\boldsymbol{\phi} \in C_c^\infty(\mathbb{R}^d; \mathbb{R}^d)$  with  $\|\boldsymbol{\phi}\|_{L^\infty(\mathbb{R}^d; \mathbb{R}^d)} \leq 1$ . Integrating by parts and using the fact that  $\sum_{T \in \mathcal{T}_h} \sum_{F \in \mathcal{F}_T} \int_F (\mathbf{v}_F \cdot \boldsymbol{\phi}) \mathbf{n}_{TF} = 0$ , we have that

$$\begin{aligned} \int_{\mathbb{R}^d} \mathbf{v}_h \cdot \partial_i \boldsymbol{\phi} &= \sum_{T \in \mathcal{T}_h} \int_T ((\nabla \boldsymbol{\phi})^\top \mathbf{v}_T)_i \\ &= - \sum_{T \in \mathcal{T}_h} \left( \int_T ((\nabla \mathbf{v}_T)^\top \boldsymbol{\phi})_i + \sum_{F \in \mathcal{F}_T} \int_F (\mathbf{v}_F \cdot \boldsymbol{\phi} - \mathbf{v}_T \cdot \boldsymbol{\phi}) (\mathbf{n}_{TF})_i \right) \\ &\leq \sum_{T \in \mathcal{T}_h} \left( \int_T \sum_{j=1}^d |(\nabla \mathbf{v}_T)_{ji}| + \sum_{F \in \mathcal{F}_T} \int_F \sum_{j=1}^d |(\mathbf{v}_F - \mathbf{v}_T)_j (\mathbf{n}_{TF})_i| \right), \end{aligned}$$

where, in order to pass to the third line, we have used  $\|\boldsymbol{\phi}\|_{L^\infty(\mathbb{R}^d; \mathbb{R}^d)} \leq 1$ . Therefore, summing over  $i \in \{1, \dots, d\}$ , observing that, for all  $T \in \mathcal{T}_h$  and all  $F \in \mathcal{F}_T$ , we have  $\sum_{i=1}^d |(\mathbf{n}_{TF})_i| \leq d^{1/2}$ , and using the Lebesgue embeddings arising from the Hölder inequality on bounded domain, leads to

$$\|\mathbf{v}_h\|_{\text{BV}} \lesssim \sum_{T \in \mathcal{T}_h} \left( |T|_d^{1/2} \|\nabla \mathbf{v}_T\|_T + \sum_{F \in \mathcal{F}_T} |F|_{d-1}^{1/2} \|\mathbf{v}_F - \mathbf{v}_T\|_F \right),$$

where  $|\cdot|_d$  denotes the  $d$ -dimensional Hausdorff measure. Moreover, using the Cauchy–Schwarz inequality together with the geometric bound  $|F|_{d-1} h_F \lesssim |T|_d$ , we obtain that

$$\|\mathbf{v}_h\|_{\text{BV}} \lesssim |\Omega|_d^{1/2} \left( \sum_{T \in \mathcal{T}_h} \left[ \|\nabla \mathbf{v}_T\|_T^2 + \sum_{F \in \mathcal{F}_T} h_F^{-1} \|\mathbf{v}_F - \mathbf{v}_T\|_F^2 \right] \right)^{1/2}.$$

Thus, using the discrete Korn inequality (3.25), it is readily inferred that

$$\|\mathbf{v}_h\|_{\text{BV}} \lesssim \|\underline{\mathbf{v}}_h\|_{\epsilon, h} \lesssim 1.$$

Owing to the Helly selection principle [99, Section 5.2.3], the sequence  $(\mathbf{v}_h)_{h \in \mathcal{H}}$  is relatively compact in  $L^1(\mathbb{R}^d; \mathbb{R}^d)$  and thus in  $L^1(\Omega; \mathbb{R}^d)$ . It only remains to prove that the sequence is also relatively compact in  $L^q(\Omega; \mathbb{R}^d)$ , with  $1 < q < +\infty$  if  $d = 2$  or  $1 < q < 6$  if  $d = 3$ .

Owing to the discrete Sobolev embeddings [71, Proposition 5.4] together with the discrete Korn inequality (3.25), it holds, with  $r = q + 1$  if  $d = 2$  and  $r = 6$  if  $d = 3$ , that

$$\|\mathbf{v}_h\|_{L^r(\Omega; \mathbb{R}^d)} \lesssim \left( \sum_{T \in \mathcal{T}_h} \left[ \|\nabla \mathbf{v}_T\|_T^2 + \sum_{F \in \mathcal{F}_T} h_F^{-1} \|\mathbf{v}_F - \mathbf{v}_T\|_F^2 \right] \right)^{1/2},$$

Thus, we can complete the proof by means of the interpolation inequality [40, Remark 2 p. 93]. For all  $h, h' \in \mathcal{H}$  we have with  $\theta := \frac{r-q}{q(r-1)} \in (0, 1)$ ,

$$\|\mathbf{v}_h - \mathbf{v}_{h'}\|_{L^q(\Omega; \mathbb{R}^d)} \leq \|\mathbf{v}_h - \mathbf{v}_{h'}\|_{L^1(\Omega; \mathbb{R}^d)}^\theta \|\mathbf{v}_h - \mathbf{v}_{h'}\|_{L^r(\Omega; \mathbb{R}^d)}^{1-\theta} \lesssim \|\mathbf{v}_h - \mathbf{v}_{h'}\|_{L^1(\Omega; \mathbb{R}^d)}^\theta.$$

Therefore, up to a subsequence,  $(\mathbf{v}_h)_{h \in \mathcal{H}}$  is a Cauchy sequence in  $L^q(\Omega; \mathbb{R}^d)$ , so it converges.  $\square$

### 3.8.3 Consistency of the discrete symmetric gradient operator

*Proof of Proposition 3.13.* 1) *Strong consistency.* We first assume that  $\mathbf{v} \in H^2(\Omega; \mathbb{R}^d)$ . Owing to the commuting property (3.18) and the approximation property (3.13a) with  $m = 1$  and  $s = 2$ , it is inferred that  $\|\mathbf{G}_{s,T}^k \mathbf{I}_T^k \mathbf{v} - \nabla_s \mathbf{v}\|_T \lesssim h \|\mathbf{v}\|_{H^2(T; \mathbb{R}^d)}$ . Squaring, summing over  $T \in \mathcal{T}_h$ , and taking the square root of the resulting inequality gives

$$\|\mathbf{G}_{s,h}^k \mathbf{I}_h^k \mathbf{v} - \nabla_s \mathbf{v}\| \lesssim h \|\mathbf{v}\|_{H^2(\Omega; \mathbb{R}^d)}. \quad (3.61)$$

If  $\mathbf{v} \in H^1(\Omega; \mathbb{R}^d)$  we reason by density, namely we take a sequence  $(\mathbf{v}_\epsilon)_{\epsilon > 0} \subset H^2(\Omega; \mathbb{R}^d)$  that converges to  $\mathbf{v}$  in  $H^1(\Omega; \mathbb{R}^d)$  as  $\epsilon \rightarrow 0$  and, using twice the triangular inequality, we write

$$\|\mathbf{G}_{s,h}^k \mathbf{I}_h^k \mathbf{v} - \nabla_s \mathbf{v}\| \leq \|\mathbf{G}_{s,h}^k \mathbf{I}_h^k (\mathbf{v} - \mathbf{v}_\epsilon)\| + \|\mathbf{G}_{s,h}^k \mathbf{I}_h^k \mathbf{v}_\epsilon - \nabla_s \mathbf{v}_\epsilon\| + \|\nabla_s (\mathbf{v} - \mathbf{v}_\epsilon)\|.$$

By (3.61), the second term in the right-hand side tends to 0 as  $h \rightarrow 0$ . Moreover, owing to the commuting property (3.18) and the  $H^1$ -boundedness of  $\pi_T^k$ , one has

$$\begin{aligned} \|\mathbf{G}_{s,h}^k \mathbf{I}_h^k (\mathbf{v} - \mathbf{v}_\epsilon)\| &= \left( \sum_{T \in \mathcal{T}_h} \|\pi_T^k \nabla_s (\mathbf{v} - \mathbf{v}_\epsilon)\|_T^2 \right)^{1/2} \leq \left( \sum_{T \in \mathcal{T}_h} \|\nabla_s (\mathbf{v} - \mathbf{v}_\epsilon)\|_T^2 \right)^{1/2} \\ &\leq \|\nabla_s (\mathbf{v} - \mathbf{v}_\epsilon)\|. \end{aligned}$$

Therefore, taking the supremum limit as  $h \rightarrow 0$  and then the supremum limit as  $\epsilon \rightarrow 0$ , concludes the proof of (3.27) (notice that the order in which the limits are taken is important).

2) *Sequential consistency.* In order to prove (3.28) we observe that, by the definitions (3.19) of  $\mathbf{G}_{s,h}^k$  and (3.17b) of  $\mathbf{G}_{s,T}^k$  one has, for all  $\boldsymbol{\tau} \in H^1(\Omega; \mathbb{R}_{\text{sym}}^{d \times d})$  and all  $\underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_h^k$ ,



$$\begin{aligned}
\int_{\Omega} \mathbf{G}_{s,h}^k \underline{\mathbf{v}}_h : \boldsymbol{\tau} &= \sum_{T \in \mathcal{T}_h} \int_T \mathbf{G}_{s,T}^k \underline{\mathbf{v}}_T : \boldsymbol{\tau} \\
&= \sum_{T \in \mathcal{T}_h} \int_T (\mathbf{G}_{s,T}^k \underline{\mathbf{v}}_T - \nabla_s \mathbf{v}_T) : (\boldsymbol{\tau} - \boldsymbol{\pi}_T^0 \boldsymbol{\tau}) + \sum_{T \in \mathcal{T}_h} \int_T (\mathbf{G}_{s,T}^k \underline{\mathbf{v}}_T - \nabla_s \mathbf{v}_T) : \boldsymbol{\pi}_T^0 \boldsymbol{\tau} \\
&\quad + \sum_{T \in \mathcal{T}_h} \int_T \nabla_s \mathbf{v}_T : \boldsymbol{\tau} \\
&= \mathfrak{I}_1 + \sum_{T \in \mathcal{T}_h} \sum_{F \in \mathcal{F}_T} \int_F (\mathbf{v}_F - \mathbf{v}_T) \cdot (\boldsymbol{\pi}_T^0 \boldsymbol{\tau}) \mathbf{n}_{TF} + \sum_{T \in \mathcal{T}_h} \int_T \nabla_s \mathbf{v}_T : \boldsymbol{\tau} \\
&= \mathfrak{I}_1 + \sum_{T \in \mathcal{T}_h} \sum_{F \in \mathcal{F}_T} \int_F (\mathbf{v}_F - \mathbf{v}_T) \cdot (\boldsymbol{\pi}_T^0 \boldsymbol{\tau} - \boldsymbol{\tau}) \mathbf{n}_{TF} - \sum_{T \in \mathcal{T}_h} \int_T \mathbf{v}_T \cdot (\nabla \cdot \boldsymbol{\tau}) \\
&\quad + \sum_{F \in \mathcal{F}_h^b} \int_F \mathbf{v}_F \cdot (\boldsymbol{\tau} \mathbf{n}_{TF}) \\
&= \mathfrak{I}_1 + \mathfrak{I}_2 - \int_{\Omega} \mathbf{v}_h \cdot (\nabla \cdot \boldsymbol{\tau}) + \int_{\Gamma} \mathbf{v}_{\Gamma,h} \cdot \boldsymbol{\gamma}_n(\boldsymbol{\tau}).
\end{aligned} \tag{3.62}$$

To obtain the third equality, we used an element-wise integration by parts together with the relation

$$\sum_{T \in \mathcal{T}_h} \sum_{F \in \mathcal{F}_T \cap \mathcal{F}_h^i} \int_F \mathbf{v}_F \cdot (\boldsymbol{\tau} \mathbf{n}_{TF}) = \sum_{F \in \mathcal{F}_h^i} \int_F \mathbf{v}_F \cdot (\boldsymbol{\tau} \mathbf{n}_{T_1 F} + \boldsymbol{\tau} \mathbf{n}_{T_2 F}) = 0,$$

where for all  $F \in \mathcal{F}_h^i$ ,  $T_1, T_2 \in \mathcal{T}_h$  are such that  $F \subset \partial T_1 \cap \partial T_2$ . Owing to (3.62), the conclusion follows once we prove that  $|\mathfrak{I}_1 + \mathfrak{I}_2| \lesssim h \|\underline{\mathbf{v}}_h\|_{\epsilon,h} \|\boldsymbol{\tau}\|_{H^1(\Omega; \mathbb{R}^{d \times d})}$ . By (3.13a) (with  $m = 0$  and  $s = 1$ ) we have  $\|\boldsymbol{\tau} - \boldsymbol{\pi}_T^0 \boldsymbol{\tau}\|_T \lesssim h_T \|\boldsymbol{\tau}\|_{H^1(T; \mathbb{R}^{d \times d})}$  and thus, using the Cauchy–Schwarz and triangle inequalities followed by the norm equivalence (3.24),

$$\begin{aligned}
|\mathfrak{I}_1| &\leq \left( \sum_{T \in \mathcal{T}_h} \|\mathbf{G}_{s,T}^k \underline{\mathbf{v}}_T - \nabla_s \mathbf{v}_T\|_T^2 \right)^{1/2} \left( \sum_{T \in \mathcal{T}_h} \|\boldsymbol{\tau} - \boldsymbol{\pi}_T^0 \boldsymbol{\tau}\|_T^2 \right)^{1/2} \\
&\lesssim h \left( \|\mathbf{G}_{s,h}^k \underline{\mathbf{v}}_h\|^2 + \|\underline{\mathbf{v}}_h\|_{\epsilon,h}^2 \right)^{1/2} \|\boldsymbol{\tau}\|_{H^1(\Omega; \mathbb{R}^{d \times d})} \lesssim h \|\underline{\mathbf{v}}_h\|_{\epsilon,h} \|\boldsymbol{\tau}\|_{H^1(\Omega; \mathbb{R}^{d \times d})}.
\end{aligned} \tag{3.63}$$

In a similar way, we obtain an upper bound for  $\mathfrak{I}_2$ . By (3.13b) (with  $m = 0$  and  $s = 1$ ), for all  $F \in \mathcal{F}_T$ , we have  $\|\boldsymbol{\tau} - \boldsymbol{\pi}_T^0 \boldsymbol{\tau}\|_F \lesssim h_T^{1/2} \|\boldsymbol{\tau}\|_{H^1(T; \mathbb{R}^{d \times d})} \lesssim h_F^{1/2} \|\boldsymbol{\tau}\|_{H^1(T; \mathbb{R}^{d \times d})}$  and thus, using the Cauchy–Schwarz inequality,

$$|\mathfrak{I}_2| \lesssim \sum_{T \in \mathcal{T}_h} \sum_{F \in \mathcal{F}_T} h_F^{1/2} \|\mathbf{v}_F - \mathbf{v}_T\|_F \|\boldsymbol{\tau}\|_{H^1(T; \mathbb{R}^{d \times d})} \lesssim h \|\underline{\mathbf{v}}_h\|_{\epsilon,h} \|\boldsymbol{\tau}\|_{H^1(\Omega; \mathbb{R}^{d \times d})}. \tag{3.64}$$

Owing to (3.63) and (3.64), the triangle inequality  $|\mathfrak{I}_1 + \mathfrak{I}_2| \leq |\mathfrak{I}_1| + |\mathfrak{I}_2|$  yields the conclusion.  $\square$

Table 3.1: Convergence results on the triangular mesh family. OCV stands for order of convergence.

$h$	$\ \nabla_s \mathbf{u} - \mathbf{G}_{s,h}^k \mathbf{u}_h\ $	OCV	$\ \boldsymbol{\pi}_h^k \mathbf{u} - \mathbf{u}_h\ $	OCV
$k = 1$				
$3.07 \cdot 10^{-2}$	$5.59 \cdot 10^{-2}$	—	$7.32 \cdot 10^{-3}$	—
$1.54 \cdot 10^{-2}$	$1.51 \cdot 10^{-2}$	1.9	$1.05 \cdot 10^{-3}$	2.81
$7.68 \cdot 10^{-3}$	$3.86 \cdot 10^{-3}$	1.96	$1.34 \cdot 10^{-4}$	2.96
$3.84 \cdot 10^{-3}$	$1.01 \cdot 10^{-3}$	1.93	$1.7 \cdot 10^{-5}$	2.98
$1.92 \cdot 10^{-3}$	$2.59 \cdot 10^{-4}$	1.96	$2.15 \cdot 10^{-6}$	2.98
$k = 2$				
$3.07 \cdot 10^{-2}$	$1.3 \cdot 10^{-2}$	—	$1.47 \cdot 10^{-3}$	—
$1.54 \cdot 10^{-2}$	$1.29 \cdot 10^{-3}$	3.35	$6.05 \cdot 10^{-5}$	4.62
$7.68 \cdot 10^{-3}$	$2.11 \cdot 10^{-4}$	2.6	$5.36 \cdot 10^{-6}$	3.48
$3.84 \cdot 10^{-3}$	$2.73 \cdot 10^{-5}$	2.95	$3.6 \cdot 10^{-7}$	3.9
$1.92 \cdot 10^{-3}$	$3.42 \cdot 10^{-6}$	3	$2.28 \cdot 10^{-8}$	3.98
$k = 3$				
$3.07 \cdot 10^{-2}$	$2.81 \cdot 10^{-3}$	—	$2.39 \cdot 10^{-4}$	—
$1.54 \cdot 10^{-2}$	$3.72 \cdot 10^{-4}$	2.93	$1.95 \cdot 10^{-5}$	3.63
$7.68 \cdot 10^{-3}$	$2.16 \cdot 10^{-5}$	4.09	$5.47 \cdot 10^{-7}$	5.14
$3.84 \cdot 10^{-3}$	$1.43 \cdot 10^{-6}$	3.92	$1.66 \cdot 10^{-8}$	5.04
$1.92 \cdot 10^{-3}$	$9.51 \cdot 10^{-8}$	3.91	$5.34 \cdot 10^{-10}$	4.96
$k = 4$				
$3.07 \cdot 10^{-2}$	$1.37 \cdot 10^{-3}$	—	$1.13 \cdot 10^{-4}$	—
$1.54 \cdot 10^{-2}$	$5.97 \cdot 10^{-5}$	4.54	$3.04 \cdot 10^{-6}$	5.24
$7.68 \cdot 10^{-3}$	$1.76 \cdot 10^{-6}$	5.07	$4.09 \cdot 10^{-8}$	6.19
$3.84 \cdot 10^{-3}$	$6.46 \cdot 10^{-8}$	4.77	$7.64 \cdot 10^{-10}$	5.74

Table 3.2: Convergence results on the locally refined mesh family. OCV stands for order of convergence.

$h$	$\ \nabla_s \mathbf{u} - \mathbf{G}_{s,h}^k \underline{\mathbf{u}}_h\ $	OCV	$\ \boldsymbol{\pi}_h^k \mathbf{u} - \mathbf{u}_h\ $	OCV
$k = 1$				
0.25	0.13	—	$1.9 \cdot 10^{-2}$	—
0.13	$2.64 \cdot 10^{-2}$	2.28	$2.54 \cdot 10^{-3}$	2.9
$6.25 \cdot 10^{-2}$	$4.97 \cdot 10^{-3}$	2.41	$3.22 \cdot 10^{-4}$	2.98
$3.12 \cdot 10^{-2}$	$9.14 \cdot 10^{-4}$	2.44	$4.12 \cdot 10^{-5}$	2.96
$1.56 \cdot 10^{-2}$	$1.67 \cdot 10^{-4}$	2.45	$5.21 \cdot 10^{-6}$	2.98
$k = 2$				
0.25	$1.88 \cdot 10^{-2}$	—	$3.79 \cdot 10^{-3}$	—
0.13	$5.05 \cdot 10^{-3}$	1.9	$3.55 \cdot 10^{-4}$	3.42
$6.25 \cdot 10^{-2}$	$6.51 \cdot 10^{-4}$	2.96	$2.92 \cdot 10^{-5}$	3.6
$3.12 \cdot 10^{-2}$	$6.83 \cdot 10^{-5}$	3.25	$1.89 \cdot 10^{-6}$	3.94
$1.56 \cdot 10^{-2}$	$6.23 \cdot 10^{-6}$	3.45	$1.19 \cdot 10^{-7}$	3.99
$k = 3$				
0.25	$7.84 \cdot 10^{-3}$	—	$1.41 \cdot 10^{-3}$	—
0.13	$1.09 \cdot 10^{-3}$	2.85	$7.5 \cdot 10^{-5}$	4.23
$6.25 \cdot 10^{-2}$	$8.22 \cdot 10^{-5}$	3.73	$3.93 \cdot 10^{-6}$	4.25
$3.12 \cdot 10^{-2}$	$5.64 \cdot 10^{-6}$	3.86	$1.45 \cdot 10^{-7}$	4.75
$1.56 \cdot 10^{-2}$	$3.44 \cdot 10^{-7}$	4.04	$5.23 \cdot 10^{-9}$	4.79
$k = 4$				
0.25	$4.35 \cdot 10^{-3}$	—	$4.68 \cdot 10^{-4}$	—
0.13	$3.65 \cdot 10^{-4}$	3.58	$3.19 \cdot 10^{-5}$	3.87
$6.25 \cdot 10^{-2}$	$1.5 \cdot 10^{-5}$	4.6	$6.02 \cdot 10^{-7}$	5.73
$3.12 \cdot 10^{-2}$	$5.78 \cdot 10^{-7}$	4.69	$1.03 \cdot 10^{-8}$	5.86

Table 3.3: Convergence results on the hexagonal mesh family. OCV stands for order of convergence.

$h$	$\ \nabla_s \mathbf{u} - \mathbf{G}_{s,h}^k \mathbf{u}_h\ $	OCV	$\ \boldsymbol{\pi}_h^k \mathbf{u} - \mathbf{u}_h\ $	OCV
$k = 1$				
$6.3 \cdot 10^{-2}$	0.22	—	$2.75 \cdot 10^{-2}$	—
$3.42 \cdot 10^{-2}$	$3.72 \cdot 10^{-2}$	2.89	$3.73 \cdot 10^{-3}$	3.27
$1.72 \cdot 10^{-2}$	$7.17 \cdot 10^{-3}$	2.4	$4.83 \cdot 10^{-4}$	2.97
$8.59 \cdot 10^{-3}$	$1.44 \cdot 10^{-3}$	2.31	$6.14 \cdot 10^{-5}$	2.97
$4.3 \cdot 10^{-3}$	$2.4 \cdot 10^{-4}$	2.59	$7.7 \cdot 10^{-6}$	3
$k = 2$				
$6.3 \cdot 10^{-2}$	$2.68 \cdot 10^{-2}$	—	$3.04 \cdot 10^{-3}$	—
$3.42 \cdot 10^{-2}$	$7.01 \cdot 10^{-3}$	2.2	$3.56 \cdot 10^{-4}$	3.51
$1.72 \cdot 10^{-2}$	$1.09 \cdot 10^{-3}$	2.71	$3.31 \cdot 10^{-5}$	3.46
$8.59 \cdot 10^{-3}$	$1.41 \cdot 10^{-4}$	2.95	$2.53 \cdot 10^{-6}$	3.7
$4.3 \cdot 10^{-3}$	$1.96 \cdot 10^{-5}$	2.85	$1.72 \cdot 10^{-7}$	3.89
$k = 3$				
$6.3 \cdot 10^{-2}$	$1.11 \cdot 10^{-2}$	—	$1.08 \cdot 10^{-3}$	—
$3.42 \cdot 10^{-2}$	$1.92 \cdot 10^{-3}$	2.87	$9.29 \cdot 10^{-5}$	4.02
$1.72 \cdot 10^{-2}$	$2.79 \cdot 10^{-4}$	2.81	$6.13 \cdot 10^{-6}$	3.95
$8.59 \cdot 10^{-3}$	$2.54 \cdot 10^{-5}$	3.45	$2.88 \cdot 10^{-7}$	4.4
$4.3 \cdot 10^{-3}$	$1.61 \cdot 10^{-6}$	3.99	$1.24 \cdot 10^{-8}$	4.55
$k = 4$				
$6.3 \cdot 10^{-2}$	$5.53 \cdot 10^{-3}$	—	$4.49 \cdot 10^{-4}$	—
$3.42 \cdot 10^{-2}$	$5.76 \cdot 10^{-4}$	3.7	$3.07 \cdot 10^{-5}$	4.39
$1.72 \cdot 10^{-2}$	$6.29 \cdot 10^{-5}$	3.22	$1.21 \cdot 10^{-6}$	4.7
$8.59 \cdot 10^{-3}$	$2.21 \cdot 10^{-6}$	4.82	$2.69 \cdot 10^{-8}$	5.48

Table 3.4: Convergence results on the Voronoi mesh family. OCV stands for order of convergence.

$h$	$\ \nabla_s \mathbf{u} - \mathbf{G}_{s,h}^k \mathbf{u}_h\ $	OCV	$\ \boldsymbol{\pi}_h^k \mathbf{u} - \mathbf{u}_h\ $	OCV
$k = 1$				
$6.5 \cdot 10^{-2}$	$8.82 \cdot 10^{-2}$	—	$1.55 \cdot 10^{-2}$	—
$3.15 \cdot 10^{-2}$	$1.49 \cdot 10^{-2}$	2.45	$2.29 \cdot 10^{-3}$	2.64
$1.61 \cdot 10^{-2}$	$3.63 \cdot 10^{-3}$	2.1	$3.01 \cdot 10^{-4}$	3.02
$9.09 \cdot 10^{-3}$	$8.68 \cdot 10^{-4}$	2.5	$3.95 \cdot 10^{-5}$	3.55
$4.26 \cdot 10^{-3}$	$2.04 \cdot 10^{-4}$	1.91	$4.97 \cdot 10^{-6}$	2.74
$k = 2$				
$6.5 \cdot 10^{-2}$	$1.43 \cdot 10^{-2}$	—	$2.63 \cdot 10^{-3}$	—
$3.15 \cdot 10^{-2}$	$4.03 \cdot 10^{-3}$	1.75	$2.53 \cdot 10^{-4}$	3.23
$1.61 \cdot 10^{-2}$	$4.78 \cdot 10^{-4}$	3.18	$2.22 \cdot 10^{-5}$	3.63
$9.09 \cdot 10^{-3}$	$6.7 \cdot 10^{-5}$	3.44	$1.45 \cdot 10^{-6}$	4.77
$4.26 \cdot 10^{-3}$	$9.08 \cdot 10^{-6}$	2.64	$9.07 \cdot 10^{-8}$	3.66
$k = 3$				
$6.5 \cdot 10^{-2}$	$7.12 \cdot 10^{-3}$	—	$9.08 \cdot 10^{-4}$	—
$3.15 \cdot 10^{-2}$	$8.34 \cdot 10^{-4}$	2.96	$6.78 \cdot 10^{-5}$	3.58
$1.61 \cdot 10^{-2}$	$7.03 \cdot 10^{-5}$	3.69	$3.18 \cdot 10^{-6}$	4.56
$9.09 \cdot 10^{-3}$	$4.17 \cdot 10^{-6}$	4.94	$9.67 \cdot 10^{-8}$	6.11
$4.26 \cdot 10^{-3}$	$2.42 \cdot 10^{-7}$	3.76	$3.15 \cdot 10^{-9}$	4.52
$k = 4$				
$6.5 \cdot 10^{-2}$	$3.25 \cdot 10^{-3}$	—	$3.68 \cdot 10^{-4}$	—
$3.15 \cdot 10^{-2}$	$2.94 \cdot 10^{-4}$	3.32	$2.14 \cdot 10^{-5}$	3.93
$1.61 \cdot 10^{-2}$	$9.86 \cdot 10^{-6}$	5.06	$4.34 \cdot 10^{-7}$	5.81
$9.09 \cdot 10^{-3}$	$3.47 \cdot 10^{-7}$	5.85	$6.74 \cdot 10^{-9}$	7.29

# Chapter 4

---

## Nonlinear poroelasticity

---

This chapter has been submitted for publication. The preprint is available in HAL:

**Analysis of a Hybrid High-Order–discontinuous Galerkin discretization  
method for nonlinear poroelasticity,**  
<https://hal.archives-ouvertes.fr/hal-01785810>, 2018.

### Contents

---

<b>4.1</b>	<b>Introduction</b>	<b>84</b>
<b>4.2</b>	<b>Continuous setting</b>	<b>86</b>
4.2.1	Notation for function spaces	86
4.2.2	Stress-strain law	87
4.2.3	Weak formulation	88
<b>4.3</b>	<b>Discrete setting</b>	<b>89</b>
4.3.1	Space mesh	89
4.3.2	Time mesh	89
4.3.3	$L^2$ -orthogonal projectors on local and broken polynomial spaces	90
4.3.4	Discrete spaces	91
4.3.4.1	Displacement	91
4.3.4.2	Pore pressure	93
<b>4.4</b>	<b>Discretization</b>	<b>93</b>
4.4.1	Nonlinear elasticity operator	93
4.4.2	Hydro-mechanical coupling	95
4.4.3	Darcy operator	97
4.4.4	Discrete problem	97
<b>4.5</b>	<b>Stability and well-posedness</b>	<b>98</b>
<b>4.6</b>	<b>Convergence analysis</b>	<b>103</b>

---

## 4.1 Introduction

In this chapter, we formulate and analyze a coupled Hybrid High-Order–discontinuous Galerkin (HHO–dG) discretization method for nonlinear poroelastic problems. We overstep the work of Chapter 2 devoted to the Biot’s model [27, 190] by considering more general nonlinear stress-strain constitutive laws as in Chapter 3. The model is valid under the assumptions of small deformations of the rock matrix, small variations of the porosity, and small relative variations of the fluid density. The interest of the poroelastic models considered here is particularly manifest in geosciences applications [123, 127, 151], where fluid flows in geological subsurface, modeled as a porous media, induce a deformation of the rock matrix. The challenge is then to design a discretization method able to (i) treat a complex geometry with polyhedral meshes and nonconforming interfaces, (ii) handle possible heterogeneities of the poromechanical parameters and nonlinearities of the stress-strain relation, and (iii) deal with the numerical instabilities encountered in this type of coupled problem. We focus here on the theoretical aspects while numerical tests assessing the convergence of the method have been presented in Section A.5; further numerical investigation will be carried out in a future work.

Let  $\Omega \subset \mathbb{R}^d$ ,  $d \in \{2, 3\}$ , denote a bounded connected polyhedral domain with boundary  $\partial\Omega$  and outward normal  $\mathbf{n}$ . Without loss of generality, we assume that the domain is scaled so that its diameter  $\text{diam}(\Omega)$  is equal to 1. For a given finite time  $t_F > 0$ , volumetric load  $\mathbf{f}$ , fluid source  $g$ , the nonlinear poroelasticity problem considered here consists in finding a vector-valued displacement field  $\mathbf{u}$  and a scalar-valued pore pressure field  $p$  solution of

$$-\nabla \cdot \boldsymbol{\sigma}(\cdot, \nabla_s \mathbf{u}) + \alpha \nabla p = \mathbf{f} \quad \text{in } \Omega \times (0, t_F), \quad (4.1a)$$

$$C_0 d_t p + \alpha d_t \nabla \cdot \mathbf{u} - \nabla \cdot (\boldsymbol{\kappa}(\cdot) \nabla p) = g \quad \text{in } \Omega \times (0, t_F), \quad (4.1b)$$

where  $\nabla_s$  denotes the symmetric gradient,  $d_t$  denotes the time derivative,  $\alpha$  is the Biot–Willis coefficient,  $C_0 \geq 0$  is the constrained specific storage coefficient, and  $\boldsymbol{\kappa} : \Omega \rightarrow \mathbb{R}_s^{d \times d}$  is the uniformly elliptic permeability tensor field which, for real numbers  $0 < \underline{\kappa} \leq \bar{\kappa}$ , satisfies for almost every (a.e.)  $\mathbf{x} \in \Omega$  and all  $\boldsymbol{\xi} \in \mathbb{R}^d$ ,

$$\underline{\kappa} |\boldsymbol{\xi}|^2 \leq \boldsymbol{\kappa}(\mathbf{x}) \boldsymbol{\xi} \cdot \boldsymbol{\xi} \leq \bar{\kappa} |\boldsymbol{\xi}|^2.$$

For the sake of simplicity, we assume in the following discussion that  $\boldsymbol{\kappa}$  is piecewise constant on a polyhedral partition  $P_\Omega$  of  $\Omega$ , an assumption typically verified in geoscience applications. In the poroelasticity theory [61], the medium is modeled as a continuous superposition of solid and fluid phases. The momentum equilibrium equation (4.1a) is based on the Terzaghi decomposition [190] of the total stress tensor into a mechanical contribution and a pore pressure contribution. Examples and assumptions for the constitutive stress-strain relation  $\boldsymbol{\sigma} : \Omega \times \mathbb{R}_s^{d \times d} \rightarrow \mathbb{R}_s^{d \times d}$  are detailed in Section 4.2.2; we refer the reader to [22, 28] for a physical and experimental investigation of the nonlinear behavior of porous solids. On the other hand, the mass conservation equation (4.1b) is derived for fully saturated porous media assuming Darcean flow. The first two terms of this equation quantify the variation of fluid content in the pores. The dimensionless coupling coefficient  $\alpha$  express the amount of fluid that can be forced into the medium by a variation of pore volume for a constant

fluid pressure, while  $C_0$  measures the amount of fluid that can be forced into the medium by pressure increments due to compressibility of the structure. The case of a solid matrix with incompressible grains corresponds to the limit value  $C_0 = 0$ . Following [181, 205], for the sake of simplicity we take  $\alpha = 1$  in what follows. To close the problem, we enforce homogeneous boundary conditions corresponding to a clamped, impermeable boundary, i.e.,

$$\mathbf{u} = \mathbf{0} \quad \text{on } \partial\Omega \times (0, t_F), \quad (4.1c)$$

$$(\kappa(\cdot)\nabla p) \cdot \mathbf{n} = 0 \quad \text{on } \partial\Omega \times (0, t_F), \quad (4.1d)$$

as well as the following initial condition which prescribes the initial fluid content:

$$C_0 p(\cdot, 0) + \nabla \cdot \mathbf{u}(\cdot, 0) = \phi^0(\cdot). \quad (4.1e)$$

In the case  $C_0 = 0$ , we also need the following compatibility conditions on  $g$  and  $\phi^0$  and zero-average constraint on  $p$ :

$$\int_{\Omega} \phi^0 = 0, \quad \int_{\Omega} g(\cdot, t) = 0, \quad \text{and} \quad \int_{\Omega} p(\cdot, t) = 0 \quad \forall t \in (0, t_F). \quad (4.1f)$$

When discretizing the poroelasticity system (4.1), the main challenges are to ensure stability and convergence under mild assumptions on the nonlinear stress-strain relation and on the permeability field, and to prevent localized pressure oscillations arising in the case of low-permeable and low-compressible porous media. Since the latter issue is in part related to the saddle point structure in the coupled equations for  $C_0 = 0$  and small  $\underline{\kappa}$ , the discrete spaces for the displacement and the pressure should satisfy an inf-sup condition. Indeed, as observed in [117, 155, 168] in the context of finite element discretizations of the linear poroelasticity problem, the inf-sup condition yields an  $L^2$ -estimate of the discrete pressure independent of  $\underline{\kappa}^{-1}$ , and allows one to prove the convergence of the approximate pressure towards the continuous pressure also in the incompressible case  $C_0 = 0$ . We notice, however, that the problem of spurious pressure oscillations is actually more involved than a simple saddle-point coupling issue. For instance, it has been recently pointed out in [176] that, even for discretization methods leading to an inf-sup stable discretization of the Stokes problem in the steady case, pressure oscillations can arise owing to a lack of monotonicity of the discrete operator. The robustness with respect to spurious oscillations has been numerically observed in Section 2.6.2 for a HHO–dG discretization of the linear poroelasticity model.

In this chapter, we present and analyze a nonconforming space discretization of problem (4.1) where the nonlinear elasticity operator is discretized using the HHO method of Chapter 3 (c.f. also [71, 74]), while the Darcy operator relies on the Symmetric Weighted Interior Penalty (SWIP) method of [77]. The proposed method has several assets: (i) it is valid in two and three space dimensions; (ii) it delivers an inf-sup stable discretization on general spatial meshes including, e.g., polyhedral elements and nonmatching interfaces; (iii) it allows one to increase the space approximation order to accelerate convergence in the presence of (locally) regular solutions. Compared to the method proposed in Chapter 2 for the linear poroelasticity problem, there are two main differences in the design. First, for a given polynomial degree  $k \geq 1$ , the symmetric gradient reconstruction sits in the full space



of tensor-valued polynomials of total degree  $\leq k$ , as opposed to symmetric gradients of vector-valued polynomials of total degree  $\leq (k + 1)$ . Following [33, 72], this modification is required to obtain optimal convergence rates when considering nonlinear stress-strain laws. Second, the right-hand side of the discrete problem of Section 4.4.4 is obtained by taking the average in time of the loading force  $\mathbf{f}$  and fluid source  $g$  over a time step instead of their value at the end of the time step. This modification allows us to prove stability and optimal error estimates under significantly weaker time regularity assumptions on data (cf. Remark 4.12 and 4.18). Finally, in Section 4.3.4 we give a new simple proof of discrete Korn's inequality, not requiring particular geometrical assumption on the mesh. The interest of these results goes beyond the specific application considered here.

The material is organized as follows. In Section 4.2 we present the assumptions on the stress-strain law and the variational formulation of the nonlinear poroelasticity problem. In Section 4.3 we define the space and time meshes and the discrete spaces for the displacement and the pressure fields. In Section 4.4 we define the discrete counterparts of the elasticity, Darcy, and hydromechanical coupling operators and formulate the discrete problem. In Section 4.5 we prove the well-posedness of the scheme by deriving an a priori estimate on the discrete solution that holds also when the specific storage coefficient vanishes. The convergence analysis of the method is carried out in Section 4.6.

## 4.2 Continuous setting

In this section we introduce the notation for function spaces, formulate the assumptions on the stress-strain law, and derive a weak formulation of problem (4.1).

### 4.2.1 Notation for function spaces

Let  $X \subset \bar{\Omega}$ . Spaces of functions, vector fields, and tensor fields defined over  $X$  are respectively denoted by italic capital, boldface Roman capital, and special Roman capital letters. The subscript "s" appended to a special Roman capital letter denotes a space of symmetric tensor fields. Thus, for example,  $L^2(X)$ ,  $\mathbf{L}^2(X)$ , and  $\mathbb{L}_s^2(X)$  denote the spaces of square integrable functions, vector fields, and symmetric tensor fields over  $X$ , respectively. For any measured set  $X$  and any  $m \in \mathbb{Z}$ , we denote by  $H^m(X)$  the usual Sobolev space of functions that have weak partial derivatives of order up to  $m$  in  $L^2(X)$ , with the convention that  $H^0(X) := L^2(X)$ , while  $C^m(X)$  and  $C_c^\infty(X)$  denote, respectively, the usual spaces of  $m$ -times continuously differentiable functions and infinitely continuously differentiable functions with compact support on  $X$ . We denote by  $(\cdot, \cdot)_X$  and  $(\cdot, \cdot)_{m,X}$  the usual scalar products in  $L^2(X)$  and  $H^m(X)$  respectively, and by  $\|\cdot\|_X$  and  $\|\cdot\|_{m,X}$  the induced norms.

For a vector space  $V$  with scalar product  $(\cdot, \cdot)_V$ , the space  $C^m(V) := C^m([0, t_F]; V)$  is spanned by  $V$ -valued functions that are  $m$ -times continuously differentiable in the time interval  $[0, t_F]$ . The space  $C^m(V)$  is a Banach space when equipped with the norm

$$\|\varphi\|_{C^m(V)} := \max_{0 \leq i \leq m} \max_{t \in [0, t_F]} \|d_t^i \varphi(t)\|_V.$$

Similarly, the Hilbert space  $H^m(V) := H^m((0, t_F); V)$  is spanned by  $V$ -valued functions of the time interval, and the norm  $\|\cdot\|_{H^m(V)}$  is induced by the scalar product

$$(\varphi, \psi)_{H^m(V)} = \sum_{j=0}^m \int_0^{t_F} (d_t^j \varphi(t), d_t^j \psi(t))_V dt \quad \forall \varphi, \psi \in H^m(V).$$

### 4.2.2 Stress-strain law

The following assumptions on the stress-strain relation are required to obtain a well-posed weak formulation of the nonlinear poroelasticity problem.

**Assumption 4.1** (Stress-strain relation). We assume that the stress function  $\sigma : \Omega \times \mathbb{R}_s^{d \times d} \rightarrow \mathbb{R}_s^{d \times d}$  is a Carathéodory function, i.e.,  $\sigma(\mathbf{x}, \cdot)$  is continuous on  $\mathbb{R}_s^{d \times d}$  for almost every  $\mathbf{x} \in \Omega$  and  $\sigma(\cdot, \boldsymbol{\tau})$  is measurable on  $\Omega$  for all  $\boldsymbol{\tau} \in \mathbb{R}_s^{d \times d}$ . Moreover, there exist real numbers  $C_{gr}, C_{cv} \in (0, +\infty)$  such that, for a.e.  $\mathbf{x} \in \Omega$  and all  $\boldsymbol{\tau}, \boldsymbol{\eta} \in \mathbb{R}_s^{d \times d}$ , the following conditions hold:

$$|\sigma(\mathbf{x}, \boldsymbol{\tau})|_{d \times d} \leq C_{gr} |\boldsymbol{\tau}|_{d \times d}, \quad (\text{growth}) \quad (4.2a)$$

$$\sigma(\mathbf{x}, \boldsymbol{\tau}) : \boldsymbol{\tau} \geq C_{cv}^2 |\boldsymbol{\tau}|_{d \times d}^2, \quad (\text{coercivity}) \quad (4.2b)$$

$$(\sigma(\mathbf{x}, \boldsymbol{\tau}) - \sigma(\mathbf{x}, \boldsymbol{\eta})) : (\boldsymbol{\tau} - \boldsymbol{\eta}) > 0 \text{ if } \boldsymbol{\eta} \neq \boldsymbol{\tau}. \quad (\text{monotonicity}) \quad (4.2c)$$

Above, we have introduced the Frobenius product such that, for all  $\boldsymbol{\tau}, \boldsymbol{\eta} \in \mathbb{R}^{d \times d}$ ,  $\boldsymbol{\tau} : \boldsymbol{\eta} := \sum_{1 \leq i, j \leq d} \tau_{ij} \eta_{ij}$  with corresponding matrix norm such that, for all  $\boldsymbol{\tau} \in \mathbb{R}^{d \times d}$ ,  $|\boldsymbol{\tau}|_{d \times d} := (\boldsymbol{\tau} : \boldsymbol{\tau})^{\frac{1}{2}}$ .

Three meaningful examples for the stress-strain relation  $\sigma : \Omega \times \mathbb{R}_s^{d \times d} \rightarrow \mathbb{R}_s^{d \times d}$  in (4.1a) are:

- The (possibly heterogeneous) *linear elasticity model* given by the usual Hooke's law

$$\sigma(\cdot, \boldsymbol{\tau}) = \lambda(\cdot) \operatorname{tr}(\boldsymbol{\tau}) \mathbf{I}_d + 2\mu(\cdot) \boldsymbol{\tau}, \quad (4.3)$$

where  $\mu : \Omega \rightarrow [\mu_*, \mu^*]$ , with  $0 < \mu_* \leq \mu^* < +\infty$ , and  $\lambda : \Omega \rightarrow \mathbb{R}_+$  are the Lamé parameters.

- The *nonlinear Hencky–Mises model* of [105, 157] corresponding to the mechanical behavior law

$$\sigma(\cdot, \boldsymbol{\tau}) = \tilde{\lambda}(\cdot, \operatorname{dev}(\boldsymbol{\tau})) \operatorname{tr}(\boldsymbol{\tau}) \mathbf{I}_d + 2\tilde{\mu}(\cdot, \operatorname{dev}(\boldsymbol{\tau})) \boldsymbol{\tau}, \quad (4.4)$$

with nonlinear Lamé scalar functions  $\tilde{\mu} : \Omega \times \mathbb{R}_+ \rightarrow [\mu_*, \mu^*]$  and  $\tilde{\lambda} : \Omega \times \mathbb{R}_+ \rightarrow \mathbb{R}_+$  depending on the deviatoric part of the strain  $\operatorname{dev}(\boldsymbol{\tau}) := \operatorname{tr}(\boldsymbol{\tau}^2) - \frac{1}{d} \operatorname{tr}(\boldsymbol{\tau})^2$ .

- The *isotropic reversible hyperelastic damage model* [50], for which the stress-strain relation reads

$$\sigma(\cdot, \boldsymbol{\tau}) = (1 - D(\cdot, \boldsymbol{\tau})) \mathfrak{C}(\cdot) \boldsymbol{\tau}. \quad (4.5)$$

where  $D : \Omega \times \mathbb{R}_s^{d \times d} \rightarrow [0, 1]$  is the scalar damage function and  $\mathfrak{C} : \Omega \rightarrow \mathbb{R}^{d^4}$  is a fourth-order symmetric and uniformly elliptic tensor field, namely, for some strictly positive constants  $\underline{C}$  and  $\overline{C}$ , it holds

$$\underline{C} |\boldsymbol{\tau}|_{d \times d}^2 \leq \mathfrak{C}(\mathbf{x}) \boldsymbol{\tau} : \boldsymbol{\tau} \leq \overline{C} |\boldsymbol{\tau}|_{d \times d}^2 \quad \forall \boldsymbol{\tau} \in \mathbb{R}^{d \times d}, \forall \mathbf{x} \in \Omega.$$

Being linear, the Cauchy stress tensor in (4.3) clearly satisfies the previous assumptions. Moreover, under some mild requirements (cf. Appendix B and [87]) on the nonlinear Lamé scalar functions  $\tilde{\mu}$  and  $\tilde{\lambda}$  in (4.4) and on the damage function  $D$  in (4.5), it can be proven that also the Hencky–Mises model and the isotropic reversible damage model satisfy Assumption 4.1.

### 4.2.3 Weak formulation

At each time  $t \in [0, t_F]$ , the natural functional spaces for the displacement  $\mathbf{u}(t)$  and pore pressure  $p(t)$  taking into account the boundary condition (4.1c) and the zero average constraint (4.1f) are, respectively,

$$U := \mathbf{H}_0^1(\Omega) \quad \text{and} \quad P := \begin{cases} H^1(\Omega) & \text{if } C_0 > 0, \\ H^1(\Omega) \cap L_0^2(\Omega) & \text{if } C_0 = 0, \end{cases}$$

with  $\mathbf{H}_0^1(\Omega) := \{\mathbf{v} \in \mathbf{H}^1(\Omega) : \mathbf{v}|_{\partial\Omega} = \mathbf{0}\}$  and  $L_0^2(\Omega) := \{q \in L^2(\Omega) : \int_{\Omega} q = 0\}$ . We consider the following weak formulation of problem (4.1): For a loading term  $\mathbf{f} \in L^2(\mathbf{L}^2(\Omega))$ , a fluid source  $g \in L^2(L^2(\Omega))$ , and an initial datum  $\phi^0 \in L^2(\Omega)$  that verify (4.1f), find  $\mathbf{u} \in L^2(U)$  and  $p \in L^2(P)$  such that, for all  $\mathbf{v} \in U$ , all  $q \in P$ , and all  $\varphi \in C_c^\infty((0, t_F))$

$$\int_0^{t_F} a(\mathbf{u}(t), \mathbf{v}) \varphi(t) dt + \int_0^{t_F} b(\mathbf{v}, p(t)) \varphi(t) dt = \int_0^{t_F} (\mathbf{f}(t), \mathbf{v})_{\Omega} \varphi(t) dt, \quad (4.6a)$$

$$\int_0^{t_F} [b(\mathbf{u}(t), q) - C_0(p(t), q)_{\Omega}] d_t \varphi(t) dt + \int_0^{t_F} c(p(t), q) \varphi(t) dt = \int_0^{t_F} (g(t), q)_{\Omega} \varphi(t) dt, \quad (4.6b)$$

$$(C_0 p(0) + \nabla \cdot \mathbf{u}(0), q)_{\Omega} = (\phi^0, q)_{\Omega}, \quad (4.6c)$$

where we have defined the nonlinear function  $a : U \times U \rightarrow \mathbb{R}$  and the bilinear forms  $b : U \times P \rightarrow \mathbb{R}$  and  $c : P \times P \rightarrow \mathbb{R}$  such that, for all  $\mathbf{v}, \mathbf{w} \in U$  and all  $q, r \in P$ ,

$$a(\mathbf{v}, \mathbf{w}) := (\boldsymbol{\sigma}(\cdot, \nabla_s \mathbf{v}), \nabla_s \mathbf{w})_{\Omega}, \quad b(\mathbf{v}, q) := -(\nabla \cdot \mathbf{v}, q)_{\Omega}, \quad c(q, r) := (\kappa(\cdot) \nabla r, \nabla q)_{\Omega}.$$

The first term in (4.6a) is well defined thanks to the growth assumption (4.2a). Moreover, owing to (4.2b) together with Korn's first inequality and Poincaré's inequality,  $a(\cdot, \cdot)$  and  $c(\cdot, \cdot)$  are coercive on  $U$  and  $P$ , respectively. The strict monotonicity assumption (4.2c) guarantees the uniqueness of the weak solution.

*Remark 4.2* (Regularity of the fluid content and of the pore pressure). Using an integration by parts in time in (4.6b), it is inferred that

$$d_t [C_0(p, q)_{\Omega} - b(\mathbf{u}, q)] + c(p, q) = (g, q)_{\Omega} \quad \forall q \in P \quad \text{in } L^2((0, t_F)). \quad (4.7)$$

Therefore, defining the fluid content  $\phi := C_0 p + \nabla \cdot \mathbf{u}$ , we have that  $t \mapsto (\phi(t), q)_{\Omega} \in H^1((0, t_F)) \subset C^0([0, t_F])$  for all  $q \in P$ , and, as a result, (4.6c) makes sense. Moreover, in the

case  $C_0 > 0$ , taking  $q = 1$  in (4.7) and owing to the definition of the bilinear form  $c$  and the homogeneous Dirichlet condition (4.1c), we infer that

$$d_t \left( C_0 \int_{\Omega} p(\cdot, t) \right) = \int_{\Omega} g(\cdot, t) \quad \text{in } L^2((0, t_F)).$$

Thus,  $t \mapsto \int_{\Omega} p(\cdot, t) \in H^1((0, t_F))$ , namely the average of the pore pressure on the medium  $\Omega$  is a continuous function in  $[0, t_F]$ .

## 4.3 Discrete setting

In this section we define the space and time meshes, recall the definition and properties of  $L^2$ -orthogonal projectors on local and broken polynomial spaces, and introduce the discrete spaces for the displacement and the pressure.

### 4.3.1 Space mesh

We consider here polygonal or polyhedral meshes corresponding to couples  $\mathcal{M}_h := (\mathcal{T}_h, \mathcal{F}_h)$ , where  $\mathcal{T}_h$  is a finite collection of polygonal elements  $T$  such that  $h := \max_{T \in \mathcal{T}_h} h_T > 0$  with  $h_T$  denoting the diameter of  $T$ , while  $\mathcal{F}_h$  is a finite collection of hyperplanar faces  $F$ . It is assumed henceforth that the mesh  $\mathcal{M}_h$  matches the geometrical requirements detailed in [86, Definition 7.2]; see also [82, Section 2]. We additionally assume here that the mesh elements are star-shaped with respect to a ball of diameter uniformly comparable to the element diameter in order to use the results of [32, Appendix A]. To avoid dealing with jumps of the permeability coefficient inside elements, we additionally assume that  $\mathcal{M}_h$  is compliant with the partition  $P_{\Omega}$  on which  $\kappa$  is piecewise constant meaning that, for every  $T \in \mathcal{T}_h$ , there exists a unique subdomain  $\omega \in P_{\Omega}$  such that  $T \subset \omega$ . For every mesh element  $T \in \mathcal{T}_h$ , we denote by  $\mathcal{F}_T$  the subset of  $\mathcal{F}_h$  containing the faces that lie on the boundary  $\partial T$  of  $T$ . For each face  $F \in \mathcal{F}_T$ ,  $\mathbf{n}_{TF}$  is the (constant) unit normal vector to  $F$  pointing out of  $T$ . Boundary faces lying on  $\partial\Omega$  and internal faces contained in  $\Omega$  are collected in the sets  $\mathcal{F}_h^b$  and  $\mathcal{F}_h^i$ , respectively.

Our focus is on the so-called  $h$ -convergence analysis, so we consider a sequence of refined meshes that is regular in the sense of [82, Definition 3]. The corresponding strictly positive regularity parameter, uniformly bounded away from zero, is denoted by  $\varrho$ . We additionally assume that, for each mesh in the sequence, all the elements are star-shaped with respect to every point of a ball of radius uniformly comparable to the diameter of the element. The mesh regularity assumption implies, in particular, that the diameter  $h_T$  of a mesh element  $T \in \mathcal{T}_h$  is uniformly comparable to the diameter  $h_F$  of each face  $F \in \mathcal{F}_T$ , and that the number of faces in  $\mathcal{F}_T$  is bounded above by an integer  $N_{\partial}$  independent of  $h$ .

### 4.3.2 Time mesh

We subdivide  $(0, t_F)$  into  $N \in \mathbb{N}^*$  uniform subintervals, and introduce the timestep  $\tau := t_F/N$  and the discrete times  $t^n := n\tau$  for all  $0 \leq n \leq N$ . We define the space of piecewise  $H^1$

functions on  $(0, t_F)$  by

$$H^1(\mathcal{T}_\tau) := \left\{ \varphi \in L^2((0, t_F)) : \varphi|_{(t^{n-1}, t^n)} \in H^1((t^{n-1}, t^n)) \text{ for all } 1 \leq n \leq N \right\}.$$

Since each  $\psi \in H^1((t^{n-1}, t^n))$  has an absolutely continuous representative in  $[t^{n-1}, t^n]$ , we can identify  $\varphi \in H^1(\mathcal{T}_\tau)$  with a left continuous function in  $(0, t_F)$ . Therefore, for any vector space  $V$  and any  $\varphi \in H^1(\mathcal{T}_\tau; V)$ , we set  $\varphi^0 := \varphi(0)$  and, for all  $1 \leq n \leq N$ ,

$$\varphi^n := \lim_{t \rightarrow (t^n)^-} \varphi(t) \in V.$$

If  $\varphi \in C^0(V)$ , this simply amounts to setting  $\varphi^n := \varphi(t^n)$ . For all  $n \geq 1$  and  $\psi \in L^1(V)$ , we define the time average of  $\psi$  in  $(t^{n-1}, t^n)$  as

$$\bar{\psi}^n := \tau^{-1} \int_{t^{n-1}}^{t^n} \psi(t) dt \in V, \quad (4.8)$$

with the convention that  $\bar{\psi}^0 = 0 \in V$ . We also let, for all  $(\varphi^i)_{0 \leq i \leq N} \in V^{N+1}$  and all  $1 \leq n \leq N$ ,

$$\delta_t \varphi^n := \frac{\varphi^n - \varphi^{n-1}}{\tau} \in V$$

denote the backward approximation of the first derivative of  $\varphi$  at time  $t^n$ .

We note a preliminary result that will be used in the convergence analysis of Section 4.6. Let  $\psi \in H^1(\mathcal{T}_\tau)$  and  $1 \leq n \leq N$ . Identifying  $\psi|_{(t^{n-1}, t^n)} \in H^1((t^{n-1}, t^n))$  with its absolutely continuous representative in  $(t^{n-1}, t^n]$ , one can use the fundamental theorem of calculus to infer that

$$\psi^n - \bar{\psi}^n = \psi(t^n) - \frac{1}{\tau} \int_{t^{n-1}}^{t^n} \left( \psi(t^n) - \int_s^{t^n} d_t \psi(t) dt \right) ds = \frac{1}{\tau} \int_{t^{n-1}}^{t^n} \int_s^{t^n} d_t \psi(t) dt ds \leq \int_{t^{n-1}}^{t^n} |d_t \psi(t)| dt.$$

Thus, applying the previous result together with the Jensen inequality yields, for all  $\varphi \in H^1(\mathcal{T}_\tau; L^2(\Omega))$ ,

$$\|\varphi^n - \bar{\varphi}^n\|_\Omega^2 \leq \int_\Omega \left( \int_{t^{n-1}}^{t^n} |d_t \varphi(\mathbf{x}, t)| dt \right)^2 d\mathbf{x} \leq \tau \int_{t^{n-1}}^{t^n} \|d_t \varphi(t)\|_\Omega^2 dt \leq \tau \|\varphi\|_{H^1((t^{n-1}, t^n); L^2(\Omega))}^2, \quad (4.9)$$

where  $\varphi(\mathbf{x}, t)$  is a shorthand notation for  $(\varphi(t))(\mathbf{x})$ . As a result of (4.9), we get

$$\sum_{n=1}^N \tau \|\varphi^n - \bar{\varphi}^n\|_\Omega^2 \leq \tau^2 \sum_{n=1}^N \|\varphi\|_{H^1((t^{n-1}, t^n); L^2(\Omega))}^2 =: \tau^2 \|\varphi\|_{H^1(\mathcal{T}_\tau; L^2(\Omega))}^2.$$

### 4.3.3 $L^2$ -orthogonal projectors on local and broken polynomial spaces

For  $X \subset \bar{\Omega}$  and  $k \in \mathbb{N}$ , we denote by  $P^k(X)$  the space spanned by the restriction to  $X$  of scalar-valued,  $d$ -variate polynomials of total degree  $k$ . The  $L^2$ -projector  $\pi_X^k : L^1(X) \rightarrow P^k(X)$  is defined such that, for all  $v \in L^1(X)$ ,

$$\int_X (\pi_X^k v - v) w = 0 \quad \forall w \in P^k(X). \quad (4.10)$$

As a projector,  $\pi_X^k$  is linear and idempotent so that, for all  $v \in P^k(X)$ ,  $\pi_X^k v = v$ . When dealing with the vector-valued polynomial space  $\mathbf{P}^k(X)$  or with the tensor-valued polynomial space  $\mathbb{P}^k(X)$ , we use the boldface notation  $\boldsymbol{\pi}_X^k$  for the corresponding  $L^2$ -orthogonal projectors acting component-wise. At the global level, we denote by  $P^k(\mathcal{T}_h)$ ,  $\mathbf{P}^k(\mathcal{T}_h)$ , and  $\mathbb{P}^k(\mathcal{T}_h)$ , respectively, the spaces of fully discontinuous scalar-valued, vector-valued, and tensor-valued broken polynomial functions on  $\mathcal{T}_h$  of total degree  $k$ , and by  $\pi_h^k$  and  $\boldsymbol{\pi}_h^k$  the  $L^2$ -projectors on  $P^k(\mathcal{T}_h)$  and  $\mathbf{P}^k(\mathcal{T}_h)$ . The following optimal approximation properties for the  $L^2$ -projector  $\pi_X^k$  follow from [69, Lemmas 3.4 and 3.6]: There exists a strictly positive real number  $C_{\text{ap}}$  independent of  $h$  such that, for all  $T \in \mathcal{T}_h$ , all  $l \in \{1, \dots, k+1\}$ , and all  $v \in H^l(T)$ ,

$$|v - \pi_T^k v|_{H^m(T)} + h_T^{\frac{1}{2}} |v - \pi_T^k v|_{H^m(\mathcal{F}_T)} \leq C_{\text{ap}} h_T^{l-m} |v|_{H^l(T)} \quad \forall m \in \{0, \dots, l\}, \quad (4.11)$$

where  $|\cdot|_{H^m(\mathcal{F}_T)}$  is the broken Sobolev seminorm on  $\mathcal{F}_T$ .

### 4.3.4 Discrete spaces

In this section we define the discrete spaces upon which the HHO method corresponding to a polynomial degree  $k \geq 1$  is built.

#### 4.3.4.1 Displacement

The discrete unknowns for the displacement are collected in the space

$$\underline{U}_h^k := \left\{ \underline{\mathbf{v}}_h = ((\mathbf{v}_T)_{T \in \mathcal{T}_h}, (\mathbf{v}_F)_{F \in \mathcal{F}_h}) : \mathbf{v}_T \in \mathbf{P}^k(T) \forall T \in \mathcal{T}_h \text{ and } \mathbf{v}_F \in \mathbf{P}^k(F) \forall F \in \mathcal{F}_h \right\}.$$

For any  $\underline{\mathbf{v}}_h \in \underline{U}_h^k$ , we denote by  $\mathbf{v}_h \in \mathbf{P}^k(\mathcal{T}_h)$  the broken polynomial vector field obtained patching element-based unknowns, so that

$$(\mathbf{v}_h)|_T = \mathbf{v}_T \quad \forall T \in \mathcal{T}_h.$$

The discrete unknowns corresponding to a function  $\mathbf{v} \in \mathbf{H}^1(\Omega)$  are obtained by means of the interpolator  $\underline{\mathbf{I}}_h^k : \mathbf{H}^1(\Omega) \rightarrow \underline{U}_h^k$  such that

$$\underline{\mathbf{I}}_h^k \mathbf{v} := ((\boldsymbol{\pi}_T^k \mathbf{v}|_T)_{T \in \mathcal{T}_h}, (\boldsymbol{\pi}_F^k \mathbf{v}|_F)_{F \in \mathcal{F}_h}). \quad (4.12)$$

For all  $T \in \mathcal{T}_h$ , we denote by  $\underline{U}_T^k$  and  $\underline{\mathbf{I}}_T^k$  the restrictions to  $T$  of  $\underline{U}_h^k$  and  $\underline{\mathbf{I}}_h^k$ , respectively and, for any  $\underline{\mathbf{v}}_h \in \underline{U}_h^k$ , we let  $\underline{\mathbf{v}}_T := (\mathbf{v}_T, (\mathbf{v}_F)_{F \in \mathcal{F}_T})$  collect the local discrete unknowns attached to  $T$ . At each time step, the displacement is sought in the following subspace of  $\underline{U}_h^k$  that strongly accounts for the homogeneous Dirichlet condition (4.1c):

$$\underline{U}_{h,\text{D}}^k := \left\{ \underline{\mathbf{v}}_h = ((\mathbf{v}_T)_{T \in \mathcal{T}_h}, (\mathbf{v}_F)_{F \in \mathcal{F}_h}) \in \underline{U}_h^k : \mathbf{v}_F = \mathbf{0} \quad \forall F \in \mathcal{F}_h^{\text{b}} \right\}.$$

We next prove a discrete version of Korn's first inequality on  $\underline{U}_{h,\text{D}}^k$  that will play a key role in the analysis. To this purpose, we endow the space  $\underline{U}_h^k$  with the discrete strain seminorm  $\|\cdot\|_{\varepsilon,h}$  defined, for all  $\underline{\mathbf{v}}_h \in \underline{U}_h^k$ , such that

$$\|\underline{\mathbf{v}}_h\|_{\varepsilon,h} := \left[ \sum_{T \in \mathcal{T}_h} \left( \|\nabla_s \mathbf{v}_T\|_T^2 + \sum_{F \in \mathcal{F}_T} \frac{\|\mathbf{v}_F - \mathbf{v}_T\|_F^2}{h_F} \right) \right]^{\frac{1}{2}}. \quad (4.13)$$

We will also need the following continuous trace inequality, whose proof follows the arguments of [76, Lemma 1.49] (where a slightly different notion of mesh faces is considered): There exists a strictly positive real number  $C_{\text{tr}}$ , independent of  $h$  but possibly depending on  $\varrho$ , such that, for all  $T \in \mathcal{T}_h$ , all  $\mathbf{v}_T \in \mathbf{H}^1(T)$ , and all  $F \in \mathcal{F}_T$ ,

$$\|\mathbf{v}_T\|_F^2 \leq C_{\text{tr}}^2 \left( \|\nabla_s \mathbf{v}_T\|_T + h_T^{-1} \|\mathbf{v}_T\|_T \right) \|\mathbf{v}_T\|_T. \quad (4.14)$$

**Proposition 4.3** (Discrete Korn's first inequality). *There is a real number  $C_K > 0$ , only depending on  $\Omega$ ,  $d$ , and  $\varrho$  such that, for all  $\underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_{h,D}^k$ ,*

$$\|\mathbf{v}_h\|_\Omega \leq C_K \|\underline{\mathbf{v}}_h\|_{\varepsilon,h}. \quad (4.15)$$

*Remark 4.4* (Strain norm). An immediate consequence of (4.15) is that the map  $\|\cdot\|_{\varepsilon,h}$  is a norm in  $\underline{\mathbf{U}}_{h,D}^k$ .

*Proof.* Let  $\underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_{h,D}^k$ . Since the divergence operator  $\nabla \cdot : \mathbb{H}_s^1(\Omega) \rightarrow \mathbf{L}^2(\Omega)$  is onto (c.f. [31, 56]), there exists  $\boldsymbol{\tau}_{\mathbf{v}_h} \in \mathbb{H}_s^1(\Omega)$  such that  $\nabla \cdot \boldsymbol{\tau}_{\mathbf{v}_h} = \mathbf{v}_h$  and  $\|\boldsymbol{\tau}_{\mathbf{v}_h}\|_{1,\Omega} \leq C_{\text{sj}} \|\mathbf{v}_h\|_\Omega$ , with  $C_{\text{sj}} > 0$  independent of  $h$ . It follows that

$$\|\mathbf{v}_h\|_\Omega^2 = (\mathbf{v}_h, \nabla \cdot \boldsymbol{\tau}_{\mathbf{v}_h})_\Omega = \sum_{T \in \mathcal{T}_h} \left( - \int_T \nabla_s \mathbf{v}_T : \boldsymbol{\tau}_{\mathbf{v}_h} + \sum_{F \in \mathcal{F}_T} \int_F (\mathbf{v}_T - \mathbf{v}_F) \cdot \boldsymbol{\tau}_{\mathbf{v}_h} \mathbf{n}_{TF} \right),$$

where we have integrated by parts element by element and used the fact that  $\boldsymbol{\tau}_{\mathbf{v}_h}$  has continuous normal components across interfaces and that boundary unknowns are set to zero in order to insert  $\mathbf{v}_F$  into the boundary term. Applying a Cauchy–Schwarz inequality on the integrals over the element, a generalized Hölder inequality with exponents  $(2, 2, +\infty)$  on the integrals over faces, using the fact that  $\|\mathbf{n}_{TF}\|_{\mathbf{L}^\infty(F)} \leq 1$ , and invoking a discrete Cauchy–Schwarz inequality on the sum over  $T \in \mathcal{T}_h$ , we infer that

$$\begin{aligned} \|\mathbf{v}_h\|_\Omega^2 &\leq \sum_{T \in \mathcal{T}_h} \left( \|\nabla_s \mathbf{v}_T\|_T \|\boldsymbol{\tau}_{\mathbf{v}_h}\|_T + \sum_{F \in \mathcal{F}_T} \frac{\|\mathbf{v}_F - \mathbf{v}_T\|_F}{h_F^{\frac{1}{2}}} h_F^{\frac{1}{2}} \|\boldsymbol{\tau}_{\mathbf{v}_h}\|_F \right) \\ &\leq \left( \sum_{T \in \mathcal{T}_h} \|\nabla_s \mathbf{v}_T\|_T^2 \right)^{\frac{1}{2}} \|\boldsymbol{\tau}_{\mathbf{v}_h}\|_\Omega + \left( \sum_{T \in \mathcal{T}_h} \sum_{F \in \mathcal{F}_T} \frac{\|\mathbf{v}_F - \mathbf{v}_T\|_F^2}{h_F} \right)^{\frac{1}{2}} \left( \sum_{T \in \mathcal{T}_h} \sum_{F \in \mathcal{F}_T} h_F \|\boldsymbol{\tau}_{\mathbf{v}_h}\|_F^2 \right)^{\frac{1}{2}} \\ &\leq \sqrt{2N_\partial} C_{\text{tr}} \left( \sum_{T \in \mathcal{T}_h} \|\nabla_s \mathbf{v}_T\|_T^2 + \sum_{T \in \mathcal{T}_h} \sum_{F \in \mathcal{F}_T} \frac{\|\mathbf{v}_F - \mathbf{v}_T\|_F^2}{h_F} \right)^{\frac{1}{2}} \|\boldsymbol{\tau}_{\mathbf{v}_h}\|_{1,\Omega} \\ &= \sqrt{2N_\partial} C_{\text{tr}} \|\underline{\mathbf{v}}_h\|_{\varepsilon,h} \|\boldsymbol{\tau}_{\mathbf{v}_h}\|_{1,\Omega}, \end{aligned}$$

where, to pass to the third line, we have estimated  $\|\boldsymbol{\tau}_{\mathbf{v}_h}\|_F$  using the continuous trace inequality (4.14) together with the fact that  $h_F \leq h_T \leq \text{diam}(\Omega) = 1$  for any  $T \in \mathcal{T}_h$  and  $F \in \mathcal{F}_T$ . Thus, invoking the boundedness of the divergence operator, we get

$$\|\mathbf{v}_h\|_\Omega^2 \leq \sqrt{2N_\partial} C_{\text{sj}} C_{\text{tr}} \|\underline{\mathbf{v}}_h\|_{\varepsilon,h} \|\mathbf{v}_h\|_\Omega,$$

which yields the conclusion with  $C_K = \sqrt{2N_\partial} C_{\text{sj}} C_{\text{tr}}$ .  $\square$

#### 4.3.4.2 Pore pressure

At each time step, the discrete pore pressure is sought in the space

$$P_h^k := \begin{cases} P^k(\mathcal{T}_h) & \text{if } C_0 > 0, \\ P_0^k(\mathcal{T}_h) := \{q_h \in P^k(\mathcal{T}_h) : \int_{\Omega} q_h = 0\} & \text{if } C_0 = 0. \end{cases}$$

For any internal face  $F \in \mathcal{F}_h^i$ , we denote by  $T_{F,1}, T_{F,2} \in \mathcal{T}_h$  the two mesh elements that share  $F$  as a face, that is to say  $F \subset \partial T_{F,1} \cap \partial T_{F,2}$  and  $T_{F,1} \neq T_{F,2}$  (the ordering of the elements is arbitrary but fixed), and we set

$$\kappa_{F,i} := (\boldsymbol{\kappa}_{|T_{F,i}} \mathbf{n}_{T_{F,i}F}) \cdot \mathbf{n}_{T_{F,i}F} \quad \text{for } i \in \{1, 2\}, \quad \kappa_F := \frac{2\kappa_{F,1}\kappa_{F,2}}{\kappa_{F,1} + \kappa_{F,2}}. \quad (4.16)$$

For all  $q_h \in P_h^k$ , we denote by  $q_T$  the restriction of  $q_h$  to an element  $T \in \mathcal{T}_h$  and we define the discrete seminorm

$$\|q_h\|_{\kappa,h} := \left( \sum_{T \in \mathcal{T}_h} \|\boldsymbol{\kappa}^{\frac{1}{2}} \nabla q_T\|_T^2 + \sum_{F \in \mathcal{F}_h^i} \frac{\kappa_F}{h_F} \|q_{T_{F,1}} - q_{T_{F,2}}\|_F^2 \right)^{\frac{1}{2}}. \quad (4.17)$$

The fact that in (4.17) boundary terms only appear on internal faces reflects the homogeneous Neumann boundary condition (4.1d).

Using the surjectivity of the divergence operator  $\nabla \cdot : \mathbf{H}_0^1(\Omega) \rightarrow L_0^2(\Omega)$  and proceeding as in the proof of the discrete Korn inequality (4.15), a discrete Poincaré-Wirtinger inequality in  $P^k(\mathcal{T}_h)$  is readily inferred, namely one has the existence of  $C_P > 0$ , only depending on  $\Omega$ ,  $d$ , and  $\varrho$  such that, for all  $q_h \in P^k(\mathcal{T}_h)$ ,

$$\|q_h - \pi_{\Omega}^0 q_h\|_{\Omega} \leq C_P \underline{\kappa}^{-\frac{1}{2}} \|q_h\|_{\kappa,h}.$$

This result ensures, in particular, that the seminorm  $\|\cdot\|_{\kappa,h}$  defined in (4.17) is a norm on  $P_0^k(\mathcal{T}_h)$ . For a proof of more general Sobolev inequalities on broken polynomial spaces, we refer the reader to [75] and [76, Section 5.1.2].

## 4.4 Discretization

In this section we define the discrete counterparts of the elasticity, hydro-mechanical coupling, and Darcy operators, and formulate the HHO–dG scheme for problem (4.6).

### 4.4.1 Nonlinear elasticity operator

The discretization of the nonlinear elasticity operator closely follows the one presented in Section 3.4. We define the local symmetric gradient reconstruction  $\mathbf{G}_{s,T}^k : \underline{\mathbf{U}}_T^k \rightarrow \mathbb{P}_s^k(T)$  such that, for a given  $\underline{\mathbf{v}}_T = (\mathbf{v}_T, (\mathbf{v}_F)_{F \in \mathcal{F}_T}) \in \underline{\mathbf{U}}_T^k$ ,  $\mathbf{G}_{s,T}^k \underline{\mathbf{v}}_T \in \mathbb{P}_s^k(T)$  solves

$$\int_T \mathbf{G}_{s,T}^k \underline{\mathbf{v}}_T : \boldsymbol{\tau} = - \int_T \mathbf{v}_T \cdot (\nabla \cdot \boldsymbol{\tau}) + \sum_{F \in \mathcal{F}_T} \int_F \mathbf{v}_F \cdot (\boldsymbol{\tau} \mathbf{n}_{TF}) \quad \forall \boldsymbol{\tau} \in \mathbb{P}_s^k(T). \quad (4.18)$$



Existence and uniqueness of  $\mathbf{G}_{s,T}^k \underline{\mathbf{v}}_T$  follow from the Riesz representation theorem in  $\mathbb{P}_s^k(T)$  for the  $L^2(T)^{d \times d}$ -inner product. This definition is motivated by the following property.

**Proposition 4.5** (Commuting property for the local symmetric gradient reconstruction). *For all  $\mathbf{v} \in \mathbf{H}^1(T)$ , it holds that*

$$\mathbf{G}_{s,T}^k \mathbf{I}_T^k \mathbf{v} = \boldsymbol{\pi}_T^k(\nabla_s \mathbf{v}). \quad (4.19)$$

*Remark 4.6* (Approximation properties of the local symmetric gradient reconstruction). The commuting property (4.19) combined with (4.11) shows that  $\mathbf{G}_{s,T}^k \mathbf{I}_T^k \mathbf{v}$  optimally approximates  $\nabla_s \mathbf{v}$  in  $\mathbb{P}_s^k(T)$ .

*Proof.* We start by noticing that it holds, for all  $\boldsymbol{\alpha} \in \mathbb{R}_s^{d \times d}$  and all  $\boldsymbol{\beta} \in \mathbb{R}^{d \times d}$ , denoting by  $\boldsymbol{\beta}_s := \frac{1}{2}(\boldsymbol{\beta} + \boldsymbol{\beta}^T)$  the symmetric part of  $\boldsymbol{\beta}$ ,

$$\boldsymbol{\alpha} : \boldsymbol{\beta} = \boldsymbol{\alpha} : \boldsymbol{\beta}_s. \quad (4.20)$$

For all  $\boldsymbol{\tau} \in \mathbb{P}^k(T)$  we can then write

$$\begin{aligned} \int_T \mathbf{G}_{s,T}^k \mathbf{I}_T^k \mathbf{v} : \boldsymbol{\tau} &= \int_T \mathbf{G}_{s,T}^k \mathbf{I}_T^k \mathbf{v} : \boldsymbol{\tau}_s \\ &= - \int_T \boldsymbol{\pi}_T^k \mathbf{v} \cdot (\nabla \cdot \boldsymbol{\tau}_s) + \sum_{F \in \mathcal{F}_T} \int_F \boldsymbol{\pi}_F^k \mathbf{v} \cdot (\boldsymbol{\tau}_s \mathbf{n}_{TF}) \\ &= - \int_T \mathbf{v} \cdot (\nabla \cdot \boldsymbol{\tau}_s) + \sum_{F \in \mathcal{F}_T} \int_F \mathbf{v} \cdot (\boldsymbol{\tau}_s \mathbf{n}_{TF}) \\ &= \int_T \nabla \mathbf{v} : \boldsymbol{\tau}_s = \int_T \nabla_s \mathbf{v} : \boldsymbol{\tau}_s = \int_T \nabla_s \mathbf{v} : \boldsymbol{\tau} = \int_T \boldsymbol{\pi}_T^k(\nabla_s \mathbf{v}) : \boldsymbol{\tau}, \end{aligned}$$

where we have used (4.20) with  $\boldsymbol{\alpha} = \mathbf{G}_{s,T}^k \mathbf{I}_T^k \mathbf{v}$  and  $\boldsymbol{\beta} = \boldsymbol{\tau}$  in the first line, the definition (4.18) of the local symmetric gradient with  $\underline{\mathbf{v}}_T = \mathbf{I}_T^k \mathbf{v}$  in the second line, and definition (4.10) after observing that  $\nabla \cdot \boldsymbol{\tau}_s \in \mathbf{P}^{k-1}(T) \subset \mathbf{P}^k(T)$  and  $\boldsymbol{\tau}_s \mathbf{n}_{TF} \in \mathbf{P}^k(F)$  for all  $F \in \mathcal{F}_T$  to remove the  $L^2$ -orthogonal projectors in the third line. In the fourth line, we have used an integration by parts, then invoked (4.20) first with  $\boldsymbol{\alpha} = \boldsymbol{\tau}_s$  and  $\boldsymbol{\beta} = \nabla \mathbf{v}$ , then with  $\boldsymbol{\alpha} = \nabla_s \mathbf{v}$  and  $\boldsymbol{\beta} = \boldsymbol{\tau}$ , and we have used the definition (4.10) of  $\boldsymbol{\pi}_T^k$  to conclude.  $\square$

From  $\mathbf{G}_{s,T}^k$ , we define the local displacement reconstruction operator  $\mathbf{r}_T^{k+1} : \underline{\mathbf{U}}_T^k \rightarrow \mathbf{P}^{k+1}(T)$  such that, for all  $\underline{\mathbf{v}}_T \in \underline{\mathbf{U}}_T^k$ ,

$$\begin{aligned} \int_T (\nabla_s \mathbf{r}_T^{k+1} \underline{\mathbf{v}}_T - \mathbf{G}_{s,T}^k \underline{\mathbf{v}}_T) : \nabla_s \mathbf{w} &= 0 \quad \forall \mathbf{w} \in \mathbf{P}^{k+1}(T), \\ \int_T \mathbf{r}_T^{k+1} \underline{\mathbf{v}}_T &= \int_T \mathbf{v}_T, \quad \int_T \nabla \times (\mathbf{r}_T^{k+1} \underline{\mathbf{v}}_T) = \sum_{F \in \mathcal{F}_T} \int_F \mathbf{n}_{TF} \times \mathbf{v}_F, \end{aligned}$$

where  $\nabla \times$  denotes the curl operator and  $\mathbf{v} \times \mathbf{w}$  the vector product of two vectors  $\mathbf{v}, \mathbf{w} \in \mathbb{R}^d$ . Optimal approximation properties for  $\mathbf{r}_T^{k+1} \mathbf{I}_T^k$  have been recently proved in [32, Appendix A]

generalizing the ones of [74, Lemma 2]. The optimal approximation properties of  $\mathbf{r}_T^{k+1} \underline{\mathbf{I}}_T^k$  are required to infer (4.46) below.

The discretization of the nonlinear elasticity operator is realized by the function  $a_h : \underline{\mathbf{U}}_h^k \times \underline{\mathbf{U}}_h^k \rightarrow \mathbb{R}$  such that, for all  $\underline{\mathbf{w}}_h, \underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_h^k$ ,

$$a_h(\underline{\mathbf{w}}_h, \underline{\mathbf{v}}_h) := \sum_{T \in \mathcal{T}_h} \left( \int_T \boldsymbol{\sigma}(\cdot, \mathbf{G}_{s,T}^k \underline{\mathbf{u}}_T) : \mathbf{G}_{s,T}^k \underline{\mathbf{v}}_T + \sum_{F \in \mathcal{F}_T} \frac{\gamma}{h_F} \int_F \mathbf{A}_{TF}^k \underline{\mathbf{u}}_T \cdot \mathbf{A}_{TF}^k \underline{\mathbf{v}}_T \right), \quad (4.21)$$

where  $\gamma > 0$  denotes a user-dependent parameter and we penalize in a least-square sense the face-based residual  $\mathbf{A}_{TF}^k : \underline{\mathbf{U}}_T^k \rightarrow \mathbf{P}^k(F)$  such that, for all  $T \in \mathcal{T}_h$ , all  $\underline{\mathbf{v}}_T \in \underline{\mathbf{U}}_T^k$ , and all  $F \in \mathcal{F}_T$ ,

$$\mathbf{A}_{TF}^k \underline{\mathbf{v}}_T := \boldsymbol{\pi}_F^k(\mathbf{r}_T^{k+1} \underline{\mathbf{v}}_T - \mathbf{v}_F) - \boldsymbol{\pi}_T^k(\mathbf{r}_T^{k+1} \underline{\mathbf{v}}_T - \mathbf{v}_T).$$

This definition ensures that  $\mathbf{A}_{TF}^k$  vanishes whenever its argument is of the form  $\underline{\mathbf{I}}_T^k \mathbf{w}$  with  $\mathbf{w} \in \mathbf{P}^{k+1}(T)$ , a crucial property to obtain high-order error estimates (cf. Theorems 2.12, 3.16). For further use, we note the following seminorm equivalence, which can be proved using the arguments of [74, Lemma 4]: For all  $\underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_h^k$ ,

$$C_{\text{eq}}^{-2} \|\underline{\mathbf{v}}_h\|_{\varepsilon,h}^2 \leq \sum_{T \in \mathcal{T}_h} \left( \|\mathbf{G}_{s,T}^k \underline{\mathbf{v}}_T\|_T^2 + \sum_{F \in \mathcal{F}_T} h_F^{-1} \|\mathbf{A}_{TF}^k \underline{\mathbf{v}}_T\|_F^2 \right) \leq C_{\text{eq}}^2 \|\underline{\mathbf{v}}_h\|_{\varepsilon,h}^2, \quad (4.22)$$

where  $C_{\text{eq}} > 0$  is independent of  $h$ , and the discrete strain seminorm  $\|\cdot\|_{\varepsilon,h}$  is defined by (4.13). By (4.2b), this implies the coercivity of  $a_h$ .

*Remark 4.7* (Choice of the stabilization parameter). The constants  $C_{\text{gr}}, C_{\text{cv}}$  appearing in (4.2) satisfy  $C_{\text{cv}}^2 \leq C_{\text{gr}}$ . Indeed, owing to (4.2b), the Cauchy–Schwarz inequality, and (4.2a), it holds for all  $\boldsymbol{\tau} \in \mathbb{R}_s^{d \times d}$ ,

$$C_{\text{cv}}^2 |\boldsymbol{\tau}|_{d \times d}^2 \leq \boldsymbol{\sigma}(\mathbf{x}, \boldsymbol{\tau}) : \boldsymbol{\tau} \leq |\boldsymbol{\sigma}(\mathbf{x}, \boldsymbol{\tau})|_{d \times d} |\boldsymbol{\tau}|_{d \times d} \leq C_{\text{gr}} |\boldsymbol{\tau}|_{d \times d}^2. \quad (4.23)$$

Thus, we choose the stabilization parameter  $\gamma$  in (4.21) such that

$$\gamma \in [C_{\text{cv}}^2, C_{\text{gr}}]. \quad (4.24)$$

For the linear elasticity model (4.3), we have  $C_{\text{gr}} = 2\mu + d\lambda$  and  $C_{\text{cv}} = \sqrt{2\mu}$ , so that a natural choice for the stabilization parameter is  $\gamma = 2\mu$ .

#### 4.4.2 Hydro-mechanical coupling

The hydro-mechanical coupling is realized by means of the bilinear form  $b_h$  on  $\underline{\mathbf{U}}_h^k \times P^k(\mathcal{T}_h)$  such that, for all  $\underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_h^k$  and all  $q_h \in P^k(\mathcal{T}_h)$ ,

$$b_h(\underline{\mathbf{v}}_h, q_h) := \sum_{T \in \mathcal{T}_h} \left( \int_T \mathbf{v}_T \cdot \nabla q_T - \sum_{F \in \mathcal{F}_T} \int_F \mathbf{v}_F \cdot q_T \mathbf{n}_{TF} \right), \quad (4.25)$$

where we remind the reader that  $q_T := q_h|_T$ . It can be checked using Cauchy–Schwarz inequalities together with the definition (4.13) of the discrete strain seminorm and discrete trace inequalities that there exists  $C_{bd} > 0$  independent of  $h$  such that

$$b_h(\underline{\mathbf{v}}_h, q_h) \leq C_{bd} \|\underline{\mathbf{v}}_h\|_{\varepsilon, h} \|q_h\|_{\Omega}.$$

Additionally, using the strongly enforced boundary condition in  $\underline{\mathbf{U}}_{h, D}^k$ , it can be proved that

$$b_h(\underline{\mathbf{v}}_h, 1) = 0, \quad \forall \underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_{h, D}^k. \quad (4.26)$$

Finally, we note the following lemma stating that the hybrid interpolator  $\underline{\mathbf{I}}_h^k : \mathbf{H}^1(\Omega) \rightarrow \underline{\mathbf{U}}_h^k$  is a Fortin operator.

**Lemma 4.8** (Fortin operator). *For all  $\mathbf{v} \in \mathbf{H}^1(\Omega)$  and all  $q_h \in P^k(\mathcal{T}_h)$ , the interpolator  $\underline{\mathbf{I}}_h^k$  satisfies*

$$\|\underline{\mathbf{I}}_h^k \mathbf{v}\|_{\varepsilon, h} \leq C_{st} |\mathbf{v}|_{1, \Omega}, \quad (4.27a)$$

$$b_h(\underline{\mathbf{I}}_h^k \mathbf{v}, q_h) = b(\mathbf{v}, q_h), \quad (4.27b)$$

where the strictly positive real number  $C_{st}$  is independent of  $h$ .

*Proof.* (i) *Proof of (4.27a).* Recalling the definitions (4.13) of the discrete strain seminorm and (4.12) of the global interpolator, we can write

$$\begin{aligned} \|\underline{\mathbf{I}}_h^k \mathbf{v}\|_{\varepsilon, h}^2 &= \sum_{T \in \mathcal{T}_h} \left( \|\nabla_s \pi_T^k \mathbf{v}\|_T^2 + \sum_{F \in \mathcal{F}_T} h_F^{-1} \|\pi_F^k \mathbf{v} - \pi_T^k \mathbf{v}\|_F^2 \right) \\ &\leq \sum_{T \in \mathcal{T}_h} \left( 2\|\nabla_s(\pi_T^k \mathbf{v} - \mathbf{v})\|_T^2 + 2\|\nabla_s \mathbf{v}\|_T^2 + \sum_{F \in \mathcal{F}_T} h_F^{-1} \|\mathbf{v} - \pi_T^k \mathbf{v}\|_F^2 \right) \leq C_{st} |\mathbf{v}|_{1, \Omega}, \end{aligned}$$

with  $C_{st} > 0$  independent of  $h$ . To pass to the second line, we have used a triangle inequality after inserting  $\pm \nabla_s \mathbf{v}$  into the first term, and we have used the linearity, idempotency, and boundedness of  $\pi_F^k$  to write  $\|\pi_F^k \mathbf{v} - \pi_T^k \mathbf{v}\|_F = \|\pi_F^k(\mathbf{v} - \pi_T^k \mathbf{v})\|_F \leq \|\mathbf{v} - \pi_T^k \mathbf{v}\|_F$ . To conclude, we have used (4.11) respectively with  $l = m = 1$  and with  $l = 1$  and  $m = 0$  to bound the first and third terms inside the summation.

(ii) *Proof of (4.27b).* Recalling the definitions (4.25) of  $b_h(\cdot, \cdot)$  and (4.12) of the global interpolator, we can write for all  $q_h \in P^k(\mathcal{T}_h)$ ,

$$\begin{aligned} b_h(\underline{\mathbf{I}}_h^k \mathbf{v}, q_h) &= \sum_{T \in \mathcal{T}_h} \left( \int_T \pi_T^k \mathbf{v} \cdot \nabla q_T - \int_F \pi_F^k \mathbf{v}|_F \cdot q_T \mathbf{n}_{TF} \right) \\ &= \sum_{T \in \mathcal{T}_h} \left( \int_T \mathbf{v} \cdot \nabla q_T - \int_F \mathbf{v} \cdot q_T \mathbf{n}_{TF} \right) = b(\mathbf{v}, q_h), \end{aligned}$$

where we have used definition (4.10) after observing that  $\nabla q_T \in \mathbf{P}^{k-1}(T) \subset \mathbf{P}^k(T)$  and  $q_T|_F \mathbf{n}_{TF} \in \mathbf{P}^k(F)$  to remove the  $L^2$ -orthogonal projectors in the second line, and integration by parts over  $T \in \mathcal{T}_h$  to conclude.  $\square$

As a result of the previous Lemma, one has the following inf-sup condition, cf. Proposition A.1 for the proof.

**Proposition 4.9.** *There is a strictly positive real number  $\beta$  independent of  $h$  such that, for all  $q_h \in P_0^k(\mathcal{T}_h)$ ,*

$$\|q_h\|_\Omega \leq \beta \sup_{\underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_{h,D}^k \setminus \{\mathbf{0}\}} \frac{b_h(\underline{\mathbf{v}}_h, q_h)}{\|\underline{\mathbf{v}}_h\|_{\varepsilon,h}}. \quad (4.28)$$

### 4.4.3 Darcy operator

The discretization of the Darcy operator is based on the Symmetric Weighted Interior Penalty method of [77], cf. also [76, Section 4.5]. For all  $F \in \mathcal{F}_h^i$  and all  $q_h \in P^k(\mathcal{T}_h)$ , we define the jump and weighted average operators such that

$$[q_h]_F := q_{T_{F,1}} - q_{T_{F,2}}, \quad \{q_h\}_F := \frac{\sqrt{\kappa_{F,2}}}{\sqrt{\kappa_{F,1}} + \sqrt{\kappa_{F,2}}} q_{T_{F,1}} + \frac{\sqrt{\kappa_{F,1}}}{\sqrt{\kappa_{F,1}} + \sqrt{\kappa_{F,2}}} q_{T_{F,2}},$$

with  $T_{F,1}, T_{F,2} \in \mathcal{T}_h$ ,  $T_{F,1} \neq T_{F,2}$ , such that  $F \subset \partial T_{F,1} \cap \partial T_{F,2}$  and  $\kappa_{F,1}, \kappa_{F,2}$  defined in (4.16). The bilinear form  $c_h$  on  $P^k(\mathcal{T}_h) \times P^k(\mathcal{T}_h)$  is defined such that, for all  $q_h, r_h \in P^k(\mathcal{T}_h)$ ,

$$\begin{aligned} c_h(r_h, q_h) &:= \int_\Omega \boldsymbol{\kappa} \nabla_h r_h \cdot \nabla_h q_h + \sum_{F \in \mathcal{F}_h^i} \frac{S_{\kappa F}}{h_F} \int_F [r_h]_F [q_h]_F \\ &\quad - \sum_{F \in \mathcal{F}_h^i} \int_F ([r_h]_F \{\boldsymbol{\kappa} \nabla_h q_h\}_F + [q_h]_F \{\boldsymbol{\kappa} \nabla_h r_h\}_F) \cdot \mathbf{n}_{TF}, \end{aligned}$$

where we have introduced the broken gradient operator  $\nabla_h$  on  $\mathcal{T}_h$  and we have denoted by  $\varsigma > \underline{\varsigma} > 0$  a user-defined penalty parameter chosen large enough to ensure the coercivity of  $c_h$  (the proof is similar to [76, Lemma 4.51]):

$$c_h(q_h, q_h) \geq (\varsigma - \underline{\varsigma})(1 + \varsigma)^{-1} \|q_h\|_{\boldsymbol{\kappa},h}^2, \quad \forall q_h \in P_h^k.$$

Since, under this condition,  $c_h$  is a symmetric positive definite bilinear form on the broken polynomial space  $P_h^k$ , we can define an associated norm by setting  $\|\cdot\|_{c,h} := c_h(\cdot, \cdot)^{\frac{1}{2}}$ .

The following consistency result can be proved adapting the arguments of [76, Chapter 4] to homogeneous Neumann boundary conditions and will be instrumental for the analysis. We define the functional spaces  $P_* := \{r \in H^1(\Omega) \cap H^2(P_\Omega) : \boldsymbol{\kappa} \nabla r \cdot \mathbf{n} = 0 \text{ on } \partial\Omega\}$  and set  $P_{*h}^k := P_* + P_h^k$ . Extending the bilinear form  $c_h$  to  $P_{*h}^k \times P_{*h}^k$ , it is inferred that, for all  $r \in P_*$ ,

$$-(\nabla \cdot (\boldsymbol{\kappa} \nabla r), q)_\Omega = c_h(r, q) \quad \forall q \in P_{*h}. \quad (4.29)$$

### 4.4.4 Discrete problem

For all  $1 \leq n \leq N$ , the discrete solution  $(\underline{\mathbf{u}}_h^n, p_h^n) \in \underline{\mathbf{U}}_{h,D}^k \times P_h^k$  at time  $t^n$  is such that, for all  $(\underline{\mathbf{v}}_h, q_h) \in \underline{\mathbf{U}}_{h,D}^k \times P^k(\mathcal{T}_h)$ ,

$$a_h(\underline{\mathbf{u}}_h^n, \underline{\mathbf{v}}_h) + b_h(\underline{\mathbf{v}}_h, p_h^n) = (\bar{\mathbf{f}}^n, \mathbf{v}_h)_\Omega, \quad (4.30a)$$

$$C_0(\delta_t p_h^n, q_h)_\Omega - b_h(\delta_t \underline{\mathbf{u}}_h^n, q_h) + c_h(p_h^n, q_h) = (\bar{g}^n, q_h)_\Omega, \quad (4.30b)$$

with  $\bar{\mathbf{f}}^n \in \mathbf{L}^2(\Omega)$  and  $\bar{g}^n \in L^2(\Omega)$  defined in Section 4.3.2. In order to start the time-stepping scheme we need to initialize the discrete fluid content. This is done by setting  $\phi_h^0$  equal to the  $L^2$ -orthogonal projection of  $\phi^0$  on  $P^k(\mathcal{T}_h)$  according to (4.6c). This implies in particular that

$$C_0(p_h^0, q_h)_\Omega - b_h(\underline{\mathbf{u}}_h^0, q_h) := (\phi^0, q_h)_\Omega \quad \forall q_h \in P^k(\mathcal{T}_h). \quad (4.30c)$$

*Remark 4.10* (Advance in time scheme). The modified backward Euler scheme obtained by taking time averages instead of pointwise evaluation of the right-hand sides in (4.30) can be interpreted as a low-order discontinuous Galerkin time-stepping method, cf. [179, 183].

Notice that other time discretizations could be used, but we have decided to focus on the backward Euler scheme to keep the proofs as simple as possible. From the practical point of view, at each time step  $n$ , the discrete nonlinear system (4.30) can be solved by the Newton method using as initial guess the solution at step  $(n - 1)$ . The size of the linear system to be solved at each Newton iteration can be reduced by statically condensing a large part of the unknowns as described in Section 2.5.

## 4.5 Stability and well-posedness

In this section we study the stability of problem (4.30) and prove its well-posedness. We start with an a priori estimate on the discrete solutions not requiring conditions on the time step  $\tau$  and robust with respect to vanishing storage coefficients and small permeability fields.

**Proposition 4.11** (A priori estimate). *Denote by  $(\underline{\mathbf{u}}_h^n, p_h^n)_{1 \leq n \leq N}$  the solution to (4.30). Under Assumption 4.1 on the stress-strain relation and the regularity on the data  $\mathbf{f}$ ,  $g$ , and  $\phi^0$  assumed in Section 4.2.3, it holds*

$$\begin{aligned} & \sum_{n=1}^N \tau \|\underline{\mathbf{u}}_h^n\|_{\varepsilon, h}^2 + \sum_{n=1}^N \tau \left( \|p_h^n - \pi_\Omega^0 p_h^n\|_\Omega^2 + C_0 \|p_h^n\|_\Omega^2 \right) + \|s_h^N\|_{c, h}^2 \leq \\ & C \left( \|\mathbf{f}\|_{L^2(L^2(\Omega))}^2 + t_F^2 \|g\|_{L^2(L^2(\Omega))}^2 + t_F \|\phi^0\|_\Omega^2 + \frac{t_F^2}{C_0} \|\pi_\Omega^0 g\|_{L^2(L^2(\Omega))}^2 + \frac{t_F}{C_0} \|\pi_\Omega^0 \phi^0\|_\Omega^2 \right). \end{aligned} \quad (4.31)$$

where  $C > 0$  denotes a real number independent of  $h$ ,  $\tau$ , the physical parameters  $C_0$  and  $\kappa$ , and the final time  $t_F$ . In (4.31), we have defined  $s_h^N := \sum_{n=1}^N \tau p_h^n$  and we have adopted the convention that  $C_0^{-1} \|\pi_\Omega^0 g\|_{L^2(L^2(\Omega))}^2 = 0$  and  $C_0^{-1} \|\pi_\Omega^0 \phi^0\|_\Omega^2 = 0$  if  $C_0 = 0$ .

*Remark 4.12.* In order to prove the a priori bound (4.31), no additional time regularity assumption on the loading term  $\mathbf{f}$  and the mass source  $g$  are needed, whereas the stability estimate of Lemma 2.7, valid for linear stress-strain relation, requires  $\mathbf{f} \in C^1(\mathbf{L}^2(\Omega))$  and  $g \in C^0(L^2(\Omega))$ . On the other hand, Proposition 4.11 gives an estimate of the discrete displacement and pressure in the  $L^2$ -norm in time, while Lemma 2.7 ensures a control in the  $L^\infty$ -norm in time. However, under additional requirements on the stress-strain law (for instance Assumption 4.15) and  $H^1$ -regularity in time of  $\mathbf{f}$ , a stronger version of (4.31) can be inferred, establishing in particular an estimate in the  $L^\infty$ -norm in time.

*Proof.* (i) Estimate of  $\|p_h^n - \pi_\Omega^0 p_h^n\|_\Omega$ . The growth property of the stress-strain function (4.2a) together with the Cauchy–Schwarz inequality, assumption (4.24), and the second inequality in (4.22) yield, for all  $1 \leq n \leq N$  and all  $\underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_{h,D}^k$ ,

$$\begin{aligned} a_h(\underline{\mathbf{u}}_h^n, \underline{\mathbf{v}}_h) &= \sum_{T \in \mathcal{T}_h} \left( \int_T \boldsymbol{\sigma}(\cdot, \mathbf{G}_{s,T}^k \underline{\mathbf{u}}_T^n) : \mathbf{G}_{s,T}^k \underline{\mathbf{v}}_T + \sum_{F \in \mathcal{F}_T} \frac{\gamma}{h_F} \int_F \boldsymbol{\Delta}_{TF}^k \underline{\mathbf{u}}_T^n \cdot \boldsymbol{\Delta}_{TF}^k \underline{\mathbf{v}}_T \right) \\ &\leq C_{\text{gr}} \sum_{T \in \mathcal{T}_h} \left( \|\mathbf{G}_{s,T}^k \underline{\mathbf{u}}_T^n\|_T \|\mathbf{G}_{s,T}^k \underline{\mathbf{v}}_T\|_T + \sum_{F \in \mathcal{F}_T} \frac{1}{h_F} \|\boldsymbol{\Delta}_{TF}^k \underline{\mathbf{u}}_T^n\|_F \|\boldsymbol{\Delta}_{TF}^k \underline{\mathbf{v}}_T\|_F \right) \\ &\leq C_{\text{gr}} C_{\text{eq}}^2 \|\underline{\mathbf{u}}_h^n\|_{\varepsilon,h} \|\underline{\mathbf{v}}_h\|_{\varepsilon,h}. \end{aligned} \quad (4.32)$$

Using the inf-sup condition (4.28), (4.26), and the mechanical equilibrium equation (4.30a), we get, for any  $1 \leq n \leq N$ ,

$$\|p_h^n - \pi_\Omega^0 p_h^n\|_\Omega \leq \beta \sup_{\underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_{h,D}^k \setminus \{\mathbf{0}\}} \frac{b_h(\underline{\mathbf{v}}_h, p_h^n - \pi_\Omega^0 p_h^n)}{\|\underline{\mathbf{v}}_h\|_{\varepsilon,h}} = \beta \sup_{\underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_{h,D}^k \setminus \{\mathbf{0}\}} \frac{(\bar{\mathbf{f}}^n, \underline{\mathbf{v}}_h)_\Omega - a_h(\underline{\mathbf{u}}_h^n, \underline{\mathbf{v}}_h)}{\|\underline{\mathbf{v}}_h\|_{\varepsilon,h}}.$$

Therefore, owing to the discrete Korn inequality (4.15) and to (4.32), we infer from the previous bound that

$$\|p_h^n - \pi_\Omega^0 p_h^n\|_\Omega \leq \beta \left( C_K \|\bar{\mathbf{f}}^n\|_\Omega + C_{\text{gr}} C_{\text{eq}}^2 \|\underline{\mathbf{u}}_h^n\|_{\varepsilon,h} \right). \quad (4.33)$$

(ii) *Energy balance.* For all  $1 \leq n \leq N$ , summing (4.30b) at times  $1 \leq i \leq n$  and taking  $q_h = \tau^2 p_h^n$  as a test function yields

$$\tau C_0 (p_h^n, p_h^n)_\Omega - \tau b_h(\underline{\mathbf{u}}_h^n, p_h^n) + \sum_{i=1}^n \tau^2 c_h(p_h^i, p_h^n) = \sum_{i=1}^n \tau^2 (\bar{\mathbf{g}}^i, p_h^n)_\Omega + \tau (\phi^0, p_h^n)_\Omega. \quad (4.34)$$

Moreover, using the linearity of  $c_h$  and the formula  $2x(x-y) = x^2 + (x-y)^2 - y^2$ , the third term in the left-hand side of (4.34) can be rewritten as

$$\sum_{i=1}^n \tau^2 c_h(p_h^i, p_h^n) = \tau c_h \left( \sum_{i=1}^n \tau p_h^i, p_h^n \right) = \tau c_h(s_h^n, \delta_t s_h^n) = \frac{1}{2} \left( \|s_h^n\|_{c,h}^2 + \|\delta_t s_h^n\|_{c,h}^2 - \|s_h^{n-1}\|_{c,h}^2 \right),$$

where we have set  $s_h^0 := 0$ ,  $s_h^n := \sum_{i=1}^n \tau p_h^i$  for any  $1 \leq n \leq N$ , and observed that  $p_h^n = \delta_t s_h^n$ . Therefore, summing (4.34) and (4.30a) at discrete time  $n$  with  $\underline{\mathbf{v}}_h = \tau \underline{\mathbf{u}}_h^n$ , leads to

$$\tau a_h(\underline{\mathbf{u}}_h^n, \underline{\mathbf{u}}_h^n) + \tau C_0 \|p_h^n\|_\Omega^2 + \frac{1}{2} \left( \|s_h^n\|_{c,h}^2 - \|s_h^{n-1}\|_{c,h}^2 \right) \leq \tau (\bar{\mathbf{f}}^n, \underline{\mathbf{u}}_h^n)_\Omega + \sum_{i=1}^n \tau^2 (\bar{\mathbf{g}}^i, p_h^n)_\Omega + \tau (\phi^0, p_h^n)_\Omega.$$

Summing the previous relation for  $1 \leq n \leq N$ , telescoping out the appropriate summands, and using the coercivity property (4.2b), assumption (4.24), and the first inequality in (4.22), we get

$$\frac{C_{\text{cv}}^2}{C_{\text{eq}}^2} \sum_{n=1}^N \tau \|\underline{\mathbf{u}}_h^n\|_{\varepsilon,h}^2 + C_0 \sum_{n=1}^N \tau \|p_h^n\|_\Omega^2 + \frac{1}{2} \|s_h^N\|_{c,h}^2 \leq \sum_{n=1}^N \tau (\bar{\mathbf{f}}^n, \underline{\mathbf{u}}_h^n)_\Omega + \sum_{n=1}^N \tau (G^n + \phi^0, p_h^n)_\Omega, \quad (4.35)$$

with the notation  $G^n := \sum_{i=1}^n \tau \bar{g}^i = \int_0^{t^n} g(t) dt$ . We denote by  $\mathcal{R}$  the right-hand side of (4.35) and proceed to find a suitable upper bound.

(iii) *Upper bound for  $\mathcal{R}$ .* For the first term in the right-hand side of (4.35), using the Cauchy–Schwarz, discrete Korn (4.15), and Young inequalities, we obtain

$$\begin{aligned} \sum_{n=1}^N \tau (\bar{\mathbf{f}}^n, \mathbf{u}_h^n)_\Omega &\leq C_K \left( \sum_{n=1}^N \tau \|\bar{\mathbf{f}}^n\|_\Omega^2 \right)^{\frac{1}{2}} \left( \sum_{n=1}^N \tau \|\underline{\mathbf{u}}_h^n\|_{\varepsilon,h}^2 \right)^{\frac{1}{2}} \\ &\leq \frac{C_K^2 C_{\text{eq}}^2}{C_{\text{cv}}^2} \|\mathbf{f}\|_{L^2(\mathbf{L}^2(\Omega))}^2 + \frac{C_{\text{cv}}^2}{4C_{\text{eq}}^2} \sum_{n=1}^N \tau \|\underline{\mathbf{u}}_h^n\|_{\varepsilon,h}^2, \end{aligned} \quad (4.36)$$

where we have used the Jensen inequality to infer that

$$\sum_{n=1}^N \tau \|\bar{\mathbf{f}}^n\|_\Omega^2 = \sum_{n=1}^N \frac{1}{\tau} \int_\Omega \left( \int_{t^{n-1}}^{t^n} \mathbf{f}(\mathbf{x}, t) dt \right)^2 d\mathbf{x} \leq \sum_{n=1}^N \int_{t^{n-1}}^{t^n} \|\mathbf{f}(t)\|_\Omega^2 dt = \|\mathbf{f}\|_{L^2(\mathbf{L}^2(\Omega))}^2.$$

We estimate the second term in the right-hand side of (4.35) by splitting it into two contributions as follows:

$$\sum_{n=1}^N \tau (G^n + \phi^0, p_h^n)_\Omega = \sum_{n=1}^N \tau (G^n + \phi^0, p_h^n - \pi_\Omega^0 p_h^n)_\Omega + \sum_{n=1}^N \tau (\pi_\Omega^0 (G^n + \phi^0), p_h^n)_\Omega := \mathfrak{I}_1 + \mathfrak{I}_2,$$

where we have used its definition (4.10) to move  $\pi_\Omega^0$  from  $p_h^n$  to  $(G^n + \phi^0)$  in the second term. Owing to the Cauchy–Schwarz, triangle, Jensen, and Young inequalities, and using (4.33), we have

$$\begin{aligned} |\mathfrak{I}_1| &\leq \left( \sum_{n=1}^N \tau \|G^n + \phi^0\|_\Omega^2 \right)^{\frac{1}{2}} \left( \sum_{n=1}^N \tau \|p_h^n - \pi_\Omega^0 p_h^n\|_\Omega^2 \right)^{\frac{1}{2}} \\ &\leq \sqrt{2} \beta \left( \sum_{n=1}^N \tau \int_\Omega \left( \int_0^{t^n} g(\mathbf{x}, t) dt \right)^2 d\mathbf{x} + t_F \|\phi^0\|_\Omega^2 \right)^{\frac{1}{2}} \left( \sum_{n=1}^N \tau (C_K \|\bar{\mathbf{f}}^n\|_\Omega + C_{\text{gr}} C_{\text{eq}}^2 \|\underline{\mathbf{u}}_h^n\|_{\varepsilon,h})^2 \right)^{\frac{1}{2}} \\ &\leq 2\beta \left( t_F \|g\|_{L^2(\mathbf{L}^2(\Omega))} + t_F^{\frac{1}{2}} \|\phi^0\|_\Omega \right) \left[ C_K \|\mathbf{f}\|_{L^2(\mathbf{L}^2(\Omega))} + C_{\text{gr}} C_{\text{eq}}^2 \left( \sum_{n=1}^N \tau \|\underline{\mathbf{u}}_h^n\|_{\varepsilon,h}^2 \right)^{\frac{1}{2}} \right] \\ &\leq t_F \beta^2 \left( 1 + \frac{C_{\text{gr}}^2 C_{\text{eq}}^6}{C_{\text{cv}}^2} \right) \left( t_F^{\frac{1}{2}} \|g\|_{L^2(\mathbf{L}^2(\Omega))} + \|\phi^0\|_\Omega \right)^2 + C_K^2 \|\mathbf{f}\|_{L^2(\mathbf{L}^2(\Omega))}^2 + \frac{C_{\text{cv}}^2}{4C_{\text{eq}}^2} \sum_{n=1}^N \tau \|\underline{\mathbf{u}}_h^n\|_{\varepsilon,h}^2. \end{aligned} \quad (4.37)$$

Owing to the compatibility condition (4.1f) and the linearity of the  $L^2$ -projector,  $\mathfrak{I}_2 = 0$  if  $C_0 = 0$ . Otherwise, using again the Cauchy–Schwarz, triangle, Jensen, and Young

inequalities, leads to

$$\begin{aligned}
|\mathfrak{I}_2| &\leq \left( \sum_{n=1}^N \tau \|\pi_\Omega^0 G^n + \pi_\Omega^0 \phi^0\|_\Omega^2 \right)^{\frac{1}{2}} \left( \sum_{n=1}^N \tau \|p_h^n\|_\Omega^2 \right)^{\frac{1}{2}} \\
&\leq \sqrt{2} \left( t_F^2 \|\pi_\Omega^0 g\|_{L^2(L^2(\Omega))}^2 + t_F \|\pi_\Omega^0 \phi^0\|_\Omega^2 \right)^{\frac{1}{2}} \left( \sum_{n=1}^N \tau \|p_h^n\|_\Omega^2 \right)^{\frac{1}{2}} \\
&\leq \frac{2t_F^2}{3C_0} \|\pi_\Omega^0 g\|_{L^2(L^2(\Omega))}^2 + \frac{2t_F}{3C_0} \|\pi_\Omega^0 \phi^0\|_\Omega^2 + \frac{3C_0}{4} \sum_{n=1}^N \tau \|p_h^n\|_\Omega^2.
\end{aligned} \tag{4.38}$$

Finally, from (4.36), (4.37), (4.38), it follows that

$$\begin{aligned}
\mathcal{R} &\leq \frac{C_{cv}^2}{2C_{eq}^2} \sum_{n=1}^N \tau \|\underline{\mathbf{u}}_h^n\|_{\varepsilon,h}^2 + \frac{3C_0}{4} \sum_{n=1}^N \tau \|p_h^n\|_\Omega^2 + \frac{2t_F^2}{3C_0} \|\pi_\Omega^0 g\|_{L^2(L^2(\Omega))}^2 + \frac{2t_F}{3C_0} \|\pi_\Omega^0 \phi^0\|_\Omega^2 \\
&\quad + \frac{C_K^2}{C_{cv}^2} (C_{eq}^2 + C_{cv}^2) \|\mathbf{f}\|_{L^2(\mathbf{L}^2(\Omega))}^2 + \frac{t_F \beta^2}{C_{cv}^2} (4C_{gr}^2 C_{eq}^6 + C_{cv}^2) \left( t_F^{\frac{1}{2}} \|g\|_{L^2(L^2(\Omega))} + \|\phi^0\|_\Omega \right)^2.
\end{aligned} \tag{4.39}$$

(iv) *Conclusion.* Passing the first two terms in the right-hand side of (4.39) to the left-hand side of (4.35) and multiplying both sides by a factor 4, we obtain

$$\frac{2C_{cv}^2}{C_{eq}^2} \sum_{n=1}^N \tau \|\underline{\mathbf{u}}_h^n\|_{\varepsilon,h}^2 + C_0 \sum_{n=1}^N \tau \|p_h^n\|_\Omega^2 + 2\|s_h^N\|_{c,h}^2 \leq 4C, \tag{4.40}$$

where we have denoted by  $C$  the last four summands in the right-hand side of (4.39). In order to conclude we apply again (4.33) to obtain a bound of the  $L^2$ -norm of the discrete pressure independent of the storage coefficient  $C_0$ . Indeed, owing to (4.33) and  $C_{cv}^2 \leq C_{gr}$ , it is inferred that

$$\frac{C_{cv}^2}{2C_{gr}^2 C_{eq}^6} \sum_{n=1}^N \tau \|p_h^n - \pi_\Omega^0 p_h^n\|_\Omega^2 \leq \frac{C_{cv}^2}{C_{eq}^2} \sum_{n=1}^N \tau \|\underline{\mathbf{u}}_h^n\|_{\varepsilon,h}^2 + \frac{\beta^2 C_K^2}{C_{cv}^2 C_{eq}^6} \|\mathbf{f}\|_{L^2(\mathbf{L}^2(\Omega))}^2.$$

Summing the previous relation to (4.40) yields

$$\begin{aligned}
\frac{C_{cv}^2}{C_{eq}^2} \sum_{n=1}^N \tau \|\underline{\mathbf{u}}_h^n\|_{\varepsilon,h}^2 + \frac{C_{cv}^2}{2C_{gr}^2 C_{eq}^6} \sum_{n=1}^N \tau \|p_h^n - \pi_\Omega^0 p_h^n\|_\Omega^2 + C_0 \sum_{n=1}^N \tau \|p_h^n\|_\Omega^2 + 2\|s_h^N\|_{c,h}^2 \leq \\
\frac{8t_F^2}{3C_0} \|\pi_\Omega^0 g\|_{L^2(L^2(\Omega))}^2 + \frac{8t_F}{3C_0} \|\pi_\Omega^0 \phi^0\|_\Omega^2 + C_K^2 \left( \frac{4C_{eq}^2}{C_{cv}^2} + \frac{\beta^2}{C_{eq}^6 C_{cv}^2} + 4 \right) \|\mathbf{f}\|_{L^2(\mathbf{L}^2(\Omega))}^2 \\
+ 4t_F \beta^2 \left( \frac{4C_{gr}^2 C_{eq}^6}{C_{cv}^2} + 1 \right) \left( t_F^{\frac{1}{2}} \|g\|_{L^2(L^2(\Omega))} + \|\phi^0\|_\Omega \right)^2.
\end{aligned}$$

Thus, multiplying both sides of the previous relation by  $\max\{C_{eq}^2 C_{cv}^{-2}, 2C_{gr}^2 C_{eq}^6 C_{cv}^{-2}, 1\}$  gives (4.31).  $\square$



*Remark 4.13* (A priori bound for  $C_0 = 0$ ). When  $C_0 = 0$ , the a priori bound (4.31) reads

$$\sum_{n=1}^N \tau \|\underline{\mathbf{u}}_h^n\|_{\varepsilon,h}^2 + \sum_{n=1}^N \tau \|p_h^n\|_{\Omega}^2 + \|s_h^N\|_{c,h}^2 \leq C \left( \|\mathbf{f}\|_{L^2(\mathbf{L}^2(\Omega))}^2 + t_F^2 \|g\|_{L^2(L^2(\Omega))}^2 + t_F \|\phi^0\|_{\Omega}^2 \right).$$

The conventions  $C_0^{-1} \|\pi_{\Omega}^0 g\|_{L^2(L^2(\Omega))}^2 = 0$  and  $C_0^{-1} \|\pi_{\Omega}^0 \phi^0\|_{\Omega}^2 = 0$  if  $C_0 = 0$  are justified since the term  $\mathfrak{Z}_2$  in point (3) of the previous proof vanishes in this case thanks to the compatibility condition (4.1f).

We next proceed to discuss the existence and uniqueness of the discrete solutions. The proof of the following theorem hinges on the arguments of [64, Theorem 3.3].

**Theorem 4.14** (Existence and uniqueness). *Let Assumption 4.1 hold and let  $(\mathcal{M}_h)_{h \in \mathcal{H}}$  be a regular mesh sequence. Then, for all  $h \in \mathcal{H}$  and all  $N \in \mathbb{N}^*$ , there exists a unique solution  $(\underline{\mathbf{u}}_h^n, p_h^n)_{1 \leq n \leq N} \in (\underline{\mathbf{U}}_{h,D}^k \times P_h^k)^N$  to (4.30).*

*Proof.* We define the linear stress-strain function  $\sigma^{\text{lin}} : \mathbb{R}_s^{d \times d} \rightarrow \mathbb{R}_s^{d \times d}$  such that, for all  $\tau \in \mathbb{R}_s^{d \times d}$ ,

$$\sigma^{\text{lin}}(\tau) = \frac{C_{\text{cv}}^2}{2} \tau + \frac{C_{\text{cv}}^2}{2d} \text{tr}(\tau) \mathbf{I}_d,$$

where  $C_{\text{cv}}$  is the coercivity constant of  $\sigma$  (see (4.2b)), and we denote by  $a_h^{\text{lin}}$  the bilinear form obtained by replacing  $\sigma$  with  $\sigma^{\text{lin}}$  in (4.21). We consider the following auxiliary linear problem:c4: For all  $1 \leq n \leq N$ , find  $(\underline{\mathbf{y}}_h^n, p_h^n) \in \underline{\mathbf{U}}_{h,D}^k \times P_h^k$  such that

$$\begin{aligned} a_h^{\text{lin}}(\underline{\mathbf{y}}_h^n, \underline{\mathbf{v}}_h) + b_h(\underline{\mathbf{v}}_h, p_h^n) &= (\bar{\mathbf{f}}^n, \mathbf{v}_h)_{\Omega} & \forall \underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_{h,D}^k, \\ C_0(\delta_t p_h^n, q_h)_{\Omega} - b_h(\delta_t \underline{\mathbf{y}}_h^n, q_h) + c_h(p_h^n, q_h) &= (\bar{g}^n, q_h)_{\Omega} & \forall q_h \in P_h^k, \end{aligned} \quad (4.41)$$

with initial condition as in (4.30c). Since the previous system is linear and square and its solution satisfies the a priori estimate of Proposition 4.11, it is readily inferred that problem (4.41) admits a unique solution.

Now we observe that, thanks to the norm equivalence (4.22),  $a_h^{\text{lin}}(\cdot, \cdot)$  is a scalar product on  $\underline{\mathbf{U}}_{h,D}^k$ , and we define the mapping  $\Phi : \underline{\mathbf{U}}_{h,D}^k \rightarrow \underline{\mathbf{U}}_{h,D}^k$  such that, for all  $\underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_{h,D}^k$ ,

$$a_h^{\text{lin}}(\Phi(\underline{\mathbf{v}}_h), \underline{\mathbf{w}}_h) = a_h(\underline{\mathbf{v}}_h, \underline{\mathbf{w}}_h), \quad \forall \underline{\mathbf{w}}_h \in \underline{\mathbf{U}}_{h,D}^k.$$

We want to show that  $\Phi$  is an isomorphism. Let  $\underline{\mathbf{v}}_h, \underline{\mathbf{z}}_h \in \underline{\mathbf{U}}_{h,D}^k$  be such that  $\Phi(\underline{\mathbf{v}}_h) = \Phi(\underline{\mathbf{z}}_h)$ . If  $\underline{\mathbf{v}}_h \neq \underline{\mathbf{z}}_h$ , owing to the norm equivalence (4.22) and the fact that  $\|\cdot\|_{\varepsilon,h}$  is a norm on  $\underline{\mathbf{U}}_{h,D}^k$ , there is at least one  $T \in \mathcal{T}_h$  such that  $\mathbf{G}_{s,T}^k \underline{\mathbf{v}}_T \neq \mathbf{G}_{s,T}^k \underline{\mathbf{z}}_T$  or  $\Delta_{TF}^k \underline{\mathbf{v}}_T \neq \Delta_{TF}^k \underline{\mathbf{z}}_T$  for some  $F \in \mathcal{F}_T$ . In both cases, owing to the definition of  $a_h$  and the strict monotonicity assumption (4.2c), it holds

$$0 < a_h(\underline{\mathbf{v}}_h, \underline{\mathbf{v}}_h - \underline{\mathbf{z}}_h) - a_h(\underline{\mathbf{z}}_h, \underline{\mathbf{v}}_h - \underline{\mathbf{z}}_h) = a_h^{\text{lin}}(\Phi(\underline{\mathbf{v}}_h) - \Phi(\underline{\mathbf{z}}_h), \underline{\mathbf{v}}_h - \underline{\mathbf{z}}_h) = 0.$$

Thus, we infer by contradiction that  $\underline{\mathbf{v}}_h = \underline{\mathbf{z}}_h$  and, as a result,  $\Phi$  is injective. In order to prove that  $\Phi$  is also onto, we recall the following result: If  $(E, (\cdot, \cdot)_E)$  is a Euclidean space

and  $\Psi : E \rightarrow E$  is a continuous map such that  $\frac{(\Psi(x), x)_E}{\|x\|_E} \rightarrow +\infty$  as  $\|x\|_E \rightarrow +\infty$ , then  $\Psi$  is surjective. Since  $(\underline{U}_{h,D}^k, a_h^{\text{lin}}(\cdot, \cdot))$  is a Euclidean space and the coercivity (4.2b) of  $\sigma$  together with the definition of  $\sigma^{\text{lin}}$  yields  $a_h^{\text{lin}}(\Phi(\underline{y}_h), \underline{y}_h) \geq a_h^{\text{lin}}(\underline{y}_h, \underline{y}_h)$  for all  $\underline{y}_h \in \underline{U}_{h,D}^k$ , we deduce that  $\Phi$  is an isomorphism. Let, for all  $1 \leq n \leq N$ ,  $(\underline{y}_h^n, p_h^n) \in \underline{U}_{h,D}^k \times P_h^k$  be the unique solution to problem (4.41). By the surjectivity and injectivity of  $\Phi$ , for all  $1 \leq n \leq N$ , there exists a unique  $\underline{u}_h^n \in \underline{U}_{h,D}^k$  such that  $\Phi(\underline{u}_h^n) = (\underline{y}_h^n, p_h^n)$ . By definition of  $\Phi$  and  $(\underline{y}_h^n)_{1 \leq n \leq N}$ ,  $(\underline{u}_h^n, p_h^n)_{1 \leq n \leq N}$  is therefore the unique solution of the discrete problem (4.30).  $\square$

## 4.6 Convergence analysis

In this section we study the convergence of problem (4.30) and prove optimal error estimates under the following additional assumptions on the stress-strain function  $\sigma$ .

**Assumption 4.15** (Stress-strain relation II). There exist real numbers  $C_{\text{lp}}, C_{\text{mn}} \in (0, +\infty)$  such that, for a.e.  $\mathbf{x} \in \Omega$ , and all  $\boldsymbol{\tau}, \boldsymbol{\eta} \in \mathbb{R}_s^{d \times d}$ ,

$$|\sigma(\mathbf{x}, \boldsymbol{\tau}) - \sigma(\mathbf{x}, \boldsymbol{\eta})|_{d \times d} \leq C_{\text{lp}} |\boldsymbol{\tau} - \boldsymbol{\eta}|_{d \times d}, \quad (\text{Lipschitz continuity}) \quad (4.42a)$$

$$(\sigma(\mathbf{x}, \boldsymbol{\tau}) - \sigma(\mathbf{x}, \boldsymbol{\eta})) : (\boldsymbol{\tau} - \boldsymbol{\eta}) \geq C_{\text{mn}}^2 |\boldsymbol{\tau} - \boldsymbol{\eta}|_{d \times d}^2. \quad (\text{strong monotonicity}) \quad (4.42b)$$

*Remark 4.16* (Lipschitz continuity and strong monotonicity). It is readily seen, by taking  $\boldsymbol{\eta} = \mathbf{0}$  in (4.42), that Lipschitz continuity and strong monotonicity imply respectively the growth and coercivity properties of Assumption 4.1. Therefore, recalling (4.23), it is inferred that the constants appearing in (4.2a), (4.2b), (4.42a), and (4.42b) satisfy

$$C_{\text{mn}}^2 \leq C_{\text{cv}}^2 \leq C_{\text{gr}} \leq C_{\text{lp}}. \quad (4.43)$$

It was proved in [13, Lemma 4.1] that the stress-strain relation for the Hencky–Mises model is strongly monotone and Lipschitz-continuous. Also the isotropic damage model satisfies Assumption 4.15 if the damage function in (4.5) is, for instance, such that

$$D(\mathbf{x}, |\boldsymbol{\tau}|) = 1 - (1 + |\mathfrak{C}(\mathbf{x})\boldsymbol{\tau}|_{d \times d})^{-\frac{1}{2}} \quad \forall \mathbf{x} \in \Omega.$$

In order to prove a convergence rate of  $(k + 1)$  in space for both the displacement and pressure errors, we assume from this point on that the permeability tensor field  $\boldsymbol{\kappa}$  is constant on  $\Omega$ , and that the following elliptic regularity holds (which is the case, e.g., when  $\Omega$  is convex [115, 149]): There is a real number  $C_{\text{el}} > 0$  only depending on  $\Omega$  such that, for all  $\psi \in L_0^2(\Omega)$ , the unique function  $\zeta \in P$  solution of the homogeneous Neumann problem

$$-\nabla \cdot (\boldsymbol{\kappa} \nabla \zeta) = \psi \quad \text{in } \Omega, \quad \boldsymbol{\kappa} \nabla \zeta \cdot \mathbf{n} = 0 \quad \text{on } \partial\Omega,$$

is such that

$$\|\zeta\|_{H^2(\Omega)} \leq C_{\text{el}} \boldsymbol{\kappa}^{-\frac{1}{2}} \|\psi\|_{\Omega}. \quad (4.44)$$

Let  $(\underline{u}_h^n, p_h^n)_{1 \leq n \leq N}$  be the solution to (4.30). We consider, for all  $1 \leq n \leq N$ , the discrete error components defined as

$$\underline{e}_h^n := \underline{u}_h^n - \underline{I}_h^k \bar{\mathbf{u}}^n, \quad \epsilon_h^n := p_h^n - \widehat{p}_h^n, \quad (4.45)$$

where the global elliptic projection  $\widehat{p}_h^n \in P_h^k$  is defined as the solution to

$$c_h(\widehat{p}_h^n, q_h) = c_h(\bar{p}^n, q_h) \quad \forall q_h \in P_h^k \quad \text{and} \quad \int_{\Omega} \widehat{p}_h^n = \int_{\Omega} \bar{p}^n.$$

Before proving the convergence of the scheme, we recall two preliminary approximation results for the projector  $\underline{I}_h^k$  and the projection  $\widehat{p}_h^n$  that have been proved in [33, Theorem 16] and [30, Lemma 11], respectively. There is a strictly positive constant  $C_{\text{pj}}$  depending only on  $\Omega$ ,  $k$ , and the mesh regularity parameter, such that,

- Assuming (4.42) and  $\mathbf{u} \in L^2(\mathbf{U} \cap \mathbf{H}^{k+2}(\mathcal{T}_h))$  with  $\boldsymbol{\sigma}(\cdot, \nabla_s \mathbf{u}) \in L^2(\mathbb{H}_s^{k+1}(\mathcal{T}_h))$ , for a.e.  $t \in (0, t_F)$  and all  $\underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_{h,D}^k$ , it holds

$$\begin{aligned} & \left| a_h(\underline{I}_h^k \mathbf{u}(\cdot, t), \underline{\mathbf{v}}_h) + (\nabla \cdot \boldsymbol{\sigma}(\cdot, \nabla_s \mathbf{u}(\cdot, t)), \mathbf{v}_h) \right| \leq \\ & C_{\text{pj}} h^{k+1} \left( |\mathbf{u}(\cdot, t)|_{\mathbf{H}^{k+2}(\mathcal{T}_h)} + |\boldsymbol{\sigma}(\cdot, \nabla_s \mathbf{u}(\cdot, t))|_{\mathbb{H}^{k+1}(\mathcal{T}_h)} \right) \|\underline{\mathbf{v}}_h\|_{\boldsymbol{\varepsilon}, h}. \end{aligned} \quad (4.46)$$

- Assuming the elliptic regularity (4.44) and  $\bar{p}^n \in P \cap H^{k+1}(\mathcal{T}_h)$ , for all  $1 \leq n \leq N$ , it holds

$$h \|\widehat{p}_h^n - \bar{p}^n\|_{c,h} + \underline{\kappa}^{\frac{1}{2}} \|\widehat{p}_h^n - \bar{p}^n\|_{\Omega} \leq C_{\text{pj}} h^{k+1} \bar{\kappa}^{\frac{1}{2}} |\bar{p}^n|_{H^{k+1}(\mathcal{T}_h)}. \quad (4.47)$$

Now we have all the ingredients to estimate the discrete errors defined in (4.45).

**Theorem 4.17** (Error estimate). *Let  $(\mathbf{u}, p)$  denote the unique solution to (4.6), for which we assume*

$$\begin{aligned} \mathbf{u} & \in H^1(\mathcal{T}_\tau; \mathbf{U}) \cap L^2(\mathbf{H}^{k+2}(\mathcal{T}_h)), & \boldsymbol{\sigma}(\cdot, \nabla_s \mathbf{u}) & \in L^2(\mathbb{H}_s^{k+1}(\mathcal{T}_h)), \\ p & \in L^2(P \cap H^{k+1}(\mathcal{T}_h)), & \phi & \in H^1(\mathcal{T}_\tau; L^2(\Omega)), \end{aligned}$$

with  $\phi = C_0 p + \nabla \cdot \mathbf{u}$ . If  $C_0 > 0$ , we further assume  $\pi_{\Omega}^0 p \in H^1(\mathcal{T}_\tau; P^0(\Omega)) =: H^1(\mathcal{T}_\tau)$ . Then, under Assumption 4.15 and the elliptic regularity (4.44), it holds

$$\sum_{n=1}^N \tau \|\underline{\boldsymbol{\varepsilon}}_h^n\|_{\boldsymbol{\varepsilon}, h}^2 + \sum_{n=1}^N \tau \left( \|\epsilon_h^n - \pi_{\Omega}^0 \epsilon_h^n\|_{\Omega}^2 + C_0 \|\epsilon_h^n\|_{\Omega}^2 \right) + \|z_h^N\|_{c,h}^2 \leq C \left( h^{2k+2} C_1 + \tau^2 C_2 \right), \quad (4.48)$$

where  $C$  is a strictly positive constant independent of  $h$ ,  $\tau$ ,  $C_0$ ,  $\kappa$ , and  $t_F$ , and, for the sake of brevity, we have defined  $z_h^N := \sum_{n=1}^N \tau \epsilon_h^n$  and introduced the bounded quantities

$$\begin{aligned} C_1 & := |\mathbf{u}|_{L^2(\mathbf{H}^{k+2}(\mathcal{T}_h))}^2 + |\boldsymbol{\sigma}(\cdot, \nabla_s \mathbf{u})|_{L^2(\mathbb{H}^{k+1}(\mathcal{T}_h))}^2 + (1 + C_0) \frac{\bar{\kappa}}{\underline{\kappa}} |p|_{L^2(H^{k+1}(\mathcal{T}_h))}^2, \\ C_2 & := \|\mathbf{u}\|_{H^1(\mathcal{T}_\tau; \mathbf{H}^1(\Omega))}^2 + \|\phi\|_{H^1(\mathcal{T}_\tau; L^2(\Omega))}^2 + C_0 \|\pi_{\Omega}^0 p\|_{H^1(\mathcal{T}_\tau)}^2. \end{aligned}$$

**Remark 4.18** (Time regularity). In order to prove the previous error estimate, we only require the displacement  $\mathbf{u}$  and the fluid content  $\phi$  solving problem (4.6) to be piecewise  $H^1$ -regular in  $(0, t_F)$ , whereas 2.12 is established under the much stronger regularity  $\mathbf{u} \in C^2(\mathbf{H}^1(\Omega))$  and, if  $C_0 > 0$ ,  $p \in C^2(L^2(\Omega))$ . Moreover, the assumptions  $\phi \in H^1(\mathcal{T}_\tau; L^2(\Omega))$  and, if  $C_0 > 0$ ,  $\pi_{\Omega}^0 p \in H^1(\mathcal{T}_\tau)$  are consistent with the time regularity results observed in Remark 4.2.

*Proof.* (i) *Estimate of  $\|\underline{e}_h^n\|_{\varepsilon,h}^2$ .* First we observe that, owing to (4.1a) and the definition of  $b_h$  given in (4.25), for all  $\underline{v}_h \in \underline{U}_{h,D}^k$  and all  $1 \leq n \leq N$ , we have

$$(\bar{\mathbf{f}}^n, \mathbf{v}_h)_\Omega = -(\nabla \cdot \bar{\boldsymbol{\sigma}}^n(\cdot, \nabla_s \mathbf{u}), \mathbf{v}_h)_\Omega + (\nabla \bar{p}^n, \mathbf{v}_h)_\Omega = a_h(\mathbf{I}_h^k \bar{\mathbf{u}}^n, \underline{v}_h) + b_h(\underline{v}_h, \bar{p}_h^n) - \mathcal{R}^n(\underline{v}_h), \quad (4.49)$$

where the residual linear form  $\mathcal{R}^n : \underline{U}_{h,D}^k \rightarrow \mathbb{R}$  is defined such that, for all  $\underline{v}_h \in \underline{U}_{h,D}^k$ ,

$$\mathcal{R}^n(\underline{v}_h) := a_h(\mathbf{I}_h^k \bar{\mathbf{u}}^n, \underline{v}_h) + (\nabla \cdot \bar{\boldsymbol{\sigma}}^n(\cdot, \nabla_s \mathbf{u}), \mathbf{v}_h)_\Omega + b_h(\underline{v}_h, \bar{p}_h^n) - (\nabla \bar{p}^n, \mathbf{v}_h)_\Omega. \quad (4.50)$$

Using the norm equivalence (4.22), the strong monotonicity (4.42b) of  $\boldsymbol{\sigma}$  along with assumption (4.24) on the stabilization parameter, the discrete mechanical equilibrium (4.30a), and (4.49), yields

$$\begin{aligned} \frac{C_{mn}^2}{C_{eq}^2} \|\underline{e}_h^n\|_{\varepsilon,h}^2 &= C_{mn}^2 \sum_{T \in \mathcal{T}_h} \left( \|\mathbf{G}_{s,T}^k \underline{\mathbf{e}}_T^n\|_T^2 + \sum_{F \in \mathcal{F}_T} \frac{1}{h_F} \|\mathbf{A}_{TF}^k \underline{\mathbf{e}}_T^n\|_F^2 \right) \\ &\leq a_h(\underline{\mathbf{u}}_h^n, \underline{\mathbf{e}}_h^n) - a_h(\mathbf{I}_h^k \bar{\mathbf{u}}^n, \underline{\mathbf{e}}_h^n) \\ &= (\bar{\mathbf{f}}^n, \underline{\mathbf{e}}_h^n)_\Omega - b_h(\underline{\mathbf{e}}_h^n, p_h^n) - a_h(\mathbf{I}_h^k \bar{\mathbf{u}}^n, \underline{\mathbf{e}}_h^n) = -b_h(\underline{\mathbf{e}}_h^n, \epsilon_h^n) - \mathcal{R}^n(\underline{\mathbf{e}}_h^n). \end{aligned}$$

Thus, owing to the previous relation and defining the dual norm

$$\|\mathcal{R}^n\|_{\varepsilon,h,*} := \sup_{\underline{v}_h \in \underline{U}_{h,D}^k \setminus \{\mathbf{0}\}} \frac{\mathcal{R}^n(\underline{v}_h^n)}{\|\underline{v}_h^n\|_{\varepsilon,h}},$$

we have that

$$\frac{C_{mn}^2}{C_{eq}^2} \|\underline{e}_h^n\|_{\varepsilon,h}^2 + b_h(\underline{\mathbf{e}}_h^n, \epsilon_h^n) \leq \|\mathcal{R}^n\|_{\varepsilon,h,*} \|\underline{e}_h^n\|_{\varepsilon,h} \leq \frac{C_{eq}^2}{2C_{mn}^2} \|\mathcal{R}^n\|_{\varepsilon,h,*}^2 + \frac{C_{mn}^2}{2C_{eq}^2} \|\underline{e}_h^n\|_{\varepsilon,h}^2,$$

where the conclusion follows from Young's inequality. Hence, rearranging, we arrive at

$$\frac{C_{mn}^2}{2C_{eq}^2} \|\underline{e}_h^n\|_{\varepsilon,h}^2 + b_h(\underline{\mathbf{e}}_h^n, \epsilon_h^n) \leq \frac{C_{eq}^2}{2C_{mn}^2} \|\mathcal{R}^n\|_{\varepsilon,h,*}^2 \quad (4.51)$$

(ii) *Estimate of  $C_0 \|\epsilon_h^n\|_\Omega^2$ .* Using (4.1b), the fact that  $(\nabla \cdot \mathbf{u}(t), 1)_\Omega = 0$  to insert  $\pi_\Omega^0 q_h$ , and the consistency property (4.29), we infer that, for all  $q_h \in P_h^k$  and all  $1 \leq i \leq N$ ,

$$\begin{aligned} (\bar{g}^i, q_h)_\Omega &= (C_0 \bar{\mathbf{d}}_t^i, q_h)_\Omega + (\nabla \cdot (\bar{\mathbf{d}}_t^i), q_h - \pi_\Omega^0 q_h)_\Omega - (\nabla \cdot (\boldsymbol{\kappa} \nabla \bar{p}^i), q_h)_\Omega \\ &= \tau^{-1} \int_{t^{i-1}}^{t^i} dt \left[ C_0(p(t), q_h)_\Omega + (\nabla \cdot \mathbf{u}(t), q_h - \pi_\Omega^0 q_h)_\Omega \right] dt + c_h(\bar{p}^i, q_h) \quad (4.52) \\ &= \delta_t \left[ C_0(p^i, q_h)_\Omega + (\nabla \cdot \mathbf{u}^i, q_h - \pi_\Omega^0 q_h)_\Omega \right] + c_h(\bar{p}^i, q_h). \end{aligned}$$

Therefore, using the discrete mass conservation equation (4.30b), (4.26), the Fortin property (4.27b), the definition of the elliptic projection  $\tilde{p}_h^n$ , and (4.52), we obtain

$$\begin{aligned}
& C_0(\delta_t \epsilon_h^i, q_h)_\Omega - b_h(\delta_t \underline{e}_h^i, q_h) + c_h(\epsilon_h^i, q_h) \\
&= (\bar{g}^i, q_h)_\Omega - C_0(\delta_t \tilde{p}_h^i, q_h)_\Omega + b_h(\delta_t(\underline{\mathbf{I}}_h^k \bar{\mathbf{u}}^i), q_h - \pi_\Omega^0 q_h) - c_h(\tilde{p}_h^i, q_h) \\
&= (\bar{g}^i, q_h)_\Omega - \delta_t \left[ C_0(\tilde{p}_h^i, q_h)_\Omega - (\nabla \cdot \bar{\mathbf{u}}^i, q_h - \pi_\Omega^0 q_h)_\Omega \right] - c_h(\tilde{p}_h^i, q_h) \\
&= \delta_t \left[ C_0(p^i - \tilde{p}_h^i, q_h)_\Omega + (\nabla \cdot \mathbf{u}^i - \nabla \cdot \bar{\mathbf{u}}^i, q_h - \pi_\Omega^0 q_h)_\Omega \right] \\
&= \delta_t \left[ C_0(\bar{p}^i - \tilde{p}_h^i, q_h)_\Omega + C_0(\pi_\Omega^0(p^i - \bar{p}^i), q_h)_\Omega + (\phi^i - \bar{\phi}^i, q_h - \pi_\Omega^0 q_h)_\Omega \right],
\end{aligned} \tag{4.53}$$

where, in order to pass to the last line, we have inserted  $\pm \bar{p}^i$  into the first term inside brackets in the third line, we have defined, according to (4.1e),  $\phi^i := C_0 p^i + \nabla \cdot \mathbf{u}^i$  for all  $0 \leq i \leq N$ , and we have used the definition of the global  $L^2$ -projector  $\pi_\Omega^0$ . Moreover, setting  $\tilde{p}_h^0 := 0$ , it follows from the initial condition (4.30c), the boundary condition (4.1c), and (4.26) that

$$C_0(\epsilon_h^0, q_h)_\Omega - b_h(\underline{e}_h^0, q_h) = (\phi^0, q_h)_\Omega = (\phi^0, q_h - \pi_\Omega^0 q_h)_\Omega + C_0(\pi_\Omega^0 p^0, q_h)_\Omega. \tag{4.54}$$

For all  $1 \leq n \leq N$ , summing (4.53) for  $1 \leq i \leq n$  with the choice  $q_h = \tau \epsilon_h^n$ , using (4.54), and proceeding as in the second step of the proof of Proposition 4.11, leads to

$$\begin{aligned}
& C_0 \|\epsilon_h^n\|_\Omega^2 - b_h(\underline{e}_h^n, \epsilon_h^n) + \frac{1}{2\tau} \left( \|z_h^n\|_{c,h}^2 - \|z_h^{n-1}\|_{c,h}^2 \right) \\
& \leq C_0(\bar{p}^n - \tilde{p}_h^n, \epsilon_h^n)_\Omega + C_0(\pi_\Omega^0(p^n - \bar{p}^n), \epsilon_h^n)_\Omega + (\phi^n - \bar{\phi}^n, \epsilon_h^n - \pi_\Omega^0 \epsilon_h^n)_\Omega,
\end{aligned} \tag{4.55}$$

where  $z_h^n := \sum_{i=1}^n \tau \epsilon_h^i$  if  $n \geq 1$  and  $z_h^0 := 0$ . We bound the first term in the right-hand side of (4.55) applying the Cauchy–Schwarz and Young inequalities followed by the approximation result (4.47), yielding

$$\begin{aligned}
C_0(\bar{p}^n - \tilde{p}_h^n, \epsilon_h^n)_\Omega & \leq C_{\text{pj}}^2 h^{2(k+1)} \frac{\bar{K}}{K} C_0 |\bar{p}^n|_{H^{k+1}(\mathcal{T}_h)}^2 + \frac{C_0}{4} \|\epsilon_h^n\|_\Omega^2 \\
& \leq C_{\text{pj}}^2 \frac{h^{2(k+1)}}{\tau} \left( \frac{\bar{K}}{K} \right) C_0 |p|_{L^2((t^{n-1}, t^n); H^{k+1}(\mathcal{T}_h))}^2 + \frac{C_0}{4} \|\epsilon_h^n\|_\Omega^2,
\end{aligned} \tag{4.56}$$

where, in order to pass to the second line, we have used the Cauchy–Schwarz inequality and adopted the notation  $|\cdot|_{L^2((t^{n-1}, t^n); H^m(\mathcal{T}_h))} := \|\cdot\|_{H^m(\mathcal{T}_h)} \| \cdot \|_{L^2((t^{n-1}, t^n))}$ , for any  $m \in \mathbb{N}$ . We estimate the second and third terms using the Cauchy–Schwarz and the Young inequalities together with the time approximation result (4.9) as follows:

$$\begin{aligned}
C_0(\pi_\Omega^0(p^n - \bar{p}^n), \epsilon_h^n)_\Omega & \leq C_0 \tau \|\pi_\Omega^0 p\|_{H^1((t^{n-1}, t^n))}^2 + \frac{C_0}{4} \|\epsilon_h^n\|_\Omega^2, \\
(\phi^n - \bar{\phi}^n, \epsilon_h^n - \pi_\Omega^0 \epsilon_h^n)_\Omega & \leq \eta \tau \|\phi\|_{H^1((t^{n-1}, t^n); L^2(\Omega))}^2 + \frac{1}{4\eta} \|\epsilon_h^n - \pi_\Omega^0 \epsilon_h^n\|_\Omega^2,
\end{aligned} \tag{4.57}$$

with  $\eta$  denoting a positive real number that will be fixed later on in the proof. The relation obtained by plugging (4.56) and (4.57) into (4.55) reads

$$\begin{aligned} & \frac{C_0}{2} \|\epsilon_h^n\|_{\Omega}^2 - b_h(\underline{e}_h^n, \epsilon_h^n) + \frac{1}{2\tau} \left( \|z_h^n\|_{c,h}^2 - \|z_h^{n-1}\|_{c,h}^2 \right) - \frac{1}{4\eta} \|\epsilon_h^n - \pi_{\Omega}^0 \epsilon_h^n\|_{\Omega}^2 \leq \\ & C_{\text{pj}}^2 \frac{h^{2(k+1)}}{\tau} \left( \frac{\bar{K}}{K} \right) C_0 |p|_{L^2((t^{n-1}, t^n); H^{k+1}(\mathcal{T}_h))}^2 + C_0 \tau \|\pi_{\Omega}^0 p\|_{H^1((t^{n-1}, t^n))}^2 + \eta \tau \|\phi\|_{H^1((t^{n-1}, t^n); L^2(\Omega))}^2. \end{aligned} \quad (4.58)$$

(iii) *Estimate of  $\|\epsilon_h^n - \pi_{\Omega}^0 \epsilon_h^n\|_{\Omega}^2$ .* We proceed as in the first step of the proof of Proposition 4.11. Using the inf-sup condition (4.28), (4.26) followed by the definition (4.45) of the pressure error, the linearity of  $b_h$ , the mechanical equilibrium equation (4.30a), and (4.49), we get, for all  $1 \leq n \leq N$ ,

$$\begin{aligned} \|\epsilon_h^n - \pi_{\Omega}^0 \epsilon_h^n\|_{\Omega} & \leq \beta \sup_{\underline{v}_h \in \underline{U}_{h,D}^k \setminus \{\mathbf{0}\}} \frac{b_h(\underline{v}_h, \epsilon_h^n - \pi_{\Omega}^0 \epsilon_h^n)}{\|\underline{v}_h\|_{\varepsilon,h}} \\ & = \beta \sup_{\underline{v}_h \in \underline{U}_{h,D}^k \setminus \{\mathbf{0}\}} \frac{b_h(\underline{v}_h, p_h^n - \widehat{p}_h^n)}{\|\underline{v}_h\|_{\varepsilon,h}} \\ & = \beta \sup_{\underline{v}_h \in \underline{U}_{h,D}^k \setminus \{\mathbf{0}\}} \frac{(\underline{f}^n, \underline{v}_h)_{\Omega} - a_h(\underline{u}_h^n, \underline{v}_h) - b_h(\underline{v}_h, \widehat{p}_h^n)}{\|\underline{v}_h\|_{\varepsilon,h}} \\ & = \beta \sup_{\underline{v}_h \in \underline{U}_{h,D}^k \setminus \{\mathbf{0}\}} \frac{a_h(\underline{I}_h^k \overline{\underline{u}}^n, \underline{v}_h) - a_h(\underline{u}_h^n, \underline{v}_h) - \mathcal{R}^n(\underline{v}_h)}{\|\underline{v}_h\|_{\varepsilon,h}}. \end{aligned} \quad (4.59)$$

Moreover, the Lipschitz continuity of the stress-strain function (4.42a), the Cauchy–Schwarz inequality, assumption (4.24) on the stabilization parameter  $\gamma$  together with (4.43), and the second inequality in (4.22), lead to

$$\begin{aligned} & a_h(\underline{I}_h^k \overline{\underline{u}}^n, \underline{v}_h) - a_h(\underline{u}_h^n, \underline{v}_h) \\ & = \sum_{T \in \mathcal{T}_h} \left( \int_T (\sigma(\cdot, \mathbf{G}_{s,T}^k \underline{I}_T^k \overline{\underline{u}}^n) - \sigma(\cdot, \mathbf{G}_{s,T}^k \underline{u}_T^n)) : \mathbf{G}_{s,T}^k \underline{v}_T + \sum_{F \in \mathcal{F}_T} \frac{\gamma}{h_F} \int_F \underline{\Delta}_{TF}^k (\underline{I}_T^k \overline{\underline{u}}^n - \underline{u}_T^n) \cdot \underline{\Delta}_{TF}^k \underline{v}_T \right) \\ & \leq C_{\text{lp}} C_{\text{eq}}^2 \|\underline{e}_h^n\|_{\varepsilon,h} \|\underline{v}_h\|_{\varepsilon,h}. \end{aligned} \quad (4.60)$$

Therefore, plugging the previous bound into the last line of (4.59), yields

$$\|\epsilon_h^n - \pi_{\Omega}^0 \epsilon_h^n\|_{\Omega} \leq \beta C_{\text{lp}} C_{\text{eq}}^2 \|\underline{e}_h^n\|_{\varepsilon,h} + \beta \|\mathcal{R}^n\|_{\varepsilon,h,*}.$$

Squaring and rearranging the previous relation and recalling that, owing to (4.43),  $C_{\text{mn}}^2 \leq C_{\text{lp}}$ , it is inferred that

$$\frac{\|\epsilon_h^n - \pi_{\Omega}^0 \epsilon_h^n\|_{\Omega}^2}{2\beta^2 C_{\text{mn}}^{-2} C_{\text{lp}}^2 C_{\text{eq}}^6} \leq \frac{C_{\text{mn}}^2}{C_{\text{eq}}^2} \|\underline{e}_h^n\|_{\varepsilon,h}^2 + \frac{\|\mathcal{R}^n\|_{\varepsilon,h,*}^2}{C_{\text{mn}}^2 C_{\text{eq}}^6}. \quad (4.61)$$

(iv) *Estimate of the dual norm of the residual.* We split the residual linear form  $\mathcal{R}^n$  defined in (4.50) into three contributions  $\mathcal{R}^n := \mathcal{R}_1^n + \mathcal{R}_2^n + \mathcal{R}_3^n$ , defined, for all  $\underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_{h,D}^k$  and all  $1 \leq n \leq N$ , such that

$$\mathcal{R}_1^n(\underline{\mathbf{v}}_h) := a_h(\underline{\mathbf{I}}_h^k \bar{\mathbf{u}}^n, \underline{\mathbf{v}}_h) - \frac{1}{\tau} \int_{t^{n-1}}^{t^n} a_h(\underline{\mathbf{I}}_h^k \mathbf{u}(t), \underline{\mathbf{v}}_h) dt, \quad (4.62a)$$

$$\mathcal{R}_2^n(\underline{\mathbf{v}}_h) := (\nabla \cdot \bar{\boldsymbol{\sigma}}^n(\cdot, \nabla_s \mathbf{u}), \mathbf{v}_h)_\Omega + \frac{1}{\tau} \int_{t^{n-1}}^{t^n} a_h(\underline{\mathbf{I}}_h^k \mathbf{u}(t), \underline{\mathbf{v}}_h) dt, \quad (4.62b)$$

$$\mathcal{R}_3^n(\underline{\mathbf{v}}_h) := b_h(\underline{\mathbf{v}}_h, \bar{p}_h^n) - (\nabla \bar{p}^n, \mathbf{v}_h)_\Omega. \quad (4.62c)$$

The first contribution can be bounded proceeding as in (4.60), then using the stability property (4.27a) of the interpolator  $\underline{\mathbf{I}}_h^k$ , the Cauchy–Schwarz inequality, and a Poincaré–Wirtinger inequality on the time interval  $(t^{n-1}, t^n)$ . By doing so, we get

$$\begin{aligned} \mathcal{R}_1^n(\underline{\mathbf{v}}_h) &= \frac{1}{\tau} \int_{t^{n-1}}^{t^n} a_h(\underline{\mathbf{I}}_h^k \bar{\mathbf{u}}^n, \underline{\mathbf{v}}_h) - a_h(\underline{\mathbf{I}}_h^k \mathbf{u}(t), \underline{\mathbf{v}}_h) dt \\ &\leq \frac{1}{\tau} \int_{t^{n-1}}^{t^n} (C_{\text{Ip}} C_{\text{eq}}^2 \|\underline{\mathbf{I}}_h^k(\bar{\mathbf{u}}^n - \mathbf{u}(t))\|_{\varepsilon,h} \|\underline{\mathbf{v}}_h\|_{\varepsilon,h}) dt \\ &\leq \frac{C_{\text{Ip}} C_{\text{eq}}^2 C_{\text{st}}}{\tau} \|\underline{\mathbf{v}}_h\|_{\varepsilon,h} \int_{t^{n-1}}^{t^n} \|\bar{\mathbf{u}}^n - \mathbf{u}(t)\|_{1,\Omega} dt \\ &\leq \frac{C_{\text{Ip}} C_{\text{eq}}^2 C_{\text{st}}}{\sqrt{\tau}} \|\underline{\mathbf{v}}_h\|_{\varepsilon,h} \|\bar{\mathbf{u}}^n - \mathbf{u}\|_{L^2((t^{n-1}, t^n); \mathbf{H}^1(\Omega))} \\ &\leq C_{\text{Ip}} C_{\text{eq}}^2 C_{\text{st}} C_{\text{ap}} \sqrt{\tau} \|\underline{\mathbf{v}}_h\|_{\varepsilon,h} \|\mathbf{u}\|_{H^1((t^{n-1}, t^n); \mathbf{H}^1(\Omega))}. \end{aligned} \quad (4.63)$$

We estimate the residual linear form  $\mathcal{R}_2^n$  defined in (4.62b) using the consistency property (4.46) and the Cauchy–Schwarz inequality, obtaining

$$\begin{aligned} \mathcal{R}_2^n(\underline{\mathbf{v}}_h) &= \frac{1}{\tau} \int_{t^{n-1}}^{t^n} (\nabla \cdot \boldsymbol{\sigma}(\cdot, \nabla_s \mathbf{u}(t)), \underline{\mathbf{v}}_h)_\Omega + a_h(\underline{\mathbf{I}}_h^k \mathbf{u}(t), \underline{\mathbf{v}}_h) dt \\ &\leq C_{\text{Pj}} \frac{h^{k+1}}{\tau} \|\underline{\mathbf{v}}_h\|_{\varepsilon,h} \int_{t^{n-1}}^{t^n} (|\mathbf{u}(t)|_{\mathbf{H}^{k+2}(\mathcal{T}_h)} + |\boldsymbol{\sigma}(\cdot, \nabla_s \mathbf{u}(t))|_{\mathbb{H}^{k+1}(\mathcal{T}_h)}) dt \\ &\leq C_{\text{Pj}} \frac{h^{k+1}}{\sqrt{\tau}} \|\underline{\mathbf{v}}_h\|_{\varepsilon,h} (|\mathbf{u}|_{L^2((t^{n-1}, t^n); \mathbf{H}^{k+2}(\mathcal{T}_h))} + |\boldsymbol{\sigma}(\cdot, \nabla_s \mathbf{u})|_{L^2((t^{n-1}, t^n); \mathbb{H}^{k+1}(\mathcal{T}_h))}). \end{aligned} \quad (4.64)$$

Finally, the third term in (4.62) can be bounded integrating by parts element-wise and using the Cauchy–Schwarz inequality, the trace inequality (4.14), the consistency property (4.47),

and again the Cauchy–Schwarz inequality, namely

$$\begin{aligned}
\mathcal{R}_3^n(\underline{\mathbf{v}}_h) &\leq \sum_{T \in \mathcal{T}_h} \int_T \nabla \cdot \mathbf{v}_T (\bar{p}^n - \widehat{p}_h^n) + \sum_{F \in \mathcal{F}_T} \int_F (\mathbf{v}_F - \mathbf{v}_T) \cdot (\bar{p}^n - \widehat{p}_h^n) \mathbf{n}_{TF} \\
&\leq C_{\text{pj}} h^{k+1} \left( \frac{\bar{K}}{\underline{K}} \right)^{\frac{1}{2}} |\bar{p}^n|_{H^{k+1}(\mathcal{T}_h)} \|\underline{\mathbf{v}}_h\|_{\varepsilon, h} \\
&\leq C_{\text{pj}} \frac{h^{k+1}}{\sqrt{\tau}} \|\underline{\mathbf{v}}_h\|_{\varepsilon, h} \left( \frac{\bar{K}}{\underline{K}} \right)^{\frac{1}{2}} |p|_{L^2((t^{n-1}, t^n); H^{k+1}(\mathcal{T}_h))}.
\end{aligned} \tag{4.65}$$

Therefore, combining (4.63), (4.64), and (4.65), it is inferred that

$$\|\mathcal{R}^n\|_{\varepsilon, h, *}^2 = \|\mathcal{R}_1^n + \mathcal{R}_2^n + \mathcal{R}_3^n\|_{\varepsilon, h, *}^2 \leq 4C_{\text{lp}}^2 C_{\text{eq}}^4 C_{\text{st}}^2 C_{\text{ap}}^2 \tau \|\mathbf{u}\|_{H^1((t^{n-1}, t^n); \mathbf{H}^1(\Omega))}^2 + 4C_{\text{pj}}^2 \tau^{-1} h^{2(k+1)} \widetilde{C}_1^n, \tag{4.66}$$

with

$$\widetilde{C}_1^n := \|\mathbf{u}\|_{L^2((t^{n-1}, t^n); \mathbf{H}^{k+2}(\mathcal{T}_h))}^2 + |\boldsymbol{\sigma}(\cdot, \nabla_s \mathbf{u})|_{L^2((t^{n-1}, t^n); \mathbb{H}^{k+1}(\mathcal{T}_h))}^2 + \left( \frac{\bar{K}}{\underline{K}} \right) |p|_{L^2((t^{n-1}, t^n); H^{k+1}(\mathcal{T}_h))}^2.$$

(v) *Conclusion.* Adding (4.51) to (4.58) with  $\eta = 4\beta^2 C_{\text{mn}}^{-2} C_{\text{lp}}^2 C_{\text{eq}}^6$ , using (4.61) and (4.66), summing the resulting equation over  $1 \leq n \leq N$ , and multiplying both sides by  $2\tau$ , we obtain

$$\sum_{n=1}^N \tau \left( \frac{C_{\text{mn}}^2}{2C_{\text{eq}}^2} \|\underline{\mathbf{e}}_h^n\|_{\varepsilon, h}^2 + C_0 \|\epsilon_h^n\|_{\Omega}^2 + \frac{\|\epsilon_h^n - \pi_{\Omega}^0 \epsilon_h^n\|_{\Omega}^2}{8\beta^2 C_{\text{mn}}^{-2} C_{\text{lp}}^2 C_{\text{eq}}^6} \right) + \|z_h^N\|_{c, h}^2 \leq \widetilde{C} \left( h^{2(k+1)} C_1 + \tau^2 C_2 \right), \tag{4.67}$$

with  $\widetilde{C} := \max \left\{ 1, C_{\text{pj}}^2, 2C_{\text{pj}}^2 C_{\text{mn}}^{-2} (2C_{\text{eq}}^2 + C_{\text{eq}}^{-6}), 2C_{\text{lp}}^2 C_{\text{st}}^2 C_{\text{ap}}^2 (2C_{\text{eq}}^6 + C_{\text{eq}}^{-2}), 4\beta^2 C_{\text{lp}}^2 C_{\text{eq}}^6 C_{\text{mn}}^{-2} \right\}$  and

$$\begin{aligned}
C_1 &:= \sum_{n=1}^N \left( \widetilde{C}_1^n + C_0 \frac{\bar{K}}{\underline{K}} |p|_{L^2((t^{n-1}, t^n); H^{k+1}(\mathcal{T}_h))}^2 \right), \\
&= \|\mathbf{u}\|_{L^2(\mathbf{H}^{k+2}(\mathcal{T}_h))}^2 + |\boldsymbol{\sigma}(\cdot, \nabla_s \mathbf{u})|_{L^2(\mathbb{H}^{k+1}(\mathcal{T}_h))}^2 + (1 + C_0) \frac{\bar{K}}{\underline{K}} |p|_{L^2(H^{k+1}(\mathcal{T}_h))}^2 \\
C_2 &:= \sum_{n=1}^N \left( \|\mathbf{u}\|_{H^1((t^{n-1}, t^n); \mathbf{H}^1(\Omega))}^2 + \|\phi\|_{H^1((t^{n-1}, t^n); L^2(\Omega))}^2 + C_0 \|\pi_{\Omega}^0 p\|_{H^1((t^{n-1}, t^n))}^2 \right) \\
&= \|\mathbf{u}\|_{H^1(\mathcal{T}_\tau; \mathbf{H}^1(\Omega))}^2 + \|\phi\|_{H^1(\mathcal{T}_\tau; L^2(\Omega))}^2 + C_0 \|\pi_{\Omega}^0 p\|_{H^1(\mathcal{T}_\tau)}^2.
\end{aligned}$$

Finally, multiplying both sides of (4.67) by  $2(C_{\text{pj}}^{-2} + 1)\widetilde{C}$  yields (4.48) with  $C = 2(C_{\text{pj}}^{-2} + 1)\widetilde{C}^2$ .  $\square$





# Chapter 5

---

## Poroelasticity with uncertain coefficients

---

This chapter contains some preliminary results on the numerical solution of the Biot problem with random poroelastic coefficients in the context of uncertainty quantification. It collects part of the ongoing work carried out during the nine month internship at the BRGM that took place from January to September 2018.

### Contents

---

<b>5.1</b>	<b>Introduction</b>	<b>112</b>
<b>5.2</b>	<b>The Biot model</b>	<b>113</b>
5.2.1	Fluid mass balance	113
5.2.2	Momentum balance	115
5.2.3	Constrained specific storage coefficient	115
<b>5.3</b>	<b>Probabilistic framework</b>	<b>116</b>
5.3.1	Notations and basic results	117
5.3.2	Uncertain poroelastic coefficients	118
<b>5.4</b>	<b>The Biot problem with random coefficients</b>	<b>120</b>
5.4.1	Strong and weak formulation	120
5.4.2	Well-posedness	122
<b>5.5</b>	<b>Discrete setting</b>	<b>126</b>
5.5.1	Polynomial chaos expansion	126
5.5.2	Sparse Pseudo-Spectral Projection	127
<b>5.6</b>	<b>Test cases</b>	<b>131</b>
5.6.1	Point injection	132
5.6.2	Poroelastic footing	133
5.6.3	Convergence	134
5.6.4	Sensitivity analysis	135

---

## 5.1 Introduction

The aim of this chapter is to study poroelasticity problems where the model coefficients are uncertain. The interest of this type of hydro-mechanical coupled models is particularly manifest in geosciences applications [123, 126, 150, 151], where subsurface fluid flows induce a deformation of the rock matrix. We rely here on the linear Biot model [27, 190] that describes the Darcy fluid flow in saturated porous media under the assumptions of small deformations and small variations of the porosity and the fluid density. This model depends on physical parameters that are poorly known justifying a stochastic description, due to space heterogeneities, measurement inaccuracies, and sometimes the ill-posedness of inverse problems inherent in parameter estimation techniques. Although there is an extensive literature on the poroelasticity model (we mention, in particular, the comprehensive textbooks [61, 68, 198]) and its numerical approximation, to our knowledge few works have addressed the role of uncertainty in Biot's formulation of hydromechanical coupling. In [52, 122] the authors include uncertainty in one-dimensional consolidation analysis by incorporating heterogeneity in the consolidation coefficient. The stochastic consolidation model has been extended in [159] to nonlinear uncertain soil parameters. Variability in the initial pore pressure and in the heterogeneous hydraulic mobility has been considered in [62] and [100], respectively. We also mention [65], where a stochastic Galerkin approach is proposed to solve the poroelasticity equations with randomness in all material parameters and tested on a one-dimensional problem.

Uncertainty Quantification (UQ) methods have been developed in the last decades to take into account the effect of random input quantities of a model on the quantity of interest. They enable to obtain information on the model output that is richer than in a deterministic context, since they provide the statistical moments (mean and variance), the probability distribution, and the sensitivity analysis. Among the techniques designed for UQ in numerical models, the stochastic spectral methods have received a considerable attention. The principle of these methods is to decompose a random variable on suitable approximation bases. In particular, Polynomial Chaos (PC) expansion represents the solution as a finite sum of orthogonal polynomials. PC were initially introduced by Wiener [201] and applied by Ghanem and Spanos [112] in solid mechanics and by Le Maître and Knio [137] in fluid mechanics. They were subsequently used to treat a large variety of problems, including elliptic models (*e.g.* [10, 148]), flow and transport in porous media (*e.g.* [110, 111]), thermal problems (*e.g.* [119, 141]), and hyperbolic systems (see, *e.g.* [192]). In this chapter, we focus on the non-intrusive spectral projection method, that only requires to use a deterministic solver as a black box to construct the spectral expansion of the solution. Specifically, (i) a sparse grid allows to efficiently sample the parametric space, (ii) the numerical code is run for each sample (offline stage), and (iii) the model outputs at each sample are collected and assembled in order to construct the spectral expansion of the solution.

The numerical solution of the poroelasticity coupled system requires a discretization method able to (i) treat a complex geometry with polyhedral meshes and nonconforming interfaces, (ii) handle possible heterogeneities of the poromechanical parameters, and (iii) prevent localized pressure oscillations arising in the case of low-permeable and low-

compressible porous media. We choose to discretize the poroelasticity problem using the Hybrid High-Order–discontinuous Galerkin coupled method developed in Chapter 2, which satisfies these requirements. The contribution of this work is threefold. First, a probabilistic framework is introduced to study the Biot model with uncertain coefficients. Special attention is paid to providing a physically admissible set of poroelastic parameters. Second, the well-posedness of the stochastic Biot model is proven at the continuous level. Third, a non-intrusive polynomial chaos approach is implemented in order to investigate the effect of the random poroelastic coefficients on the displacement and the pressure fields for two test cases.

The material is organized as follows. In Section 5.2 we present the poroelasticity model and identify the relations among the poromechanical coefficients allowing to express the constrained specific storage coefficient as a function of the other parameters and deterministic primary variables. In Section 5.3 we introduce the probabilistic framework and investigate how the uncertainty of the model parameters propagates on the storage coefficient. In Section 5.4 we present the linear poroelasticity problem with random coefficients in strong and weak form. We then give the assumptions on the random model parameter yielding the well-posedness of the variational problem. In Section 5.5 we outline the features of the polynomial chaos expansion and of the pseudo-spectral projection method. This procedure allows to reduce the stochastic problem to a finite set of parametric deterministic problems discretized by the HHO–dG scheme. Finally, we present numerical results in Section 5.6 to illustrate the performance of the method and to perform the sensitivity analysis.

## 5.2 The Biot model

In this section we introduce the model problem and present some preliminary relations between the coefficients characterizing the porous medium. The linear poroelasticity model, usually referred as Biot model, consists of two coupled governing equations; one describing the mass balance of the fluid and the other expressing the momentum equilibrium of the porous medium.

### 5.2.1 Fluid mass balance

The fluid mass conservation law in a *fully saturated* porous medium reads

$$d_t(\rho_f \varphi) + \nabla \cdot (\rho_f \varphi \mathbf{v}) = \rho_f g, \quad (5.1)$$

where  $d_t$  denotes the time derivative,  $\rho_f$  [kg/m<sup>3</sup>] is the fluid density,  $\varphi$  [–] the soil porosity,  $\mathbf{v}$  [ms<sup>–1</sup>] the velocity field, and  $g$  [s<sup>–1</sup>] the fluid source. The state equation for *slightly compressible* fluids (cf. [61, Section 3.1] and [206]) yields

$$d_t \rho_f = \frac{\rho_f}{K_f} d_t p, \quad (5.2)$$

with  $p$  [Pa] the pore pressure and  $K_f > 0$  [Pa] the fluid bulk modulus. Moreover, under the assumptions of an *isotropic and isothermal* conditions, *infinitesimal strains*, and *small*

relative variations of porosity, it is inferred, following [61, Section 4.1], that the change of the porosity is caused by a fluid and a mechanical effect as follows:

$$d_t \varphi = \frac{1}{M} d_t p + \alpha d_t (\nabla \cdot \mathbf{u}), \quad (5.3)$$

where  $\mathbf{u}$  [m] denotes the displacement field. We rely on the definitions of the Biot-Willis coefficient  $\alpha \in (0, 1)$  [–] and the Biot tangent modulus  $M > 0$  [Pa<sup>–1</sup>] given in [61, 68]:

$$\alpha := 1 - \frac{K_d}{K_m}, \quad M := \frac{K_m}{\alpha - \varphi}, \quad (5.4)$$

with  $K_d, K_m > 0$  [Pa] denoting the bulk moduli of the drained medium and the solid matrix, respectively. The coefficient  $\alpha$  quantifies the amount of fluid that can be forced into the medium by a variation of the pore volume for a constant fluid pressure, while  $M$  measures the amount of fluid that can be forced into the medium by pressure increments due to the compressibility of the structure. The case of a solid matrix with incompressible grains ( $K_m \rightarrow +\infty$ ) corresponds to the limit value  $1/M = 0$ . From (5.2) and (5.3), it follows that the variation of fluid content in the medium is given by

$$d_t(\rho_f \varphi) = \varphi d_t \rho_f + \rho_f d_t \varphi = \rho_f (c_0 d_t p + \alpha d_t (\nabla \cdot \mathbf{u})).$$

where, according to (5.2), (5.3), and (5.4) the constrained specific storage coefficient is defined by

$$c_0 := \frac{\alpha - \varphi}{K_m} + \frac{\varphi}{K_f}. \quad (5.5)$$

Therefore, plugging the previous relation into (5.1) and assuming that the fluid density is *uniform* in the medium, we obtain

$$c_0 d_t p + \alpha d_t (\nabla \cdot \mathbf{u}) + \nabla \cdot (\varphi \mathbf{v}) = g, \quad (5.6)$$

The fluid velocity is related to the pore pressure through the well-known Darcy's law (see for instance [68, Section 4.1.2]). Consistent with the small perturbations hypothesis adopted here, Darcy's law can be considered in its simplest form

$$\varphi \mathbf{v} = -\frac{\mathbb{K}}{\mu_f} \nabla p,$$

where  $\mathbb{K}$  [m<sup>2</sup>] is the tensor-valued intrinsic permeability and  $\mu_f$  [Pa · s] is the fluid viscosity. Thus, defining the hydraulic mobility as  $\kappa := \mathbb{K}/\mu_f$ , the mass conservation equation (5.6) becomes

$$c_0 d_t p + \alpha d_t (\nabla \cdot \mathbf{u}) - \nabla \cdot (\kappa \nabla p) = g. \quad (5.7)$$

For the sake of simplicity, we also assume in what follows that the mobility field is isotropic and strictly positive, namely  $\kappa = \kappa \mathbf{I}_d$ , with  $\kappa > 0$ .

### 5.2.2 Momentum balance

The momentum conservation equation under the *quasi-static* assumption, namely when the inertia effects in the elastic structure are negligible, reads

$$-\nabla \cdot \tilde{\boldsymbol{\sigma}} = \mathbf{f}, \quad (5.8)$$

where  $\tilde{\boldsymbol{\sigma}}$  [Pa] denotes the total stress tensor and  $\mathbf{f}$  [Pa/m] is the loading force (*e.g.* gravity). Owing to the Terzaghi's decomposition [190], the stress tensor is modeled as the sum of an effective term (mechanical effect) and a pressure term (fluid effect), yielding

$$\tilde{\boldsymbol{\sigma}} = \boldsymbol{\sigma}(\nabla_s \mathbf{u}) - \alpha p \mathbf{I}_d, \quad (5.9)$$

where  $\mathbf{I}_d$  is the identity matrix. The symmetric part of the gradient of the displacement field  $\nabla_s \mathbf{u}$  measures the strain accordingly to the small deformation hypothesis. In the context of *linear isotropic* poroelasticity, the soil mechanical behavior is described through the Cauchy strain-stress relation defined, for all  $\boldsymbol{\epsilon} \in \mathbb{R}_{\text{sym}}^{d \times d}$ , by

$$\boldsymbol{\sigma}(\boldsymbol{\epsilon}) = 2\mu\boldsymbol{\epsilon} + \lambda \text{tr}(\boldsymbol{\epsilon})\mathbf{I}_d = 2\mu \mathbf{dev}(\boldsymbol{\epsilon}) + K \text{tr}(\boldsymbol{\epsilon})\mathbf{I}_d, \quad (5.10)$$

where  $\mu > 0$  [Pa] and  $\lambda > 0$  [Pa] are the Lamé's coefficients,  $K = 2\mu/d + \lambda$  [Pa] is the bulk modulus, and

$$\text{tr}(\boldsymbol{\epsilon}) := \sum_{i=1}^d \epsilon_{ii}, \quad \mathbf{dev}(\boldsymbol{\epsilon}) := \boldsymbol{\epsilon} - \frac{\text{tr}(\boldsymbol{\epsilon})\mathbf{I}_d}{d}$$

are the trace and deviator operator, respectively. As noted in [22, 28], physical and experimental investigations suggest that the mechanical behavior of porous solids may be nonlinear. More general stress-strain relations could be used in place of the linear law (5.10), as we have done in Chapter 3 and 4. Plugging (5.10) and (5.9) into (5.8), leads to

$$-\nabla \cdot (2\mu\nabla_s \mathbf{u} + (\lambda\nabla \cdot \mathbf{u} - \alpha p)\mathbf{I}_d) = \mathbf{f}. \quad (5.11)$$

Finally, we mention the Gassmann equation (cf. [103, 150]), that relates the bulk moduli  $K$  and  $K_d$  to the Biot–Willis and storage coefficients  $\alpha, c_0$ :

$$\alpha^2 = c_0(K - K_d). \quad (5.12)$$

### 5.2.3 Constrained specific storage coefficient

Investigating the relations between the poroelastic coefficients, we propose to express the specific storage coefficient  $c_0$  as a function of other physical parameters. Using (5.4) to express the drained bulk modulus  $K_d$  as a function of  $\alpha$  and  $K_m$ , then plugging into (5.12), leads to

$$K_m = \frac{Kc_0 - \alpha^2}{c_0(1 - \alpha)}. \quad (5.13)$$

Plugging the previous equation into (5.5) and rearranging, yields the following second-order equation in  $c_0$ :

$$Kc_0^2 - \left( \alpha + \alpha\varphi + \frac{\varphi K}{K_f} - \varphi \right) c_0 + \frac{\alpha^2 \varphi}{K_f} = 0. \quad (5.14)$$

The parameters  $\varphi$ ,  $K_f$ ,  $K$ , and  $\alpha$  are then used to evaluate  $c_0$  by solving the previous equation. Some conditions need to be prescribed in order to avoid non-physical solutions. Owing to the definition of the Biot–Willis coefficient as a volume change ratio, it holds that  $0 \leq \varphi \leq \alpha \leq 1$ . As observed in [206, Section 3], a stricter lower bound can be imposed, namely

$$\frac{3\varphi}{2 + \varphi} \leq \alpha. \quad (5.15)$$

**Lemma 5.1** (Existence). *Let  $K_f, K > 0$ , and  $\varphi, \alpha$  satisfy condition (5.15). Then, there exists  $c_0 \in \mathbb{R}^+$  solution of (5.14).*

*Proof.* We prove existence by assessing the positivity of the discriminant  $\mathcal{D}$  associated to (5.14). Computing  $\mathcal{D}$  and rearranging, leads to

$$\begin{aligned} \mathcal{D} &= \left( \alpha + \alpha\varphi + \frac{\varphi K}{K_f} - \varphi \right)^2 - \frac{4\alpha^2 \varphi K}{K_f} \\ &= \varphi^2 \left( \frac{K}{K_f} - 1 \right)^2 + 2\alpha\varphi (1 + \varphi - 2\alpha) \left( \frac{K}{K_f} - 1 \right) + \alpha^2 (1 - \varphi)^2 \\ &= [\varphi \left( \frac{K}{K_f} - 1 \right) + \alpha(1 + \varphi - 2\alpha)]^2 + \alpha^2 (1 - \varphi)^2 - \alpha^2 (1 + \varphi - 2\alpha)^2 \\ &= [\varphi \left( \frac{K}{K_f} - 1 \right) + \alpha(1 + \varphi - 2\alpha)]^2 + 4\alpha^2 (1 - \alpha)(\alpha - \varphi). \end{aligned}$$

Since  $0 \leq \varphi < \alpha \leq 1$  owing to (5.15), the second term in the previous sum is positive and, as a result,  $\mathcal{D} \geq 0$ .  $\square$

*Remark 5.2.* The previous lemma yields the existence of two real solutions  $c_0^- \leq c_0^+$ . We consider  $c_0^+$  as the unique solution to (5.14) because, for admissible values of  $(K, K_f, \varphi, \alpha)$ , we might have  $c_0^- < \varphi/K_f$  violating (5.5). For instance, this is the case if  $\alpha = 2\varphi/(1+\varphi)$  for any  $\varphi \leq 3/4$  and  $K, K_f > 0$ .

### 5.3 Probabilistic framework

In many applications, only limited information about the poroelastic coefficients in (5.7) and (5.11) is available. In the context of geomechanics, even when it is possible to carry out a large number of measurements, the actual knowledge of the soil properties typically suffers from inaccuracies. Indeed, due to the presence of different layers, the physical properties can have strong variations that are difficult to estimate. Accounting for uncertainties is therefore a fundamental issue that we propose to treat in a probabilistic framework. In this section we present the probabilistic setting and we perform some preliminary investigations on the possible choices to parametrize the model coefficients. In this chapter we adopt the notations of [10, 12, 24, 53]. Thus,  $\Omega$  denotes the probability space for the parametrization of uncertainty, while the physical domain is denoted by  $D$ .

### 5.3.1 Notations and basic results

Let  $n \in \mathbb{N}$  and  $X \subset \mathbb{R}^n$  a measurable set. Spaces of functions, vector fields, and tensor fields defined over  $X$  are respectively denoted by italic capital, boldface Roman capital, and special Roman capital letters. Thus, for example,  $L^2(X)$ ,  $\mathbf{L}^2(X)$ ,  $\mathbb{L}^2(X)$  denote the spaces of square integrable functions, vector fields, and tensor fields over  $X$  respectively. We introduce an abstract probability space  $(\Theta, \mathcal{B}, \mathcal{P})$ , where  $\Theta$  is the set of possible outcomes,  $\mathcal{B}$  a  $\sigma$ -algebra of events, and  $\mathcal{P} : \mathcal{B} \rightarrow [0, 1]$  a probability measure. For any random variable  $h : \Theta \rightarrow \mathbb{R}$  defined on the abstract probability space, the expectation of  $h$  is

$$E(h) := \int_{\Theta} h(\theta) d\mathcal{P}(\theta).$$

We assume hereafter that all random quantities are second-order ones, namely they belong to

$$L^2_{\mathcal{P}}(\Theta) := \left\{ h : \Theta \rightarrow \mathbb{R} : E(h^2) < +\infty \right\}.$$

We further introduce a parametrization of the random input data using a random vector  $\xi = (\xi_1, \dots, \xi_N) : \Theta \rightarrow \Omega := [0, 1]^N$  such that, for all  $i \in \{1, \dots, N\}$ ,  $\xi_i \in L^2_{\mathcal{P}}(\Theta)$  and has zero mean. For convenience, the random variables  $\xi_i$  are assumed independent, so that the probability law  $\mathcal{P}_{\xi}$  of  $\xi$  factorizes. We denote by  $\mathcal{B}_{\xi}$  the Borel  $\sigma$ -algebra on  $\Omega$  and by  $(\Omega, \mathcal{B}_{\xi}, \mathcal{P}_{\xi})$  the image probability space. If, furthermore, each random variable  $\xi_i$  possesses the density function  $\rho_i : [0, 1] \rightarrow [0, +\infty)$ , we may define the space of second-order random variables on  $(\Omega, \mathcal{B}_{\xi}, \mathcal{P}_{\xi})$  as the weighted Lebesgue space  $L^2_{\rho}(\Omega)$ , where the weight function

$$\rho(\xi) := \prod_{i=1}^n \rho_i(\xi_i)$$

is the joint density function of  $\xi$ . The expectation operator on the image space is denoted using brackets and is related to the expectation on  $(\Theta, \mathcal{B}, \mathcal{P})$  through the identity

$$\langle h \rangle := \int_{\Omega} h(\omega) \rho(\omega) d\omega = \int_{\Theta} h(\xi(\theta)) d\mathcal{P}(\theta) = E(h \circ \xi).$$

The variance operator of  $h \in (\Omega, \mathcal{B}_{\xi}, \mathcal{P}_{\xi})$  is then defined by

$$\text{Var}(h) := \left\langle (h - \langle h \rangle)^2 \right\rangle.$$

Let  $D \subset \mathbb{R}^d$ ,  $d \in \{2, 3\}$ , denote the spatial domain occupied by the porous medium. For a given vector space  $V(D)$  of real-valued functions on  $D$ , equipped with the norm  $\|\cdot\|_{V(D)}$ , we also define the Bochner space of second-order random fields by

$$L^2_{\rho}(\Omega, V(D)) := \left\{ v : D \times \Omega \rightarrow \mathbb{R} : \|v\|_{L^2_{\rho}(\Omega, V(D))} := \left\langle \|v\|_{V(D)}^2 \right\rangle^{\frac{1}{2}} < \infty \right\}.$$



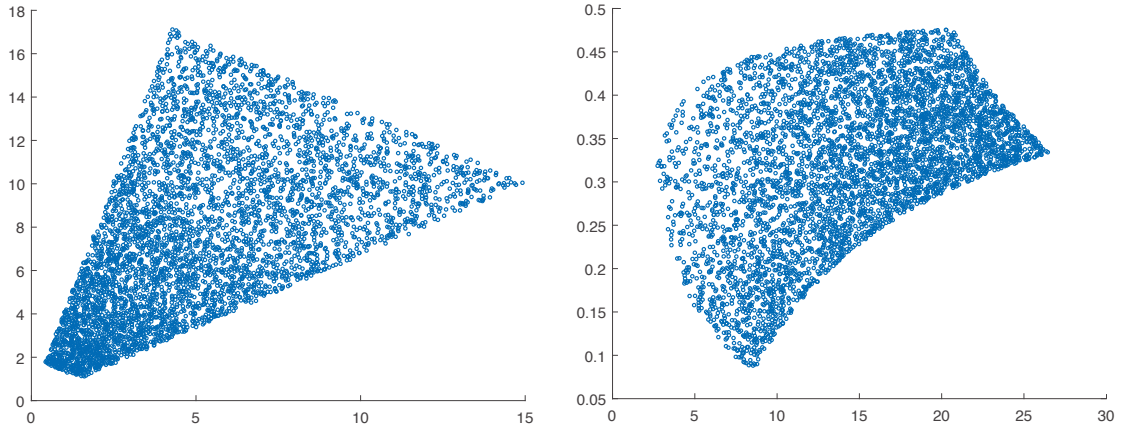


Figure 5.1: Left: Scatter plot of 5000 realizations  $(\mu, \lambda)$  for  $K \sim \mathcal{U}([1, 20])$  GPa and  $\nu \sim \mathcal{U}([0.1, 0.4])$ . Right: Scatter plot of 5000 realizations  $(K, \nu)$  for  $\mu \sim \mathcal{U}([1, 10])$  GPa and  $\lambda \sim \mathcal{U}([2, 20])$  GPa.

### 5.3.2 Uncertain poroelastic coefficients

In order to account for the uncertainty of the poroelastic material we model the Lamé's constants, the Biot–Willis coefficient, and the hydraulic mobility as spatially varying random fields depending on the finite-dimensional noise  $\xi$ . We represent the random fields  $(\mu, \lambda, \alpha, \kappa) : (\mathbf{x}, \xi) \in D \times \Omega \rightarrow \mathbb{R}^4$ , as affine combinations of the parameters  $(\xi_1, \dots, \xi_N) \in \Omega$ :

$$\begin{aligned} \mu(\mathbf{x}, \xi) &:= \mu_0(\mathbf{x}) + \sum_{i=1}^N \mu_i(\mathbf{x}) \xi_i, & \lambda(\mathbf{x}, \xi) &:= \lambda_0(\mathbf{x}) + \sum_{i=1}^N \lambda_i(\mathbf{x}) \xi_i, \\ \alpha(\mathbf{x}, \xi) &:= \alpha_0(\mathbf{x}) + \sum_{i=1}^N \alpha_i(\mathbf{x}) \xi_i, & \kappa(\mathbf{x}, \xi) &:= \kappa_0(\mathbf{x}) + \sum_{i=1}^N \kappa_i(\mathbf{x}) \xi_i. \end{aligned} \quad (5.16a)$$

One well-established approach yielding such functional dependence on the random parametrization is the truncated Karhunen–Loéve expansion [142, 180]. According to the results presented in Section 5.2.3, the constrained specific storage coefficient  $c_0$  can be expressed in terms of the elastic moduli  $\lambda, \mu$  and the coupling Biot coefficient  $\alpha$  (see also [60, 114] for theoretical and empirical investigations on the storage coefficient). Therefore, we let

$$c_0(\mathbf{x}, \xi) := c_0(\mu(\mathbf{x}, \xi), \lambda(\mathbf{x}, \xi), \alpha(\mathbf{x}, \xi)). \quad (5.16b)$$

In practice, the choice of the elastic parameters describing the mechanical properties of the medium is mostly one of convenience. For instance, one might be interested in the Young modulus  $E$  (e.g. in [128]), or the bulk modulus  $K$ , and the Poisson ratio  $\nu$  (e.g. in [60, 173]) instead of  $\mu$  and  $\lambda$ . For the sake of simplicity, here we choose the Lamé's coefficients as primary variables to avoid dealing with nonlinear parametrizations of the model coefficient inducing a warping of the  $(\mu, \lambda)$  admissible set as depicted in Figure 5.1. However, using the truncated Karhunen–Loéve or Polynomial Chaos expansions, it is possible to obtain a parametrization of the model parameter in the affine combination form (5.16a) even if the

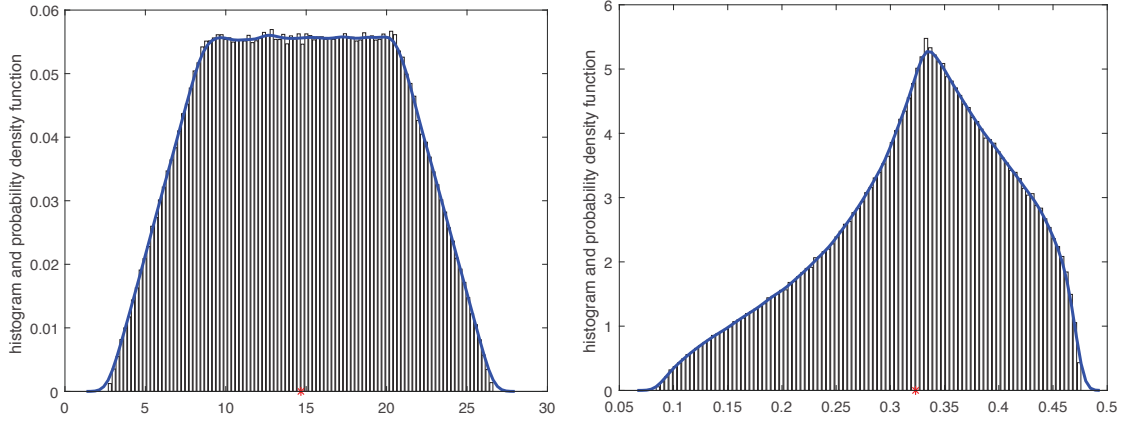


Figure 5.2: Distribution of the bulk modulus  $K$  (left) and the Poisson ratio  $\nu$  (right) corresponding to  $10^6$  realizations for  $\mu \sim \mathcal{U}([1, 10])$  GPa and  $\lambda \sim \mathcal{U}([2, 20])$  GPa.

distributions of other elastic moduli are given as entry data. In order to assess the effect of the perturbation of the Lamé's coefficients  $(\mu, \lambda)$  as in (5.16a) on the couple  $(K, \nu)$ , one can use the relations

$$K = \frac{2}{d}\mu + \lambda \quad \text{and} \quad \nu = \frac{\lambda}{2(\mu + \lambda)}.$$

We notice that taking  $(\mu, \lambda)$  as uniformly distributed primary mechanical variables, necessarily yields non-uniform distribution of  $(K, \nu)$  as illustrated in Figure 5.2.

In what follows, we investigate the effect of the parametrization (5.16b) of  $c_0$  on its probability distribution. We consider the poromechanical coefficients illustrated in Table 5.1. They are inspired from [173, Table 1], [206, Section 3], and [61, Table 4.1], and corresponds to two different layers of an ideal sandstone soil. Since in geophysical applications the fluid bulk modulus and the porosity are usually measurable with sufficient precision, we assume that  $K_f$  and  $\varphi$  are deterministic quantities, whereas  $\lambda$ ,  $\mu$ , and  $\alpha$  are random variables with uniform distribution, whose ranges are reported in Table 5.1. As a particular case of (5.16a),

	Water filled sandstone	Water filled deep sand
$K_f$	2.2 GPa	2.2 GPa
$\varphi$	10%	30%
$\mu$	$\mathcal{U}([1, 10])$ GPa	$\mathcal{U}([5, 80])$ KPa
$\lambda$	$\mathcal{U}([2, 20])$ GPa	$\mathcal{U}([12, 150])$ KPa
$\alpha$	$\mathcal{U}([3\varphi(2 + \varphi)^{-1}, 1])$	$\mathcal{U}([3\varphi(2 + \varphi)^{-1}, 1])$

Table 5.1: Poromechanical properties

we consider homogeneous random Lamé's and Biot–Willis coefficients defined such that

$$\begin{aligned} \mu(\xi) &= \mu_{\min} + \xi_1(\mu_{\max} - \mu_{\min}), \\ \lambda(\xi) &= \lambda_{\min} + \xi_2(\lambda_{\max} - \lambda_{\min}), \\ \alpha(\xi) &= \frac{3\varphi}{2 + \varphi} + \xi_3 \frac{2(1 - \varphi)}{2 + \varphi}, \end{aligned}$$

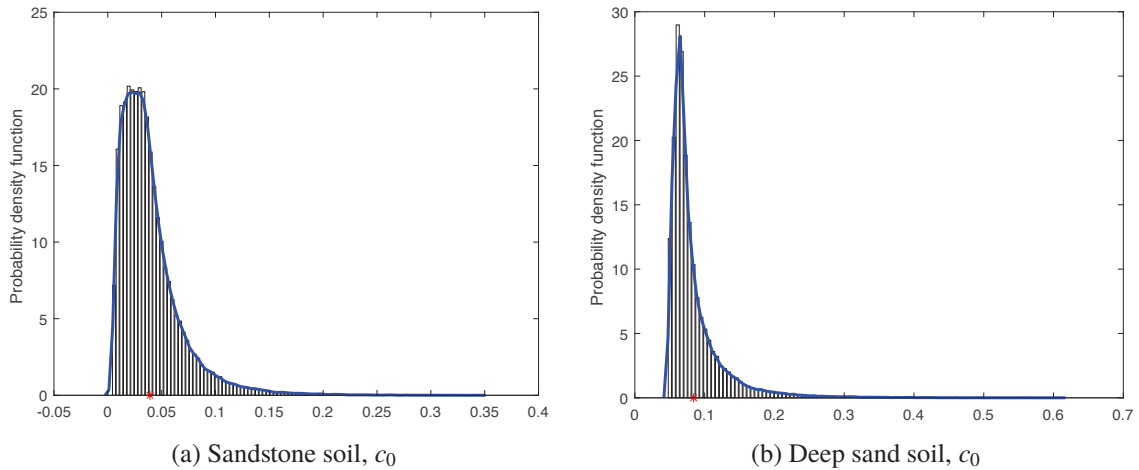


Figure 5.3: Distribution of  $10^5$  realizations of  $c_0$  obtained by solving (5.14) with uniformly distributed random input coefficients  $\lambda$ ,  $\mu$ , and  $\alpha$ .

where  $\xi \in \Omega = [0, 1]^3$  is a vector of three i.i.d. uniform variables. For each realization of  $\xi$ , we solve (5.14) in order to compute  $c_0$ . We are interested in assessing how the uncertainty on this input parameters propagates on  $c_0$ . In Figure 5.3 we plot the histograms and probability density functions of  $c_0$  corresponding to  $10^5$  realizations of  $\xi$  considering the poromechanical coefficients of Table 5.1. In particular, as expected, we observe that the computed probability density function (obtained with kernel density estimation) of  $c_0$  is far from being uniform.

## 5.4 The Biot problem with random coefficients

In this section we present the linear poroelasticity problem, discuss the assumptions on the random parameters, derive a weak formulation, and prove its well-posedness.

### 5.4.1 Strong and weak formulation

The Biot problem consists in finding a vector-valued displacement field  $\mathbf{u}$  and a scalar-valued pressure field  $p$  solutions of (5.7) and (5.11). We now consider the stochastic version obtained by taking random poroelastic coefficients as in (5.16). We assume that  $D \subset \mathbb{R}^d$  is a bounded connected polyhedral domain with boundary  $\partial D$  and outward normal  $\mathbf{n}$ . In order to close the problem, we enforce boundary conditions corresponding to a medium that is clamped on  $\Gamma_D \subset \partial D$ , traction-free on  $\Gamma_N := \partial D \setminus \Gamma_D$ , permeable with free drainage on  $\Gamma_d \subset \partial D$ , and impermeable on  $\Gamma_n := \partial D \setminus \Gamma_d$ , as well as an initial condition prescribing the initial fluid

content  $\phi_0$ . For a given finite time  $t_F > 0$ , the resulting problem is given by

$$-\nabla \cdot \boldsymbol{\sigma}(\mathbf{x}, \boldsymbol{\xi}, t) + \nabla(\alpha(\mathbf{x}, \boldsymbol{\xi})p(\mathbf{x}, \boldsymbol{\xi}, t)) = \mathbf{f}(\mathbf{x}, t), \quad (\mathbf{x}, \boldsymbol{\xi}, t) \in D \times \Omega \times (0, t_F], \quad (5.17a)$$

$$d_t \phi(\mathbf{x}, \boldsymbol{\xi}, t) - \nabla \cdot (\kappa(\mathbf{x}, \boldsymbol{\xi}) \nabla p(\mathbf{x}, \boldsymbol{\xi}, t)) = g(\mathbf{x}, t), \quad (\mathbf{x}, \boldsymbol{\xi}, t) \in D \times \Omega \times (0, t_F], \quad (5.17b)$$

$$\mathbf{u}(\mathbf{x}, \boldsymbol{\xi}, t) = \mathbf{0}, \quad (\mathbf{x}, \boldsymbol{\xi}, t) \in \Gamma_D \times \Omega \times (0, t_F], \quad (5.17c)$$

$$\boldsymbol{\sigma}(\mathbf{x}, \boldsymbol{\xi}, t) \mathbf{n} + \alpha(\mathbf{x}, \boldsymbol{\xi}) p(\mathbf{x}, \boldsymbol{\xi}, t) \mathbf{n} = \mathbf{0} \quad (\mathbf{x}, \boldsymbol{\xi}, t) \in \Gamma_N \times \Omega \times (0, t_F], \quad (5.17d)$$

$$p(\mathbf{x}, \boldsymbol{\xi}, t) = 0, \quad (\mathbf{x}, \boldsymbol{\xi}, t) \in \Gamma_d \times \Omega \times (0, t_F], \quad (5.17e)$$

$$\kappa(\mathbf{x}, \boldsymbol{\xi}) \nabla p(\mathbf{x}, \boldsymbol{\xi}, t) \cdot \mathbf{n} = 0, \quad (\mathbf{x}, \boldsymbol{\xi}, t) \in \Gamma_n \times \Omega \times (0, t_F], \quad (5.17f)$$

$$\phi(\mathbf{x}, \boldsymbol{\xi}, 0) = \phi_0(\mathbf{x}), \quad (\mathbf{x}, \boldsymbol{\xi}) \in D \times \Omega, \quad (5.17g)$$

where the stress tensor in (5.17a) and (5.17d), and the fluid content in (5.17b) and (5.17g), are defined, for all  $(\mathbf{x}, \boldsymbol{\xi}, t) \in D \times \Omega \times (0, t_F]$ , such that

$$\boldsymbol{\sigma}(\mathbf{x}, \boldsymbol{\xi}, t) = 2\mu(\mathbf{x}, \boldsymbol{\xi}) \nabla_s \mathbf{u}(\mathbf{x}, \boldsymbol{\xi}, t) + \lambda(\mathbf{x}, \boldsymbol{\xi}) (\nabla \cdot \mathbf{u}(\mathbf{x}, \boldsymbol{\xi}, t)) \mathbf{I}_d, \quad (5.17h)$$

$$\phi(\mathbf{x}, \boldsymbol{\xi}, t) = c_0(\mathbf{x}, \boldsymbol{\xi}) p(\mathbf{x}, \boldsymbol{\xi}, t) + \alpha(\mathbf{x}, \boldsymbol{\xi}) \nabla \cdot \mathbf{u}(\mathbf{x}, \boldsymbol{\xi}, t). \quad (5.17i)$$

If  $\Gamma_N = \Gamma_d = \emptyset$  and  $c_0 = 0$ , owing to (5.17b) and the homogeneous Neumann condition (5.17f), we need the following compatibility conditions on  $g$  and  $\phi_0$  and zero-average constraint on  $p$ :

$$\int_D \phi_0(\cdot) = 0, \quad \int_D g(\cdot, t) = 0, \quad \text{and} \quad \int_D p(\cdot, \cdot, t) = 0 \quad \forall t \in (0, t_F). \quad (5.17j)$$

For the sake of simplicity, we exclude the case of  $\Gamma_D$  with zero  $(d-1)$ -dimensional Hausdorff measure. Indeed, this situation simply requires an additional compatibility condition on  $\mathbf{f}$  and the prescription of the rigid-body motions of the medium. We remark that other types of boundary conditions, such as inhomogeneous ones, can also be considered up to minor modifications.

Before giving the variational formulation of problem (5.17), we introduce some notations. Let  $n \in \mathbb{N}$  and  $X \subset \mathbb{R}^n$  measurable set. For any  $m \in \mathbb{N}$ , we denote by  $H^m(X)$  the usual Sobolev space of functions that have weak partial derivatives of order up to  $m$  in  $L^2(X)$ , with the convention that  $H^0(X) := L^2(X)$ , while  $C^m(X)$  and  $C_c^\infty(X)$  denote, respectively, the usual spaces of  $m$ -times continuously differentiable functions and infinitely continuously differentiable functions with compact support on  $X$ . Since the Biot problem is unsteady, we also need to introduce, for a vector space  $V$  with scalar product  $(\cdot, \cdot)_V$ , the Bochner space  $L^2((0, t_F); V)$  spanned by square integrable  $V$ -valued functions of the time interval  $(0, t_F)$ . Similarly, the Hilbert space  $H^1(V) := H^1((0, t_F); V)$  is spanned by  $V$ -valued functions of the time interval having first-order weak derivative, and the norm  $\|\cdot\|_{H^1(V)}$  is induced by the scalar product

$$(\phi, \psi)_{H^1(V)} = \int_0^{t_F} (\phi(t), \psi(t))_V + (d_t \phi(t), d_t \psi(t))_V dt \quad \forall \phi, \psi \in H^1(V).$$

At each time  $t \in (0, t_F]$ , the natural functional spaces for the random displacement field  $\mathbf{u}(t) : D \times \Omega \rightarrow \mathbb{R}^d$  and pressure field  $p(t) : D \times \Omega \rightarrow \mathbb{R}$  taking into account the boundary

conditions (5.17c)–(5.17f) are, respectively,

$$\begin{aligned} U &:= L^2_\rho(\Omega, \mathbf{H}_{0,\Gamma_D}^1(D)), \\ P &:= \begin{cases} L^2_\rho(\Omega, H^1(D) \cap \mathbb{R}(D)^\perp) & \text{if } \Gamma_d = \Gamma_N = \emptyset \text{ and } c_0 = 0, \\ L^2_\rho(\Omega, H_{0,\Gamma_d}^1(D)) & \text{otherwise,} \end{cases} \end{aligned}$$

with  $\mathbf{H}_{0,\Gamma_D}^1(D) := \{\mathbf{v} \in \mathbf{H}^1(D) : \mathbf{v}|_{\Gamma_D} = \mathbf{0}\}$ ,  $H_{0,\Gamma_d}^1(D) := \{q \in H^1(D) : q|_{\Gamma_d} = 0\}$ , and

$$\mathbb{R}(D)^\perp := \left\{ q \in L^2(D) : \int_D q = 0 \right\}.$$

We consider the following weak formulation of problem (5.17): For a loading term  $\mathbf{f} \in L^2((0, t_F), \mathbf{L}^2(D))$ , a fluid source  $g \in L^2((0, t_F), L^2(D))$ , and an initial datum  $\phi_0 \in L^2(D)$  that verify (5.17j), find  $\mathbf{u} \in L^2((0, t_F), U)$  and  $p \in L^2((0, t_F), P)$  such that, for all  $\mathbf{v} \in U$ , all  $q \in P$ , and all  $\psi \in C_c^\infty((0, t_F))$

$$\int_0^{t_F} a(\mathbf{u}(t), \mathbf{v}) \psi(t) dt + \int_0^{t_F} b(\mathbf{v}, p(t)) \psi(t) dt = \int_0^{t_F} \langle (\mathbf{f}(t), \mathbf{v})_D \rangle \psi(t) dt, \quad (5.18a)$$

$$\int_0^{t_F} [b(\mathbf{u}(t), q) - c(p(t), q)] d_t \psi(t) dt + \int_0^{t_F} d(p(t), q) \psi(t) dt = \int_0^{t_F} \langle (g(t), q)_D \rangle \psi(t) dt, \quad (5.18b)$$

$$c(p(0), q) - b(\mathbf{u}(0), q) = \langle (\phi_0, q)_D \rangle, \quad (5.18c)$$

where  $(\cdot, \cdot)_D$  denotes the usual inner product in  $L^2(D)$  and the bilinear forms  $a : U \times U \rightarrow \mathbb{R}$ ,  $b : U \times P \rightarrow \mathbb{R}$ , and  $c : P \times P \rightarrow \mathbb{R}$  are defined such that, for all  $\mathbf{v}, \mathbf{w} \in U$  and all  $q, r \in P$ ,

$$\begin{aligned} a(\mathbf{v}, \mathbf{w}) &:= \int_\Omega \int_D (2\mu(\mathbf{x}, \boldsymbol{\xi}) \nabla_s \mathbf{v}(\mathbf{x}, \boldsymbol{\xi}) : \nabla_s \mathbf{w}(\mathbf{x}, \boldsymbol{\xi}) + \lambda(\mathbf{x}, \boldsymbol{\xi}) \nabla \cdot \mathbf{v}(\mathbf{x}, \boldsymbol{\xi}) \nabla \cdot \mathbf{w}(\mathbf{x}, \boldsymbol{\xi})) \rho(\boldsymbol{\xi}) \, d\mathbf{x} d\boldsymbol{\xi}, \\ b(\mathbf{v}, q) &:= - \int_\Omega \int_D \alpha(\mathbf{x}, \boldsymbol{\xi}) \nabla \cdot \mathbf{v}(\mathbf{x}, \boldsymbol{\xi}) q(\mathbf{x}, \boldsymbol{\xi}) \rho(\boldsymbol{\xi}) \, d\mathbf{x} d\boldsymbol{\xi}, \\ c(q, r) &:= \int_\Omega \int_D c_0(\mathbf{x}, \boldsymbol{\xi}) q(\mathbf{x}, \boldsymbol{\xi}) r(\mathbf{x}, \boldsymbol{\xi}) \rho(\boldsymbol{\xi}) \, d\mathbf{x} d\boldsymbol{\xi}, \\ d(q, r) &:= \int_\Omega \int_D \kappa(\mathbf{x}, \boldsymbol{\xi}) \nabla r(\mathbf{x}, \boldsymbol{\xi}) \cdot \nabla q(\mathbf{x}, \boldsymbol{\xi}) \rho(\boldsymbol{\xi}) \, d\mathbf{x} d\boldsymbol{\xi}. \end{aligned}$$

Above, we have introduced the Frobenius product such that, for all  $\boldsymbol{\tau}, \boldsymbol{\eta} \in \mathbb{R}^{d \times d}$ ,  $\boldsymbol{\tau} : \boldsymbol{\eta} := \sum_{1 \leq i, j \leq d} \tau_{ij} \eta_{ij}$  with corresponding norm such that, for all  $\boldsymbol{\tau} \in \mathbb{R}^{d \times d}$ ,  $|\boldsymbol{\tau}|_{d \times d} := (\boldsymbol{\tau} : \boldsymbol{\tau})^{1/2}$ .

## 5.4.2 Well-posedness

The aim of this section is to infer a stability estimate on the displacement and pressure  $(\mathbf{u}, p) \in L^2((0, t_F), U) \times L^2((0, t_F), P)$  solving (5.18), yielding, in particular, the well-posedness of the weak problem. The existence and uniqueness of a solution to the deterministic Biot

problem has been studied in [181, 205]. The results therein establish the existence. for almost every (a.e.)  $\xi \in \Omega$ , of a solution  $(\mathbf{u}(\xi), p(\xi))$  to (5.17). Thus, proving that the mapping  $\xi \rightarrow (\mathbf{u}(\xi), p(\xi))$  is measurable and bounded in  $L^2_\rho(\Omega)$  gives the existence of a solution to (5.18).

We assume some additional conditions on the input random fields defined in (5.16) that will be needed in the proofs. First, we recall that, for all  $\mathbf{x} \in D$  and all  $\xi \in \Omega$  the coupling coefficient  $\alpha(\mathbf{x}, \xi) \in (0, 1]$  satisfies the lower bound (5.15). We assume that the reference porosity of the medium is strictly positive (otherwise the Biot problem would reduce to a decoupled one), so that we have  $0 < \underline{\alpha} := \frac{3\varphi}{2+\varphi} < \alpha(\mathbf{x}, \xi)$ . We also assume that the Lamé's coefficients satisfies the following:

**Assumption 5.3** (Elastic moduli). The shear modulus  $\mu \in L^\infty(D \times \Omega)$  is uniformly bounded from above and away from zero, *i.e.* there exist positive constants  $\underline{\mu}$  and  $\bar{\mu}$  such that

$$0 < \underline{\mu} \leq \mu(\mathbf{x}, \xi) \leq \bar{\mu} < \infty, \quad \text{a.e. in } D \times \Omega. \quad (5.19)$$

The dilatation modulus  $\lambda \in L^\infty(D \times \Omega)$  is uniformly bounded, *i.e.* there exist  $\bar{\lambda} > 0$  such that

$$0 < \lambda(\mathbf{x}, \xi) \leq \bar{\lambda} < \infty, \quad \text{a.e. in } D \times \Omega. \quad (5.20)$$

The next result establishes the coercivity and the inf-sup condition respectively of the bilinear form  $a$  and  $b$  appearing in (5.18a).

**Lemma 5.4** (Coercivity and inf-sup). *Under Assumption 5.3, the following bounds hold*

$$a(\mathbf{v}, \mathbf{v}) \geq 2\underline{\mu} C_K^{-1} \|\mathbf{v}\|_U^2, \quad \forall \mathbf{v} \in U, \quad (5.21)$$

$$\sup_{\mathbf{0} \neq \mathbf{v} \in U} \frac{b(\mathbf{v}, q)}{\|\mathbf{v}\|_U} \geq \underline{\alpha} C_{\text{is}} \|q\|_{L^2_\rho(\Omega, L^2(D))}, \quad \forall q \in L^2_\rho(\Omega, L^2(D)), \quad (5.22)$$

where  $C_K > 0$  denotes the constant in Korn's first inequality, and  $C_{\text{is}} > 0$  the inf-sup constant. In the case  $\Gamma_D = \partial D$ , (5.22) holds for all  $q \in L^2_\rho(\Omega, L^2(D) \cap \mathbb{R}(D)^\perp)$ .

*Proof.* The first bound directly follows from Assumption 5.3. Indeed, since  $|\Gamma_D|_{d-1} > 0$ , we can apply Korn's first inequality, yielding

$$2\underline{\mu} \|\nabla \mathbf{v}\|_{L^2_\rho(\Omega, L^2(D))}^2 \leq 2\underline{\mu} C_K \left\langle \|\nabla_s \mathbf{v}\|_{L^2(D)}^2 \right\rangle \leq C_K a(\mathbf{v}, \mathbf{v}), \quad \forall \mathbf{v} \in U.$$

To obtain (5.22), we use [24, Lemma 7.2] establishing the existence, for any  $q \in L^2_\rho(\Omega, L^2(D))$  (or  $q \in L^2_\rho(\Omega, L^2(D) \cap \mathbb{R}(D)^\perp)$  in the case  $\Gamma_d = \partial D$ ), of  $\mathbf{v}_q \in U$  such that  $\nabla \cdot \mathbf{v}_q = q$  and  $C_{\text{is}} \|\mathbf{v}_q\|_U \leq \|q\|_{L^2_\rho(\Omega, L^2(D))}$ , with  $C_{\text{is}} > 0$  depending on  $D$ . Thus, we conclude

$$\sup_{\mathbf{0} \neq \mathbf{v} \in U} \frac{b(\mathbf{v}, q)}{\|\mathbf{v}\|_U} \geq \frac{-b(\mathbf{v}_q, q)}{\|\mathbf{v}_q\|_U} \geq \frac{\underline{\alpha} \|q\|_{L^2_\rho(\Omega, L^2(D))}^2}{\|\mathbf{v}_q\|_U} \geq \underline{\alpha} C_{\text{is}} \|q\|_{L^2_\rho(\Omega, L^2(D))}.$$

□

In order to prove the stability estimate of Proposition 5.5 below, we infer from (5.19) and (5.20) a lower bound of the specific storage coefficient  $c_0$ . We rewrite (5.13), derived from the definition of  $\alpha$  in (5.4) and the Gassmann equation (5.12), in the form

$$\frac{(1 - \alpha)(K_m - K)}{\alpha} = K - \frac{\alpha}{c_0}.$$

Since the left-hand side of the previous relation is always non-negative (according to [61, Chapter 4] we have  $K_d \leq K \leq K_m$ ), we infer that  $\alpha/c_0 \leq K$  and, as a result,

$$c_0^{-1}(\mathbf{x}, \boldsymbol{\xi}) \leq \frac{K(\mathbf{x}, \boldsymbol{\xi})}{\alpha(\mathbf{x}, \boldsymbol{\xi})} = \frac{2\mu(\mathbf{x}, \boldsymbol{\xi}) + d\lambda(\mathbf{x}, \boldsymbol{\xi})}{d\alpha(\mathbf{x}, \boldsymbol{\xi})} \leq \frac{2\bar{\mu} + d\bar{\lambda}}{d\underline{\alpha}} \quad \text{a.e. in } D \times \Omega. \quad (5.23)$$

**Proposition 5.5** (A priori estimate). *Let  $(\mathbf{u}, p) \in L^2((0, t_F), \mathbf{U} \times P)$  solve (5.18). Under Assumption 5.3, it holds*

$$\int_0^{t_F} [a(\mathbf{u}(t), \mathbf{u}(t)) + c(p(t), p(t))] dt \leq \int_0^{t_F} \left[ \frac{C_K}{2\underline{\mu}} \|\mathbf{f}\|_{L^2(D)}^2 + \frac{2\bar{\mu} + d\bar{\lambda}}{d\underline{\alpha}} \|G\|_{L^2(D)}^2 \right] dt, \quad (5.24)$$

with  $G : (0, t_F) \rightarrow L^2(D)$  defined by

$$G(t) := \int_0^t g(s) ds + \phi_0. \quad (5.25)$$

*Proof.* From (5.18a) we infer that  $a(\mathbf{u}(t), \mathbf{v}) + b(\mathbf{v}, p(t)) = (\mathbf{f}(t), \langle \mathbf{v} \rangle)_D$ , for a.e.  $t \in (0, t_F)$  and all  $\mathbf{v} \in \mathbf{U}$ . Setting  $\mathbf{v}(\mathbf{x}, \boldsymbol{\xi}) = \mathbf{u}(\mathbf{x}, \boldsymbol{\xi}, t)$  and integrating the previous relation on  $(0, t_F)$ , yields

$$\int_0^{t_F} a(\mathbf{u}(t), \mathbf{u}(t)) dt + \int_0^{t_F} b(\mathbf{u}(t), p(t)) dt = \int_0^{t_F} (\mathbf{f}(t), \langle \mathbf{u}(t) \rangle)_D dt. \quad (5.26)$$

Adapting the argument of Remark 4.2 we ensure that  $b(\mathbf{u}(t), q) - c(p(t), q) \in H^1((0, t_F))$  for all  $q \in P$ . Thus, we can integrate by parts (5.18b) and obtain

$$d_t [c(p(s), q) - b(\mathbf{u}(s), q)] + d(p(s), q) = (g(s), \langle q \rangle)_D, \quad \text{for a.e. } s \in (0, t_F).$$

Letting  $t \in (0, t_F)$ , integrating the previous identity on  $(0, t)$ , and then taking  $q(\mathbf{x}, \boldsymbol{\xi}) = p(\mathbf{x}, \boldsymbol{\xi}, t)$ , leads to

$$c(p(t), p(t)) - b(\mathbf{u}(t), p(t)) + \int_0^t d(p(s), p(t)) ds = \int_0^t (g(s), \langle p(t) \rangle)_D ds + (\phi_0, \langle p(t) \rangle)_D.$$

We define  $z(t) = \int_0^t p(s) ds$  and observe that  $d_t z(t) = p(t)$  and  $z(0) = 0$ . Using the linearity of  $d$  and the formula  $d_t z^2(t) = 2z(t) d_t z(t)$  to rewrite the third term in the left hand side of the previous relation, and recalling (5.25) we get

$$c(p(t), p(t)) - b(\mathbf{u}(t), p(t)) + \frac{1}{2} d_t [d(z(t), z(t))] = (G(t), \langle p(t) \rangle)_D ds.$$

Then, integrating on  $(0, t_F)$  and summing the resulting identity to (5.26), gives

$$\int_0^{t_F} a(\mathbf{u}(t), \mathbf{u}(t)) dt + \int_0^{t_F} c(p(t), p(t)) dt + \frac{1}{2} d(z(t_F), z(t_F)) = \int_0^{t_F} (\mathbf{f}(t), \langle \mathbf{u}(t) \rangle)_D dt + \int_0^{t_F} (G(t), \langle p(t) \rangle)_D dt. \quad (5.27)$$

To establish the result, we now bound the right-hand side of the previous relation. First, we observe that, owing to the Jensen inequality, for all  $\mathbf{v} \in \mathbf{U}$  it holds

$$\|\langle \mathbf{v} \rangle\|_{\mathbf{L}^2(D)}^2 = \int_D \left( \int_{\Omega} \mathbf{v}(\mathbf{x}, \xi) \rho(\xi) d\xi \right)^2 dx \leq \int_{\Omega} \int_D \mathbf{v}(\mathbf{x}, \xi)^2 \rho(\xi) dx d\xi \leq \|\mathbf{v}\|_{\mathbf{U}}^2.$$

Using the previous relation and the Cauchy–Schwarz and Young inequalities followed by (5.21), it is inferred that

$$\begin{aligned} \int_0^{t_F} (\mathbf{f}(t), \langle \mathbf{u}(t) \rangle)_D dt &\leq \int_0^{t_F} (2\underline{\mu} C_K^{-1})^{-1/2} \|\mathbf{f}\|_{\mathbf{L}^2(D)} (2\underline{\mu} C_K^{-1})^{1/2} \|\mathbf{u}\|_{\mathbf{U}} dt \\ &\leq \frac{C_K}{4\underline{\mu}} \int_0^{t_F} \|\mathbf{f}\|_{\mathbf{L}^2(D)}^2 dt + \frac{1}{2} \int_0^{t_F} a(\mathbf{u}(t), \mathbf{u}(t)) dt. \end{aligned}$$

Proceeding similarly for the second term in the right-hand side of (5.27) and recalling the lower bound (5.23), one has

$$\begin{aligned} \int_0^{t_F} (G(t), \langle p(t) \rangle)_D dt &\leq \int_0^{t_F} \|c_0^{-1/2} G\|_{L^2_{\rho}(\Omega, L^2(D))} \|c_0^{1/2} p\|_{L^2_{\rho}(\Omega, L^2(D))} dt \\ &\leq \frac{2\bar{\mu} + d\bar{\lambda}}{2d\underline{\alpha}} \int_0^{t_F} \|G\|_{L^2(D)}^2 dt + \frac{1}{2} \int_0^{t_F} c(p(t), p(t)) dt. \end{aligned}$$

Plugging the previous two estimates in (5.27), multiplying by a factor 2, and using the non-negativity of  $\kappa$  to infer  $d(z(t_F), z(t_F)) \geq 0$ , yields the conclusion.  $\square$

Some remarks are in order.

*Remark 5.6* (Inf-sup condition). We observe that it is possible to derive the previous a priori estimate without (5.23). Indeed the inf-sup condition (5.22) together with (5.18a) allows to bound the second term in the right-hand side of (5.27) using the following: for a.e.  $t \in (0, t_F)$ ,

$$\begin{aligned} \|p(t)\|_{L^2_{\rho}(\Omega, L^2(D))} &\leq (\underline{\alpha} C_{\text{is}})^{-1} \sup_{\mathbf{0} \neq \mathbf{v} \in \mathbf{U}} \frac{b(\mathbf{v}, p(t))}{\|\mathbf{v}\|_{\mathbf{U}}} = (\underline{\alpha} C_{\text{is}})^{-1} \sup_{\mathbf{0} \neq \mathbf{v} \in \mathbf{U}} \frac{(\mathbf{f}(t), \langle \mathbf{v} \rangle)_D - a(\mathbf{u}(t), \mathbf{v})}{\|\mathbf{v}\|_{\mathbf{U}}} \\ &\leq \frac{1}{\underline{\alpha} C_{\text{is}}} \|\mathbf{f}(t)\|_{\mathbf{L}^2(D)} + \frac{(2\bar{\mu} + d\bar{\lambda})^{1/2}}{d^{1/2} \underline{\alpha} C_{\text{is}}} a(\mathbf{u}(t), \mathbf{u}(t))^{1/2}. \end{aligned}$$

The resulting stability estimate would have, compared to (5.24), an additional dependence on  $\underline{\alpha} C_{\text{is}}$ .



*Remark 5.7* (Quasi-incompressible media). In order to prove the stability estimate (5.24) no additional assumption on the mobility  $\kappa : D \times \Omega \rightarrow \mathbb{R}^+$  is required. Thus, Proposition 5.5 can handle the case of locally poorly permeable media (*i.e.*  $\kappa \rightarrow 0$ ). However, assuming (5.20) does not allow to consider quasi-incompressible materials for which  $\lambda \rightarrow +\infty$ . To obtain a robust estimate in the case of  $\lambda$  unbounded, we can proceed as in [130, Theorem 1]. Assuming  $f \in H^1((0, t_F), \mathbf{L}^2(D))$  and  $\kappa$  uniformly bounded away from zero, the Darcy term gives a  $L^2((0, t_F), P)$  estimate of the pressure and, as a consequence, (5.20) is not needed. We remark that, in a medium featuring very low permeability, an incompressible fluid cannot flow unless the material is compressible. Therefore, the two limit cases  $\kappa \rightarrow 0$  and  $\lambda \rightarrow +\infty$  cannot occur simultaneously.

*Remark 5.8* (Lamé's coefficients). The well-posedness of problem (5.18) holds under weaker assumption on  $\mu$  and  $\lambda$ . More precisely, one can assume instead of (5.19) and (5.20), that  $\rho$ -a.e. in  $\Omega$  there holds

$$0 < \underline{\mu}(\xi) \leq \mu(\mathbf{x}, \xi) \leq \bar{\mu}(\xi) < \infty, \quad 0 \leq \lambda(\mathbf{x}, \xi) \leq \bar{\lambda}(\xi) < \infty \quad \text{a.e. in } D \times \Omega,$$

where  $\underline{\mu}(\xi)$ ,  $\bar{\mu}(\xi)$ , and  $\bar{\lambda}(\xi)$  are second-order random variables. An assumption of this type is convenient when  $\lambda$  is an unbounded (*e.g.* Gaussian or lognormal) random variable at  $\mathbf{x} \in D$ . See [11, Lemma 1.2] and [53] for a discussion in the context of elliptic PDEs with random data.

## 5.5 Discrete setting

In this section we introduce polynomial chaos expansions and outline the construction of the Pseudo-Spectral Projection (PSP) algorithm. The PSP method is one of the so-called non-intrusive techniques, which consist in performing an ensemble of deterministic model simulations to estimate the expansion coefficients. In the numerical tests of Section 5.6, each simulation is performed by solving the discrete problem of Section 2.2.5.

### 5.5.1 Polynomial chaos expansion

Let us denote  $\{\phi_{\mathbf{k}}(\xi) : \mathbf{k} \in \mathbb{N}^N\}$  an Hilbertian basis of  $L^2_\rho(\Omega)$ , where  $\phi_{\mathbf{k}}$  is a multivariate polynomial in  $\xi$  and  $\mathbf{k} = (k_1, \dots, k_N)$  is a multi-index indicating the polynomial degree in the  $\xi_i$ 's. The total degree of  $\phi_{\mathbf{k}}$  is denoted  $|\mathbf{k}| := \sum_{i=1}^N k_i$ . The basis functions are commonly chosen to be orthogonal with respect to the inner product in  $L^2_\rho(\Omega)$  characterized by the probability density function  $\rho : \Omega \rightarrow \mathbb{R}^+$ ,

$$\int_{\Omega} \phi_{\mathbf{k}}(\xi) \phi_{\mathbf{l}}(\xi) \rho(\xi) d\xi = \delta_{\mathbf{k}, \mathbf{l}} \langle \phi_{\mathbf{k}}^2 \rangle. \quad (5.28)$$

For instance, Legendre and Hermite polynomials are used for uniform and Gaussian densities, respectively (cf. [203]). If  $X(\xi) \in L^2_\rho(\Omega)$  is a second-order random variable, then it admits the so-called Polynomial Chaos (PC) expansion [112, 136],

$$X(\xi) = \sum_{\mathbf{k} \in \mathbb{N}^N} X_{\mathbf{k}} \phi_{\mathbf{k}}(\xi), \quad X_{\mathbf{k}} := \langle X, \phi_{\mathbf{k}} \rangle = \int_{\Omega} X(\xi) \phi_{\mathbf{k}}(\xi) \rho(\xi) d\xi, \quad (5.29)$$

where the deterministic coefficients of the series  $\{X_k : k \in \mathbb{N}^N\}$ , named the *spectral modes*, are defined as the projection of  $X(\xi)$  onto the basis functions. The PC approximation  $X_{\mathcal{K}}(\xi)$  of  $X(\xi)$  is obtained by truncating the expansion above to a finite series,

$$X_{\mathcal{K}}(\xi) := \sum_{k \in \mathcal{K}} X_k \phi_k(\xi), \quad (5.30)$$

where  $\mathcal{K} \subset \mathbb{N}^N$  is the set of multi-indices related to the expansion. As a result, the expectation and variance operators applied to PC approximations simplify as follows:

$$\langle X_{\mathcal{K}} \rangle = X_0, \quad (5.31)$$

$$\text{Var}(X_{\mathcal{K}}) = \sum_{k \in \mathcal{K} \setminus \mathbf{0}} X_k^2 \langle \phi_k^2 \rangle. \quad (5.32)$$

The truncation strategy, *i.e.* the selection of the set  $\mathcal{K}$ , depends on the method used to compute the spectral modes and the error  $e_{\mathcal{K}}(\xi) := X(\xi) - X_{\mathcal{K}}(\xi)$  decreases when the PC degree increases. The convergence rate of the PC expansion of  $X$  with respect to the PC degree depends on the regularity of  $X$ . In particular, an exponential convergence is expected for analytical variables; more details about the convergence conditions can be found in [46, 98]. Several methods can be implemented to compute the spectral modes of the PC expansion [125, 136, 137] such as pseudo-spectral projection, least squares regression, and compressed sensing. These procedures mainly differ in the specific discrete norm used for minimizing the truncation error  $e_{\mathcal{K}}(\xi)$ . They are non-intrusively implemented by evaluating  $X(\xi)$  at a finite number of values of  $\xi$  (the design of experiments) and therefore require only the availability of a deterministic solver. We focus here on the non-intrusive spectral projection approach [58, 59, 171] designed to avoid internal aliasing (the violation of the orthogonality condition (5.28) at the discrete level). The spectral projection method uses a quadrature rule in order to compute the spectral modes  $X_k$  by approximating (5.29),

$$X_k = \langle X, \phi_k \rangle \simeq \sum_{q=1}^{N_p} w_q X(\xi^{(q)}) \phi_k(\xi^{(q)}), \quad (5.33)$$

where  $\{\xi^{(q)}\}$  and  $\{w_q\}$  are the  $N_p$  integration points and associated weights. The previous equation shows that the main computational burden of the non-intrusive spectral projection technique is to compute  $N_p$  model evaluations for  $X(\xi^{(j)})$ . In traditional tensor-product quadrature rules,  $N_p$  scales exponentially with the number of parameters  $N$  (the so-called "Curse of Dimensionality") and is intractable for all but very low dimensional problem. This issue has motivated the use of sparse quadrature formulas, which significantly reduce the computational cost.

### 5.5.2 Sparse Pseudo-Spectral Projection

In the numerical tests of Section 5.6, we rely on the sparse Pseudo-Spectral Projection (PSP) method that is based on the partial tensorization of nested one-dimensional quadrature rules

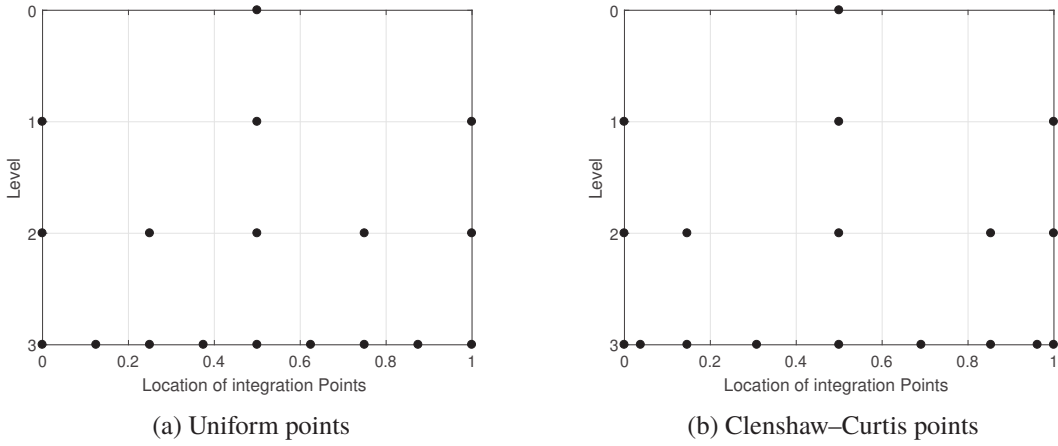


Figure 5.4: Examples of 1-D nested sequences including 1, 3, 5, and 9 quadrature points.

via Smolyak's formula [184]. Consider  $Q_0, Q_1, Q_2, \dots$  a nested sequence of 1-D quadrature formulas having increasing polynomial exactness  $p_0, p_1, p_2, \dots$ , that is

$$\int_0^1 F(\xi)\rho(\xi)d\xi = Q_l F = \sum_{q=1}^{N_p(l)} F(\xi^{(q,l)}) w_{1D}^{(q,l)}, \quad \forall F \in P^{p_l}([0, 1]),$$

where  $N_p(l)$  is the number of points in the formula  $Q_l$  and  $P^{p_l}([0, 1])$  denotes the space spanned by polynomials of total degree  $p_l$ . We call  $l$  the level of the formula  $Q_l$ . Figure 5.4 shows 1-D nested sequences over the interval  $[0, 1]$  for level  $1 \leq l \leq 3$  associated with uniform and Clenshaw–Curtis points.

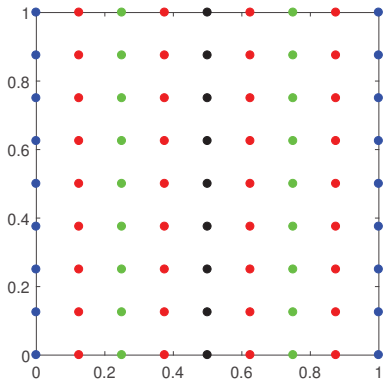
Prescribing the multi-index  $\mathbf{l} = (l_1, \dots, l_N) \in \mathbb{N}^N$ , the full tensor-product quadrature of an integrable function  $F$ , can be written as follows:

$$\mathbf{Q}_{\mathbf{l}}^{\text{FT}} F := (Q_{l_1} \otimes \dots \otimes Q_{l_N}) F = \sum_{q_1=1}^{N_p(l_1)} \dots \sum_{q_d=1}^{N_p(l_N)} F(\xi_1^{(q_1, l_1)}, \dots, \xi_N^{(q_N, l_N)}) w_{1D}^{(q_1, l_1)} \dots w_{1D}^{(q_N, l_N)}. \quad (5.34)$$

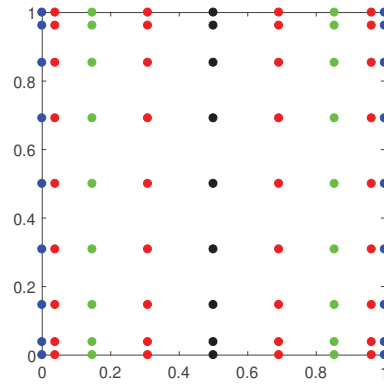
The previous quadrature rule can be used in (5.33) to compute the spectral modes  $X_k$ , using  $F = X\phi_k$ , for all  $k \in \mathcal{K}$ . The complexity of  $\mathbf{Q}_{\mathbf{l}}^{\text{FT}}$  for such nested sequences is  $N_p(\mathbf{l}) = \prod_{i=1}^N N_p(l_i)$ , and so it increases exponentially with  $N$ . Sparse grids mitigate this complexity by first introducing the 1-D difference formulas between two successive levels,  $\Delta_l^Q = Q_l - Q_{l-1}$ ,  $\Delta_0^Q = Q_0$ , such that (5.34) can be reformulated as

$$\mathbf{Q}_{\mathbf{l}}^{\text{FT}} F = \left( \left( \sum_{i_1=0}^{l_1} \Delta_{i_1}^Q \right) \otimes \dots \otimes \left( \sum_{i_N=0}^{l_N} \Delta_{i_N}^Q \right) \right) F = \sum_{i_1=0}^{l_1} \dots \sum_{i_N=0}^{l_N} (\Delta_{i_1}^Q \otimes \dots \otimes \Delta_{i_N}^Q) F = \sum_{\mathbf{i} \in \mathcal{L}_{\mathbf{l}}^{\text{FT}}} \Delta_{\mathbf{i}}^Q F,$$

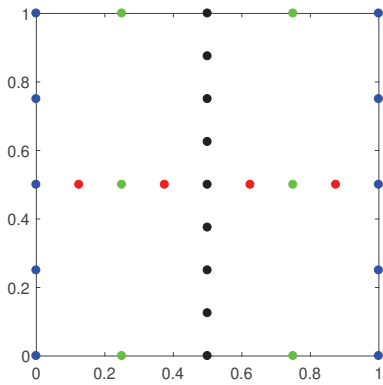
where  $\mathcal{L}_{\mathbf{l}}^{\text{FT}} = \{\mathbf{i} \in \mathbb{N}^N, i_j \leq l_j \text{ for } j = 1, \dots, d\}$  is the multi-index set of full tensorizations. The sparse quadrature rule  $\mathbf{Q}_{\mathcal{L}}$  is finally constructed by considering the summation over a



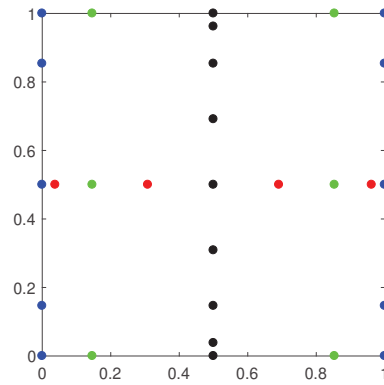
(a) Uniform points, full tensorization



(b) Clenshaw–Curtis points, full tensorization



(c) Uniform points, partial tensorization



(d) Clenshaw–Curtis points, partial tensorization

Figure 5.5: Examples of tensor-product quadrature grids for  $N = 2$  and  $\mathbf{l} = (3, 3)$ . In the full tensorization (top) the colored points correspond to the colored contributions in  $\mathcal{Q}_l^{\text{FT}} = \mathcal{Q}_0 \otimes \mathcal{Q}_3 + \Delta_1^{\mathcal{Q}} \otimes \mathcal{Q}_3 + \Delta_2^{\mathcal{Q}} \otimes \mathcal{Q}_3 + \Delta_3^{\mathcal{Q}} \otimes \mathcal{Q}_3$ . In Smolyak's tensorization (bottom) the colored points correspond to the colored contributions in  $\mathcal{Q}_{\mathcal{L}} = \mathcal{Q}_0 \otimes \mathcal{Q}_3 + \Delta_1^{\mathcal{Q}} \otimes \mathcal{Q}_2 + \Delta_2^{\mathcal{Q}} \otimes \mathcal{Q}_1 + \Delta_3^{\mathcal{Q}} \otimes \mathcal{Q}_0$ .

subset  $\mathcal{L}$  of tensorized quadrature differences:

$$\mathcal{Q}_{\mathcal{L}} F = \sum_{i \in \mathcal{L}} \Delta_i^{\mathcal{Q}} F, \quad \mathcal{L} \subset \mathcal{L}_l^{\text{FT}}.$$

The set of tensorizations  $\mathcal{L}$  must be admissible in the sense that the following condition holds [108, 109]

$$\forall \mathbf{i} = (i_1, \dots, i_N) \in \mathcal{L} : i_{1 \leq j \leq N} \geq 1 \Rightarrow \mathbf{i} - \mathbf{e}_j \in \mathcal{L},$$

where  $\{\mathbf{e}_j : j = 1, \dots, N\}$  is the set containing the canonical unit vectors of  $\mathbb{N}^N$ . This admissibility condition is necessary to preserve the telescopic property of the sum of quadrature differences. In the following, we denote by  $\mathcal{G}(\mathcal{L})$  the set of nodes in the sparse grid,

$$\mathcal{G}(\mathcal{L}) := \bigcup_{i \in \mathcal{L}} \left\{ \left( \xi^{(q_1, i_1)}, \dots, \xi^{(q_N, i_N)} \right), 1 \leq q_j \leq N_p(i_j), 1 \leq j \leq N \right\}, \quad N_p(\mathcal{L}) = |\mathcal{G}(\mathcal{L})|,$$

where  $|\cdot|$  is the cardinality of a set. In Figure 5.5 we plot the grids corresponding to full and sparse tensor-product quadrature rules in the case of  $N = 2$  and  $l = (3, 3)$ .

A quadrature rule  $\mathcal{Q}_{\mathcal{L}}$  is said to be *sufficiently exact* with respect to a polynomial multi-index set  $\mathcal{K}$  when for any couple  $(\mathbf{k}, \mathbf{k}') \in \mathcal{K} \times \mathcal{K}$ , the basis orthonormality conditions (5.28) are recovered using the discrete inner product defined, for all  $X, Y \in L^2_{\rho}(\Omega)$ , by

$$\langle X, Y \rangle_{\mathcal{Q}_{\mathcal{L}}} := \mathcal{Q}_{\mathcal{L}}(XY).$$

In the case of full tensorizations, the largest set  $\mathcal{K}^*(l)$  for which internal aliasing errors vanish can be easily determined from the degrees of polynomial exactness of the 1-D sequence, namely

$$\mathcal{K}^*(l) = \left\{ \mathbf{k} \in \mathbb{N}^N : k_i \leq \frac{p_{l_i}}{2}, i = 1, \dots, N \right\}.$$

An immediate consequence of using the sparse quadrature formula with  $\mathcal{L} \subset \mathcal{L}_l^{\text{FT}}$ , is that  $\mathcal{Q}_{\mathcal{L}}$  is not sufficiently exact with respect to  $\mathcal{K}^*(l)$ , as the sparse formula is not able to exactly integrate the product of high-order monomials. We observe that the use of sparse quadrature, while reducing the complexity ( $N_p(\mathcal{L}) \ll \prod_{j=1}^N N_p(l_j)$ ), yields also a significant reduction of the (not uniquely-defined) largest set  $\mathcal{K}^*(\mathcal{L})$  for which the quadrature is sufficiently exact.

The PSP method allows to consider polynomial spaces larger than  $\mathcal{K}^*(\mathcal{L})$ , for the same sparse grid  $\mathcal{G}(\mathcal{L})$  and without introducing internal aliasing. Let  $\{\Pi_l\}_{l \geq 1}$  denote the sequence of 1-D projection operators associated to the sequence  $\{Q_l\}_{l \geq 1}$  of 1-D nested quadrature, where

$$\Pi_l : X(\xi) \mapsto \Pi_l X(\xi) := \sum_{k=0}^{p_{l/2}} X_k \phi_k(\xi) \in P^{\frac{p_l}{2}}([0, 1]), \quad X_k = \sum_{q=1}^{N_p(l)} X(\xi^{(q,l)}) \phi_k(\xi^{(q,l)}) w_{1D}^{(q,l)}.$$

Note that  $\Pi_{l \geq 1}$  is free of internal aliasing, owing to the dependence on  $l$  of the projection spaces, which ensures that the quadrature is sufficiently exact. Consider for  $l \geq 1$  the difference of successive 1-D projection operators  $\Delta_{l>1}^{\Pi} = \Pi_l - \Pi_{l-1}$ ,  $\Delta_1^{\Pi} = \Pi_1$ . It is possible to write the full tensor-product projection corresponding to the quadrature rule given by (5.34) in terms of a sum of tensorized difference operators by observing that

$$\begin{aligned} \mathbf{\Pi}_l^{\text{FT}} X &:= \sum_{\mathbf{k} \in \mathcal{K}^*(l)} \mathcal{Q}_l^{\text{FT}}(X \phi_{\mathbf{k}})(\xi) = (\Pi_{l_1} \otimes \dots \otimes \Pi_{l_N}) X \\ &= \sum_{i_1=1}^{l_1} \dots \sum_{i_N=1}^{l_N} (\Delta_{i_1}^{\Pi} \otimes \dots \otimes \Delta_{i_N}^{\Pi}) X = \sum_{i \in \mathcal{L}_l^{\text{FT}}} (\Delta_{i_1}^{\Pi} \otimes \dots \otimes \Delta_{i_N}^{\Pi}) X. \end{aligned}$$

The sparse PSP operator  $\mathbf{\Pi}_{\mathcal{L}}$  is finally obtained by considering a summation over an admissible subset  $\mathcal{L} \subset \mathcal{L}_l^{\text{FT}}$  of tensorized difference projection operators. This results in

$$\mathbf{\Pi}_{\mathcal{L}} X = \sum_{i \in \mathcal{L}} (\Delta_{i_1}^{\Pi} \otimes \dots \otimes \Delta_{i_N}^{\Pi}) X. \quad (5.35)$$

The key point is that the sparse PSP operator in (5.35) involves a telescopic sum of differences

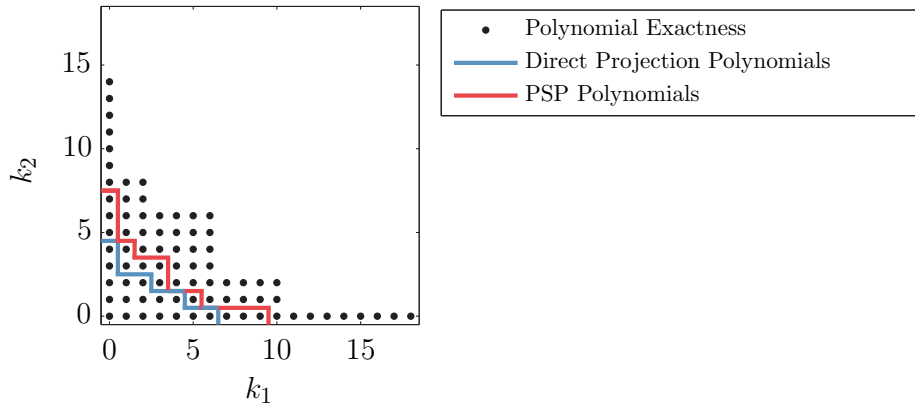


Figure 5.6: Comparison of polynomials multi-index sets  $\mathcal{K}(\mathcal{L})$  (below the red line) and  $\mathcal{K}^*(\mathcal{L})$  (below the blue line) using a quadrature rule  $\mathcal{Q}_{\mathcal{L}}$  in  $d = 2$  dimensions. Points are multi-indices  $(k_1, k_2)$  such  $F \in P^{k_1}([0, 1]) \times P^{k_2}([0, 1])$  is exactly integrated by  $\mathcal{Q}_{\mathcal{L}}$ .

in the projection onto subspaces of increasing dimensions, where each individual tensorized quadrature projection is evaluated without internal aliasing. Moreover,  $\mathbf{\Pi}_{\mathcal{L}}$  belongs to the span of  $\{\phi_{\mathbf{k}} : \mathbf{k} \in \mathcal{K}(\mathcal{L})\}$ , with  $\mathcal{K}(\mathcal{L}) = \bigcup_{i \in \mathcal{L}} \mathcal{K}^*(i)$ . We have  $\mathcal{K}(\mathcal{L}) \supseteq \mathcal{K}^*(\mathcal{L})$  and the projection space of the sparse PSP is usually significantly larger than that of the direct sparse projection. Nevertheless, the two methods have the same complexity relying both on the same sparse grid  $\mathcal{G}(\mathcal{L})$ . Note, however, that the global sparse quadrature rule  $\mathcal{Q}_{\mathcal{L}}$  is generally not sufficiently accurate with respect to  $\mathcal{K}(\mathcal{L})$ . The inclusion of the sets  $\mathcal{K}(\mathcal{L})$  and  $\mathcal{K}^*(\mathcal{L})$  is illustrated in Figure 5.6.

## 5.6 Test cases

We assess the performance of the PSP method detailed in the previous section on two model problems. The first one is an injection test designed to investigate the effect of a fluid source on the deformation of the domain, while the second one is meant to evaluate the mechanical effect of a traction term on the flux. Our numerical experiments feature a simplified model with homogeneous random coefficients, namely the dependence of the spacial variable  $\mathbf{x}$  in (5.16) is dropped. More realistic models, requiring (possibly strongly correlated) spatially varying input coefficients with non-standard distributions will be the object of future works. After presenting the description of the test cases, we numerically investigate the convergence of the probability error with respect to the level of the sparse grid. Finally, we perform a sensitivity analysis in order to rank the contribution of the input parameters on the variance of the solutions.

Our input uncertainty is parametrized by a vector  $\boldsymbol{\xi}$  of dimension  $N = 4$ . In particular, we assume that  $\mu$ ,  $\lambda$ , and  $\kappa$  have a log-uniform distribution, while  $\alpha$  is uniformly distributed. As proposed in Section 5.3.2, the specific storage coefficient  $c_0$  is evaluated from the other uncertain input parameters. The advantage of this approach is twofold: (i) it ensures that

the set of poroelastic coefficients belong to the physical admissible set, and (ii) it allows to reduce the uncertainty dimension and then the size of the sampling (*i.e.*, number of sparse grid points). Even if the experiments presented here are mainly academic model problems, we try to consider realistic sets of parameters and have a particular regard for the underlying physical phenomena. The ranges of variation of the uncertain input data are inspired from an ideal water filled sand soil with a  $\varphi = 20\%$  reference porosity. In both tests we set

$$\begin{aligned}\mu(\xi) &= 10^{2\xi_1} \text{ KPa}, \\ \lambda(\xi) &= 2 \cdot 10^{2\xi_2} \text{ KPa}, \\ \alpha(\xi) &= \frac{3\varphi}{2+\varphi} + \xi_3 \frac{2(1-\varphi)}{2+\varphi}, \\ \kappa(\xi) &= 10^{2(\xi_4-1)} \text{ m}^2 \text{ KPa}^{-1} \text{ s}^{-1},\end{aligned}\tag{5.36}$$

with uniformly distributed  $(\xi_1, \dots, \xi_4) \in [0, 1]^4$ . This choice yields

$$\begin{aligned}\langle \mu \rangle &\sim 21.5 \text{ KPa}, \quad \langle \lambda \rangle \sim 43 \text{ KPa}, \quad \langle \alpha \rangle \sim 0.64, \quad \text{and} \quad \langle \kappa \rangle \sim 0.22 \frac{\text{m}^2}{\text{KPa s}}, \\ c_v(\mu) &\sim 1.16, \quad c_v(\lambda) \sim 1.16, \quad c_v(\alpha) \sim 0.33, \quad \text{and} \quad c_v(\kappa) \sim 1.16,\end{aligned}\tag{5.37}$$

where  $c_v(\cdot) := \frac{\text{Var}(\cdot)^{1/2}}{\langle \cdot \rangle}$  denotes the coefficient of variation operator. In particular we observe that the previous average values have the same magnitude as the coefficients considered in [102, Table 2]. Our choice leads to a specific storage coefficient  $c_0$  ranging in  $(10^{-3}, 10^{-1})$ .

### 5.6.1 Point injection

The first experiment is inspired by the Barry and Mercer test considered in Section 2.6.2 (see also [14, 166]). On the boundary of the unit square domain  $D = [0, 1]^2$  we impose the following time-independent homogeneous boundary conditions:

$$\mathbf{u} \cdot \boldsymbol{\tau} = 0, \quad \mathbf{n}^T \nabla \mathbf{u} \mathbf{n} = 0, \quad p = 0,$$

where  $\boldsymbol{\tau}$  denotes the tangent vector on  $\partial D$ . The loading term  $\mathbf{f}$  in (5.17a) and the initial condition  $\phi_0$  in (5.17g) are both chosen to be zero. The evolution of the displacement and pressure fields solving (5.17) is only driven by a stationary fluid source located at  $\mathbf{x}_0 = (0.25, 0.25)$ :

$$g = 10 \delta(\mathbf{x} - \mathbf{x}_0),$$

with  $\delta$  denoting the Dirac delta function. The following results have been obtained using the HHO–dG coupled method with polynomial degree  $k = 3$  on a Cartesian mesh containing 1,024 elements. Using the static condensation procedure detailed in Section 2.5, the spacial discretization consists of 27,648 unknowns. We are interested in the stationary solution, namely the pressure and displacement fields manifesting after the initial transient phase. Therefore, we only focus on the configuration at the final time  $t_F = 1\text{s}$  obtained after 10 time steps of the implicit Euler scheme.

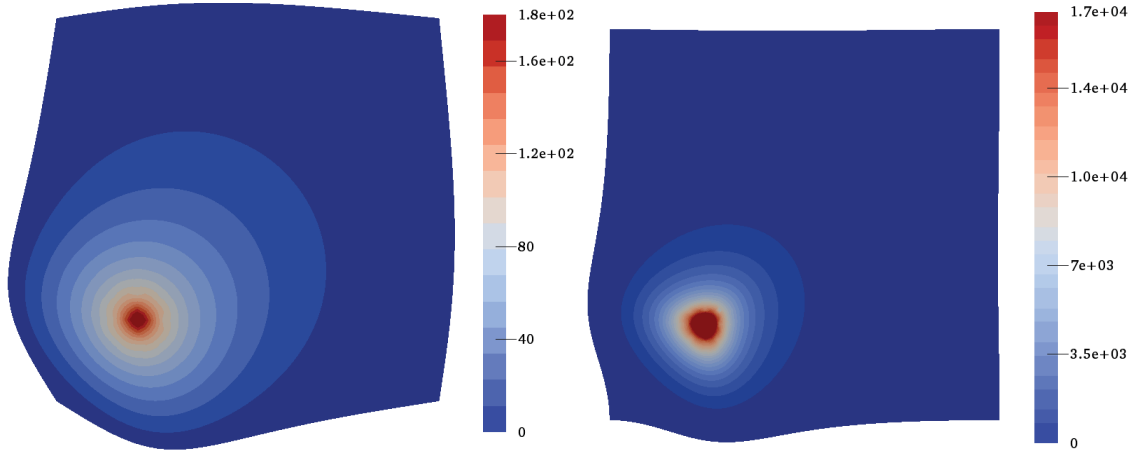


Figure 5.7: Pressure plotted on the deformed domain at  $t = t_F = 1$  (data are in KPa). Mean (Left) and Variance (Right) of the PSP method with level  $l = 5$  and  $N_p(l) = 2561$ .

Figure 5.7 represents the two first PC statistical moments given by (5.31) and (5.32) of the displacement and pressure solutions approximated using a level  $l = 5$  sparse grid. We plot the pressure field on the deformed configuration obtained by applying the displacement to the reference unit square domain. As expected, we observe a dilatation of the domain caused by the injection. The computed solutions are symmetric with respect to the diagonal according to the source location and the homogeneous boundary conditions. The pressure mean field exhibits maximum values around the source, while no perturbation is observed near the boundary. The pressure variance field has the same behavior as the mean, with higher values close to the injection point. Moreover, we notice that the magnitude of the standard deviation (*i.e.* the square root of the variance) is comparable to the mean field range.

### 5.6.2 Poroelastic footing

As second example, we consider the 2D footing problem of [102, 155]. The simulation domain is a unit square porous soil that is assumed to be free to drain and fixed at the base and at the vertical edges. A uniform load is applied on a portion of the upper boundary  $\Gamma_N = \{\mathbf{x} = (x_1, x_2) \in \partial D : x_2 = 1\}$  simulating a footing step compressing the medium. Thus, the boundary conditions data are given as follows:

$$\begin{aligned} \boldsymbol{\sigma} \mathbf{n} &= (0, -5) \text{ KPa} && \text{on } \Gamma_{N,1} := \{\mathbf{x} = (x_1, x_2) \in \partial D : 0.3 \leq x_1 \leq 0.7, x_2 = 1\}, \\ \boldsymbol{\sigma} \mathbf{n} &= \mathbf{0}, && \text{on } \Gamma_N \setminus \Gamma_{N,1}, \\ \mathbf{u} &= \mathbf{0}, && \text{on } \partial D \setminus \Gamma_N, \\ p &= 0, && \text{on } \partial D. \end{aligned}$$

Also in this case we take  $\mathbf{f} = \mathbf{0}$  and  $\phi_0 = 0$ , so that the displacement and pressure solutions are determined only by the prescription of the stress boundary condition. Thus, the situation is complementary to the one of the previous test, where the deformation of the medium was caused by a fluid term. We solve the Biot problem (5.17) with uncertain coefficients defined in (5.36). For the space discretization of the footing problem, we consider a triangular mesh



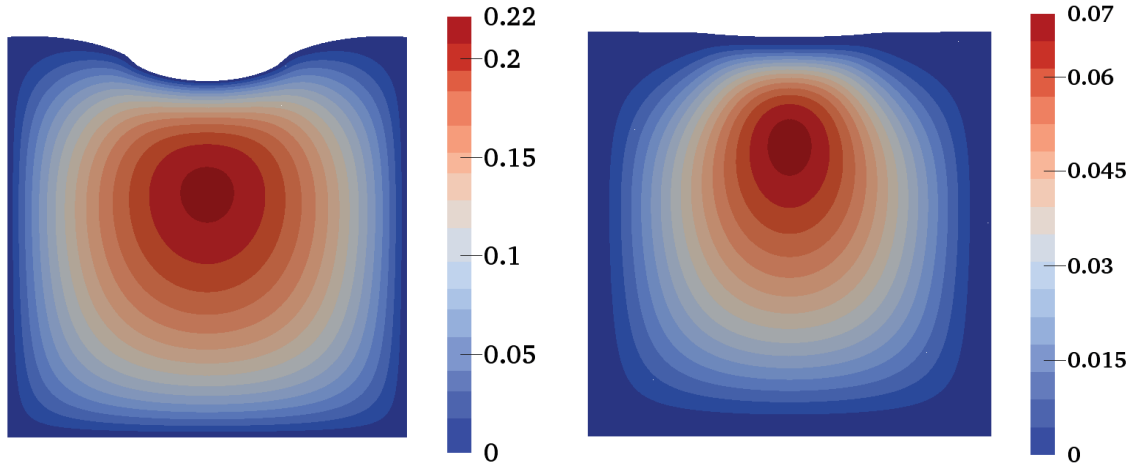


Figure 5.8: Pressure plotted on the deformed domain at  $t = 0.2$ . Mean (left) and Variance (right) computed using the PSP method with level  $l = 5$  and  $N_p(l) = 2561$ .

with 3,584 elements and we use the HHO–dG method with  $k = 2$ . For each point of the sparse grid and each time iteration, the linear system to be solved has dimension 54,720. In this test case, we are particularly interested in the pressure profile at early times. In fact, owing to the free drainage boundary condition, after a few time steps the fluid is completely squeezed out of the soil, yielding an equilibrium configuration with no pressure and the displacement balancing the Neumann condition. We consider the approximated solution at time  $t = 0.2s$  obtained after two time steps of the implicit Euler method. We mention, in passing, that the unphysical pressure oscillations observed in [155] do not occur with our choice of the discretization (this is in agreement with the results of Section 2.6.2). We plot in Figure 5.8 the mean and variance fields corresponding to the PC approximation of the solution. Concerning the mean, we observe a significant deformation where the load is applied and a maximum pressure value near the center of the domain. The variance fields exhibit a similar configuration.

### 5.6.3 Convergence

Let  $V(D)$  be a normed vector space of functions over  $D$  and  $X \in L^2_\rho(\Omega, V(D))$  a random field. We evaluate the predictive capacity of the PC expansion  $X_{\mathcal{K}}$  defined in (5.29) using the mean-squared error (MSE) computed from a validation set of  $N^* = 500$  Latin Hypercube Samples (LHS) and defined as

$$\text{MSE}(X - X_{\mathcal{K}})(\mathbf{x}) := \frac{1}{N^*} \sum_{i=1}^{N^*} (X(\mathbf{x}, \xi^i) - X_{\mathcal{K}}(\mathbf{x}, \xi^i))^2.$$

We plot in figures 5.10 and 5.11 the MSE fields of the PC solutions of the injection and footing test, respectively. Owing to the symmetry of the point injection test with respect to the diagonal, in 5.10 we only plot the first component of  $\text{MSE}(\mathbf{u} - \mathbf{u}_{\mathcal{K}})(\mathbf{x})$ . Similarly, for the footing problem we consider the second component. The results are obtained using sparse

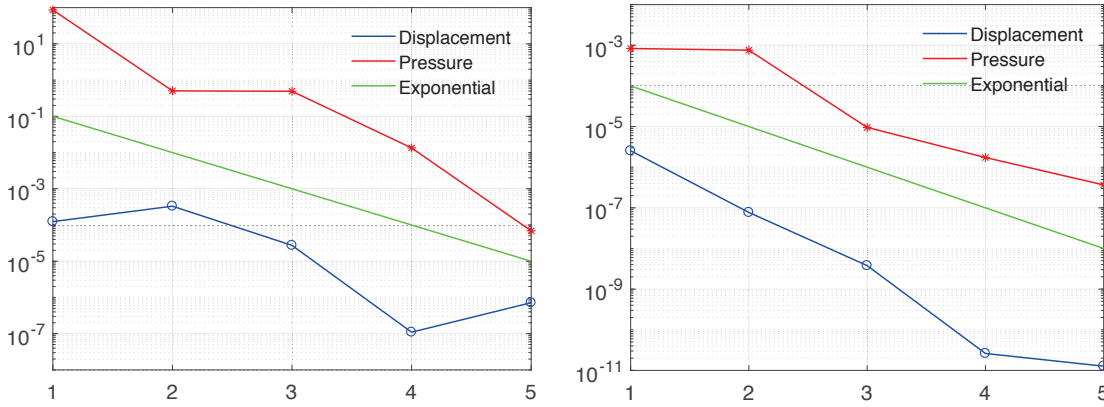


Figure 5.9: Errors  $\|\text{MSE}(\mathbf{u} - \mathbf{u}_{\mathcal{K}})\|_{L^2(D)}$  and  $\|\text{MSE}(p - p_{\mathcal{K}})\|_{L^2(D)}$  vs. level  $l$  of the Sparse Grid. Point injection test (left) and footing test (right) with model coefficients  $\mu, \lambda, \alpha, \kappa$  distributed according to (5.36).

grid of level  $l = 1, 3, 5$ . We observe a rapid decrease of the MSE field for both experiments and both the displacement and pressure. We especially note the strong decay of the maximal value with respect to the increasing sparse grid level. In Figure 5.9 we plot the  $L^2(D)$ -norm of the mean-squared error fields  $\|\text{MSE}(\mathbf{u} - \mathbf{u}_{\mathcal{K}})\|_{L^2(D)}$  and  $\|\text{MSE}(p - p_{\mathcal{K}})\|_{L^2(D)}$ . The linear trends observed in the semi-log scale of the sequential PC approximations suggest that the method achieves roughly exponential convergence with respect to the level of the sparse grid, as predicted by theory.

#### 5.6.4 Sensitivity analysis

We conclude the numerical investigation by assessing the different contributions of the input parameters on the variance of the solutions. For each of the experiments presented above, we compute the first order conditional variances

$$\text{Var}(\mathbf{u}_{\mathcal{K}}(\mathbf{x}, \boldsymbol{\xi}) | \xi_i)_{1 \leq i \leq 4}, \quad \text{and} \quad \text{Var}(p_{\mathcal{K}}(\mathbf{x}, \boldsymbol{\xi}) | \xi_i)_{1 \leq i \leq 4},$$

of the PC expansions obtained with the finer sparse grid ( $l = 5$ ). Since each  $\xi_i$  in (5.36) parametrizes only one of the poromechanical model coefficients, the previous quantities allows to evaluate and differentiate the effect of the perturbations of  $\mu, \lambda, \alpha, \kappa$  on the solutions.

Concerning the injection test, we are mainly interested in assessing the effect of the fluid source on the deformation. Therefore, we perform the sensitivity analysis considering the displacement as the quantity of interest. On the other hand, in the footing experiment we focus on the pressure. In Figure 5.12 we compare the total variance  $\text{Var}(u_{\mathcal{K},1})$  defined in (5.32) with the sum of the first-order conditional variances  $\sum_{i=1}^4 \text{Var}(u_1(\mathbf{x}, \boldsymbol{\xi}) | \xi_i)$  of the PC approximated displacement solving the injection test. In Figure 5.14 we perform the same comparison for the footing problem considering the pressure variance  $\text{Var}(p_{\mathcal{K},1})$  and its first-order conditional contributions. The difference between the two fields indicates that the contribution due to second-order conditional variances is not negligible. This is probably due to the fact that computing  $c_0$  using (5.14) yields a dependence  $c_0 = c_0(\xi_1, \xi_2, \xi_3)$ . We observe in Figures 5.13 and 5.15 that, in both experiments, the parameter with the main

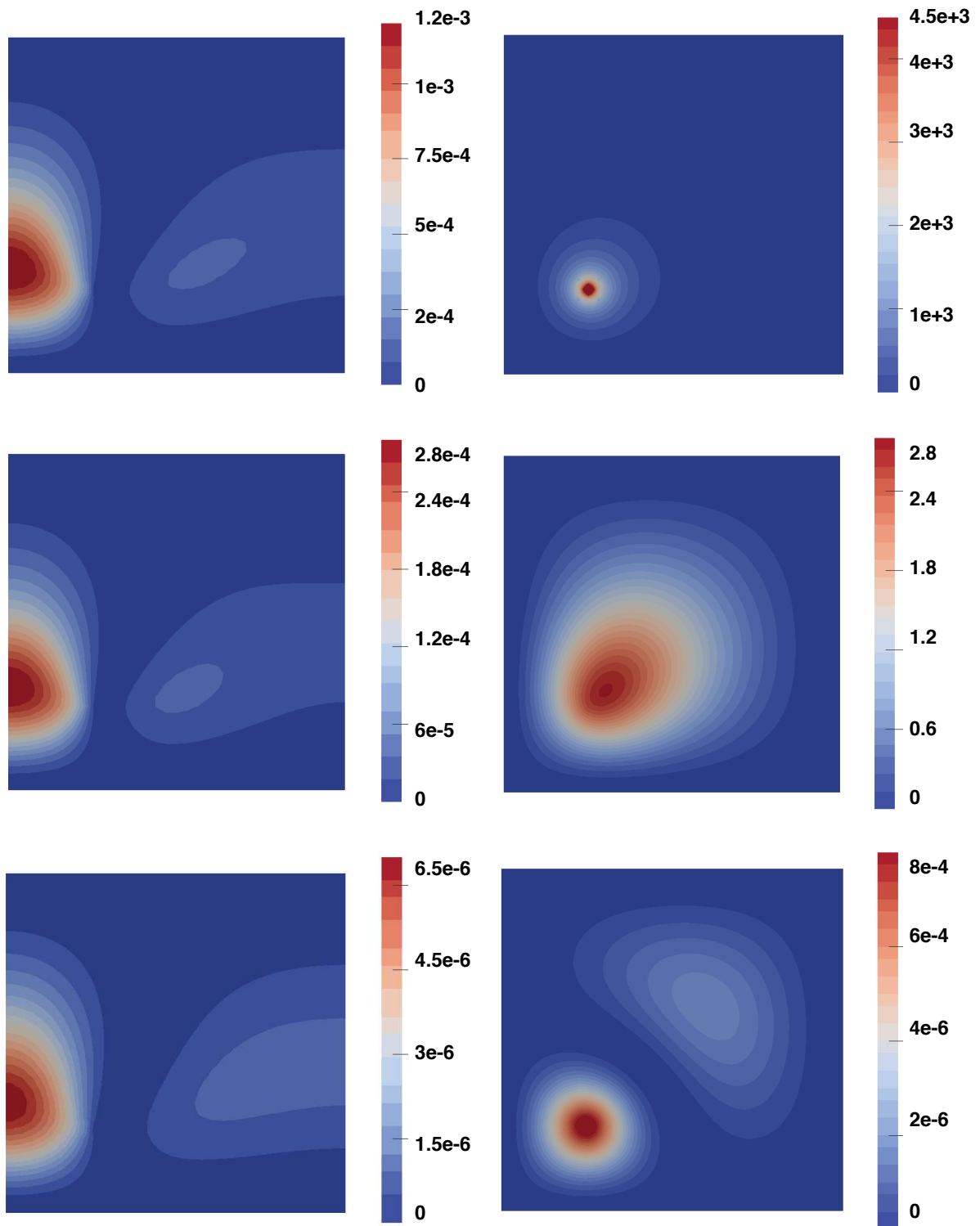


Figure 5.10: Displacement  $\text{MSE}(u_1 - u_{\mathcal{K},1})$  field (left) and pressure  $\text{MSE}(p - p_{\mathcal{K}})$  field (right) of the injection test case. Results are computed using the PSP method with  $l = 1$  (top),  $l = 3$  (middle), and  $l = 5$  (bottom).

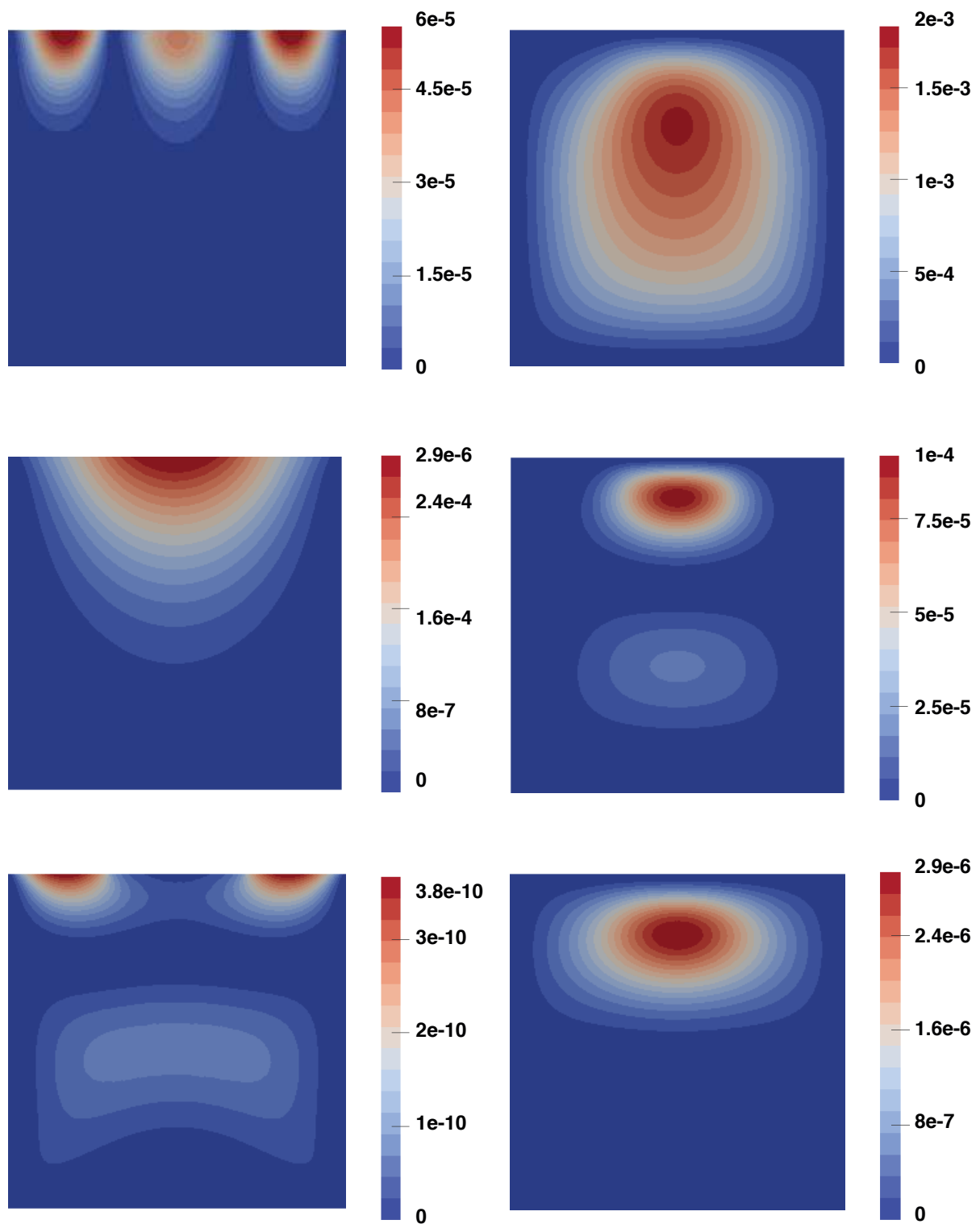


Figure 5.11: Displacement  $\text{MSE}(u_2 - u_{\mathcal{K},2})$  field (left) and pressure  $\text{MSE}(p - p_{\mathcal{K}})$  field (right) corresponding to the footing test case. Results computed using the PSP method with  $l = 1$  (top),  $l = 3$  (middle), and  $l = 5$  (bottom).

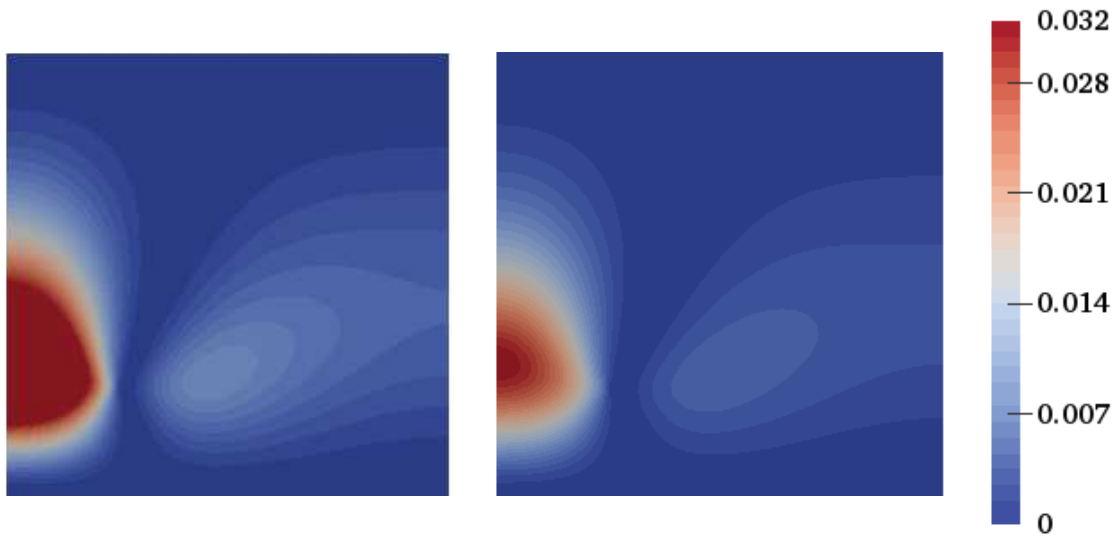


Figure 5.12: Point injection test: Comparison between the total variance  $\text{Var}(u_1(\mathbf{x}, \boldsymbol{\xi}))$  and the sum of the first-order conditional variances  $\sum_{i=1}^4 \text{Var}(u_1(\mathbf{x}, \boldsymbol{\xi})|\xi_i)$  estimated using the PSP method with  $l = 5$ .

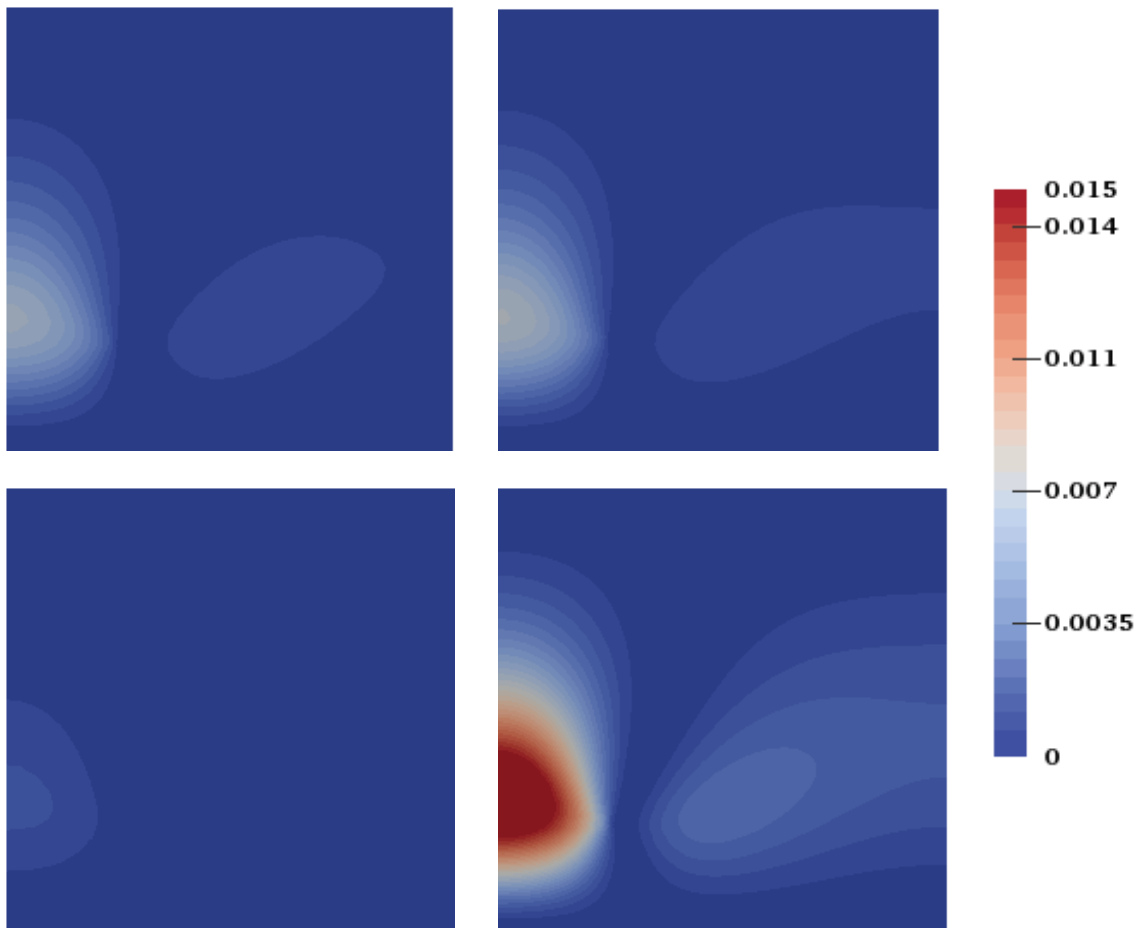


Figure 5.13: Point injection test: Conditional variances  $\text{Var}(u_1(\mathbf{x}, \boldsymbol{\xi})|\xi_1)$  (top left)  $\text{Var}(u_1(\mathbf{x}, \boldsymbol{\xi})|\xi_2)$  (top right)  $\text{Var}(u_1(\mathbf{x}, \boldsymbol{\xi})|\xi_3)$  (bottom left)  $\text{Var}(u_1(\mathbf{x}, \boldsymbol{\xi})|\xi_4)$  (bottom right) estimated using the PSP method with  $l = 5$ .

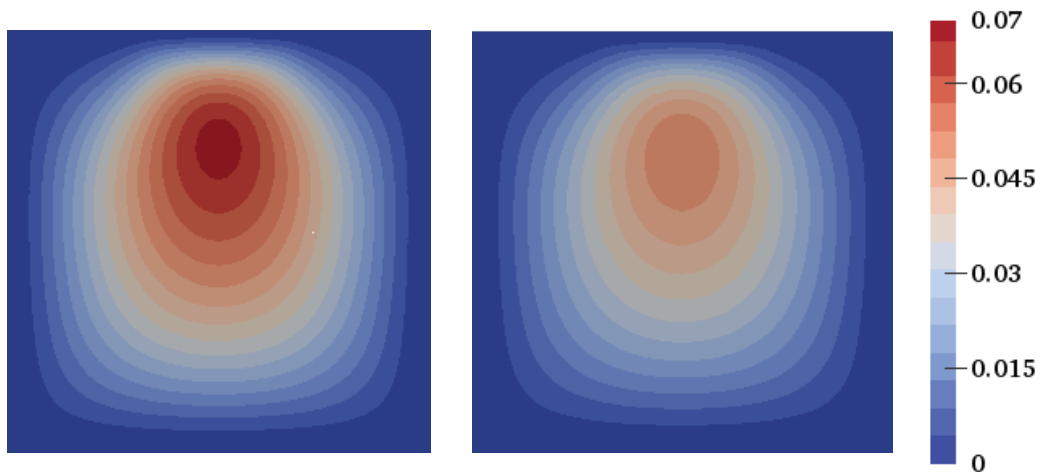


Figure 5.14: Poroelastic footing test: Comparison between the total variance  $\text{Var}(p(\mathbf{x}, \boldsymbol{\xi}))$  and the sum of the first-order conditional variances  $\sum_{i=1}^4 \text{Var}(p(\mathbf{x}, \boldsymbol{\xi})|\xi_i)$  estimated using the PSP method with  $l = 5$ .

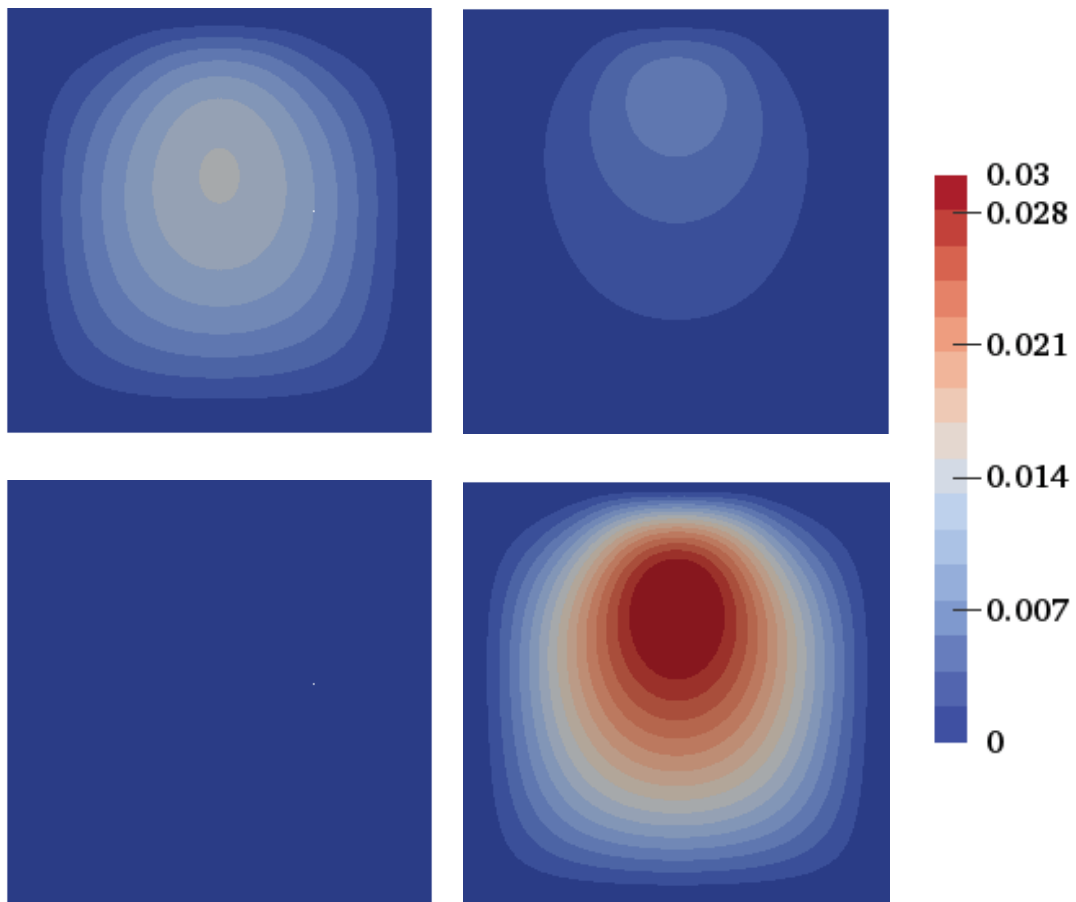


Figure 5.15: Poroelastic footing test: Conditional variances  $\text{Var}(p(\mathbf{x}, \boldsymbol{\xi})|\xi_1)$  (top left)  $\text{Var}(p(\mathbf{x}, \boldsymbol{\xi})|\xi_2)$  (top right)  $\text{Var}(p(\mathbf{x}, \boldsymbol{\xi})|\xi_3)$  (bottom left)  $\text{Var}(p(\mathbf{x}, \boldsymbol{\xi})|\xi_4)$  (bottom right) estimated using the PSP method with  $l = 5$ .

influence on the solutions is the hydraulic mobility  $\kappa$  (parametrized by  $\xi_4$ ), while the effect of the Biot–Willis coefficient  $\alpha$  (parametrized by  $\xi_3$ ) on the variance is very weak. In particular, we notice in Figure 5.15 that  $\alpha$  has no influence on the uncertainty of the pressure solving the footing problem. Thus, in this case  $\alpha$  could be fixed at a deterministic value. One possible reason may be that the coefficient of variation  $c_v(\alpha)$  computed in (5.37) is smaller than the ones relative to the other model coefficients. Finally, we remark that the conditional variances corresponding to the Lamé’s coefficients  $\mu, \lambda$  are almost the same in the case of the injection problem (cf. Figure 5.13), but are considerably different in the footing test (cf. Figure 5.15).

# Appendix A

---

## A nonconforming high-order method for nonlinear poroelasticity

---

This appendix has been published in the following conference book (see [34])

**Finite Volumes for Complex Applications VIII,**  
Hyperbolic, Elliptic and Parabolic Problems, 2017, Pages 537–545.

### Contents

---

<b>A.1</b>	<b>Introduction</b>	<b>141</b>
<b>A.2</b>	<b>Mesh and notation</b>	<b>142</b>
<b>A.3</b>	<b>Discretization</b>	<b>143</b>
A.3.1	Nonlinear elasticity operator	143
A.3.2	Darcy operator	144
A.3.3	Hydro-mechanical coupling	144
<b>A.4</b>	<b>Formulation of the method</b>	<b>145</b>
<b>A.5</b>	<b>Numerical results</b>	<b>146</b>

---

## A.1 Introduction

We consider the nonlinear poroelasticity model obtained by generalizing the linear Biot’s consolidation model of [27, 190] to nonlinear stress-strain constitutive laws. Our original motivation comes from applications in geosciences, where the support of polyhedral meshes and nonconforming interfaces is crucial.

Let  $\Omega \subset \mathbb{R}^d$ ,  $d \in \{2, 3\}$ , denote a bounded connected polyhedral domain with Lipschitz boundary  $\partial\Omega$  and outward normal  $\mathbf{n}$ . For a given finite time  $t_F > 0$ , volumetric load  $\mathbf{f}$ , fluid source  $g$ , the considered nonlinear poroelasticity problem consists in finding a vector-valued displacement field  $\mathbf{u}$  and a scalar-valued pore pressure field  $p$  solution of

$$-\nabla \cdot \boldsymbol{\sigma}(\cdot, \nabla_s \mathbf{u}) + \nabla p = \mathbf{f} \quad \text{in } \Omega \times (0, t_F), \quad (\text{A.1a})$$

$$c_0 d_t p + d_t \nabla \cdot \mathbf{u} - \nabla \cdot (\boldsymbol{\kappa}(\cdot) \nabla p) = g \quad \text{in } \Omega \times (0, t_F), \quad (\text{A.1b})$$



where  $\nabla_s$  denotes the symmetric gradient,  $c_0 \geq 0$  is the constrained specific storage coefficient, and  $\kappa : \Omega \rightarrow (0, \bar{\kappa}]$  is the scalar-valued permeability field. Equations (A.1a) and (A.1b) express, respectively, the momentum equilibrium and the fluid mass balance. For the sake of simplicity, we assume that  $\kappa$  is piecewise constant on a partition  $P_\Omega$  of  $\Omega$  into bounded disjoint polyhedra and we consider the following homogeneous boundary conditions:

$$\mathbf{u} = \mathbf{0} \quad \text{on } \partial\Omega \times (0, t_F), \quad (\text{A.1c})$$

$$(\kappa(\cdot)\nabla p) \cdot \mathbf{n} = 0 \quad \text{on } \partial\Omega \times (0, t_F). \quad (\text{A.1d})$$

The treatment of more general permeability fields and boundary conditions is possible up to minor modifications. Initial conditions are set prescribing  $\mathbf{u}(\cdot, 0) = \mathbf{u}^0$  and, if  $c_0 > 0$ ,  $p(\cdot, 0) = p^0$ . In the incompressible case  $c_0 = 0$ , we also need the following compatibility condition on  $g$  and zero-average constraint on  $p$ :

$$\int_{\Omega} g(\cdot, t) = 0 \quad \text{and} \quad \int_{\Omega} p(\cdot, t) = 0 \quad \forall t \in (0, t_F).$$

We assume that the symmetric stress tensor  $\sigma : \Omega \times \mathbb{R}_{\text{sym}}^{d \times d} \rightarrow \mathbb{R}_{\text{sym}}^{d \times d}$  is a Caratheodory function such that there exist real numbers  $\bar{\sigma}, \underline{\sigma} \in (0, +\infty)$  and, for a.e.  $\mathbf{x} \in \Omega$ , and all  $\boldsymbol{\tau}, \boldsymbol{\eta} \in \mathbb{R}_{\text{sym}}^{d \times d}$ , the following conditions hold:

$$\|\sigma(\mathbf{x}, \boldsymbol{\tau}) - \sigma(\mathbf{x}, \mathbf{0})\|_{d \times d} \leq \bar{\sigma} \|\boldsymbol{\tau}\|_{d \times d}, \quad (\text{growth}) \quad (\text{A.2a})$$

$$\sigma(\mathbf{x}, \boldsymbol{\tau}) : \boldsymbol{\tau} \geq \underline{\sigma} \|\boldsymbol{\tau}\|_{d \times d}^2, \quad (\text{coercivity}) \quad (\text{A.2b})$$

$$(\sigma(\mathbf{x}, \boldsymbol{\tau}) - \sigma(\mathbf{x}, \boldsymbol{\eta})) : (\boldsymbol{\tau} - \boldsymbol{\eta}) \geq 0, \quad (\text{monotonicity}) \quad (\text{A.2c})$$

where  $\boldsymbol{\tau} : \boldsymbol{\eta} := \sum_{i,j=1}^d \tau_{i,j} \eta_{i,j}$  and  $\|\boldsymbol{\tau}\|_{d \times d}^2 = \boldsymbol{\tau} : \boldsymbol{\tau}$ .

## A.2 Mesh and notation

Denote by  $\mathcal{H} \subset \mathbb{R}_*^+$  a countable set having 0 as unique accumulation point. We consider refined mesh sequences  $(\mathcal{T}_h)_{h \in \mathcal{H}}$  where each  $\mathcal{T}_h$  is a finite collection of disjoint open polyhedral elements  $T$  with boundary  $\partial T$  such that  $\bar{\Omega} = \bigcup_{T \in \mathcal{T}_h} \bar{T}$  and  $h = \max_{T \in \mathcal{T}_h} h_T$  with  $h_T$  diameter of  $T$ . We assume that mesh regularity holds in the sense of [76, Definition 1.38] and that, for all  $h \in \mathcal{H}$ ,  $\mathcal{T}_h$  is compatible with the partition  $P_\Omega$ , so that jumps of the permeability coefficient do not occur inside mesh elements. Mesh faces are hyperplanar subsets of  $\bar{\Omega}$  with positive  $(d-1)$ -dimensional Hausdorff measure and disjoint interiors. Interfaces are collected in the set  $\mathcal{F}_h^i$ , boundary faces in  $\mathcal{F}_h^b$ , and we assume that  $\mathcal{F}_h := \mathcal{F}_h^i \cup \mathcal{F}_h^b$  is such that  $\bigcup_{T \in \mathcal{T}_h} \partial T = \bigcup_{F \in \mathcal{F}_h} F$ . For all  $T \in \mathcal{T}_h$ ,  $\mathcal{F}_T := \{F \in \mathcal{F}_h : F \subset \partial T\}$  denotes the set of faces contained in  $\partial T$  and, for all  $F \in \mathcal{F}_T$ ,  $\mathbf{n}_{TF}$  is the unit normal to  $F$  pointing out of  $T$ .

For  $X \subset \bar{\Omega}$ , we denote by  $\|\cdot\|_X$  the norm in  $L^2(X; \mathbb{R})$ ,  $L^2(X; \mathbb{R}^d)$ , and  $L^2(X; \mathbb{R}^{d \times d})$ . For  $l \geq 0$ , the space  $\mathbb{P}_d^l(X; \mathbb{R})$  is spanned by the restriction to  $X$  of polynomials of total degree  $l$ . On regular mesh sequences, we have the following optimal approximation property for the  $L^2$ -projector  $\pi_X^l : L^1(X; \mathbb{R}) \rightarrow \mathbb{P}_d^l(X; \mathbb{R})$ : There exists  $C_{\text{ap}} > 0$  such that, for all  $h \in \mathcal{H}$ , all  $T \in \mathcal{T}_h$ , all  $s \in \{1, \dots, l+1\}$ , and all  $v \in H^s(T; \mathbb{R})$ ,

$$|v - \pi_T^l v|_{H^m(T; \mathbb{R})} \leq C_{\text{ap}} h_T^{s-m} |v|_{H^s(T; \mathbb{R})} \quad \forall m \in \{0, \dots, s\}. \quad (\text{A.3})$$

In what follows, for an integer  $l \geq 0$ , we denote by  $\mathbb{P}_d^l(\mathcal{T}_h; \mathbb{R})$ ,  $\mathbb{P}_d^l(\mathcal{T}_h; \mathbb{R}^d)$ , and  $\mathbb{P}_d^l(\mathcal{T}_h; \mathbb{R}^{d \times d})$ , respectively, the space of scalar-, vector-, and tensor-valued broken polynomials of total degree  $l$  on  $\mathcal{T}_h$ . The space of broken vector-valued polynomials of total degree  $l$  on the mesh skeleton is denoted by  $\mathbb{P}_d^l(\mathcal{F}_h; \mathbb{R}^d)$ .

### A.3 Discretization

In this section we define the discrete counterparts of the elasticity and Darcy operators and of the hydro-mechanical coupling terms.

#### A.3.1 Nonlinear elasticity operator

The discretization of the nonlinear elasticity operator is based on the Hybrid High-Order method of [74]. Let a polynomial degree  $k \geq 1$  be fixed. The degrees of freedom (DOFs) for the displacement are collected in the space  $\underline{\mathbf{U}}_h^k := \mathbb{P}_d^k(\mathcal{T}_h; \mathbb{R}^d) \times \mathbb{P}_d^k(\mathcal{F}_h; \mathbb{R}^d)$ . To account for the Dirichlet condition (A.1c) we define the subspace

$$\underline{\mathbf{U}}_{h,D}^k := \left\{ \underline{\mathbf{v}}_h := ((\mathbf{v}_T)_{T \in \mathcal{T}_h}, (\mathbf{v}_F)_{F \in \mathcal{F}_h}) \in \underline{\mathbf{U}}_h^k : \mathbf{v}_F = \mathbf{0} \quad \forall F \in \mathcal{F}_h^b \right\},$$

equipped with the discrete strain norm  $\|\cdot\|_{\epsilon,h}$  defined in (2.13). The DOFs corresponding to a function  $\mathbf{v} \in H_0^1(\Omega; \mathbb{R}^d)$  are obtained by means of the reduction map  $\underline{I}_h^k : H_0^1(\Omega; \mathbb{R}^d) \rightarrow \underline{\mathbf{U}}_{h,D}^k$  such that  $\underline{I}_h^k \mathbf{v} := ((\boldsymbol{\pi}_T^k \mathbf{v})_{T \in \mathcal{T}_h}, (\boldsymbol{\pi}_F^k \mathbf{v})_{F \in \mathcal{F}_h})$ . Using the  $H^1$ -stability of the  $L^2$ -projector and the trace inequality [76, Lemma 1.49], we infer the existence of  $C_{\text{st}} > 0$  independent of  $h$  such that, for all  $\mathbf{v} \in H_0^1(\Omega; \mathbb{R}^d)$ ,

$$\|\underline{I}_h^k \mathbf{v}\|_{\epsilon,h} \leq C_{\text{st}} \|\mathbf{v}\|_{H^1(\Omega; \mathbb{R}^d)}. \quad (\text{A.4})$$

For all  $T \in \mathcal{T}_h$ , we denote by  $\underline{\mathbf{U}}_T^k$  and  $\underline{I}_T^k$  the restrictions to  $T$  of  $\underline{\mathbf{U}}_h^k$  and  $\underline{I}_h^k$ , and we define the local symmetric gradient reconstruction  $\mathbf{G}_{s,T}^k : \underline{\mathbf{U}}_T^k \rightarrow \mathbb{P}_d^k(T; \mathbb{R}_{\text{sym}}^{d \times d})$  as the unique solution of the pure traction problem: For a given  $\underline{\mathbf{v}}_T = (\mathbf{v}_T, (\mathbf{v}_F)_{F \in \mathcal{F}_T}) \in \underline{\mathbf{U}}_T^k$ , find  $\mathbf{G}_{s,T}^k \underline{\mathbf{v}}_T \in \mathbb{P}_d^k(T; \mathbb{R}_{\text{sym}}^{d \times d})$  such that, for all  $\boldsymbol{\tau} \in \mathbb{P}_d^k(T; \mathbb{R}_{\text{sym}}^{d \times d})$ ,

$$\int_T \mathbf{G}_{s,T}^k \underline{\mathbf{v}}_T : \boldsymbol{\tau} = - \int_T \mathbf{v}_T \cdot (\boldsymbol{\nabla} \cdot \boldsymbol{\tau}) + \sum_{F \in \mathcal{F}_T} \int_F \mathbf{v}_F \cdot (\boldsymbol{\tau} \mathbf{n}_{TF}). \quad (\text{A.5})$$

The definition of  $\mathbf{G}_{s,T}^k$  is justified by the commuting property (4.19) that, combined with (A.3), shows that  $\mathbf{G}_{s,T}^k \underline{I}_T^k$  has optimal approximation properties in  $\mathbb{P}_d^k(T; \mathbb{R}_{\text{sym}}^{d \times d})$ . From  $\mathbf{G}_{s,T}^k$  we define the local displacement reconstruction operator  $\mathbf{r}_T^{k+1} : \underline{\mathbf{U}}_T^k \rightarrow \mathbb{P}_d^{k+1}(T; \mathbb{R}^d)$  as in (2.5).

The discretization of the nonlinear elasticity operator is realized by the function  $a_h : \underline{\mathbf{U}}_h^k \times \underline{\mathbf{U}}_h^k \rightarrow \mathbb{R}$  defined such that, for all  $\underline{\mathbf{w}}_h, \underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_h^k$ ,

$$a_h(\underline{\mathbf{w}}_h, \underline{\mathbf{v}}_h) := \sum_{T \in \mathcal{T}_h} \left( \int_T \boldsymbol{\sigma}(\cdot, \mathbf{G}_{s,T}^k \underline{\mathbf{u}}_T) : \mathbf{G}_{s,T}^k \underline{\mathbf{v}}_T + \sum_{F \in \mathcal{F}_T} \frac{\gamma}{h_F} \int_F \Delta_{TF}^k \underline{\mathbf{u}}_T \cdot \Delta_{TF}^k \underline{\mathbf{v}}_T \right), \quad (\text{A.6})$$

where we penalize in a least-square sense the face-based residual  $\Delta_{TF}^k : \underline{U}_T^k \rightarrow \mathbb{P}_d^k(F; \mathbb{R}^d)$  such that, for all  $T \in \mathcal{T}_h$ , all  $\underline{v}_T \in \underline{U}_T^k$ , and all  $F \in \mathcal{F}_T$ ,

$$\Delta_{TF}^k \underline{v}_T := \boldsymbol{\pi}_F^k(\mathbf{r}_T^{k+1} \underline{v}_T - \mathbf{v}_F) - \boldsymbol{\pi}_T^k(\mathbf{r}_T^{k+1} \underline{v}_T - \mathbf{v}_T).$$

This definition ensures that  $\Delta_{TF}^k$  vanishes whenever its argument is of the form  $\underline{I}_T^k \mathbf{w}$  with  $\mathbf{w} \in \mathbb{P}_d^{k+1}(T; \mathbb{R}^d)$ , a crucial property to obtain high-order error estimates (cf. Theorems 2.12, 3.16, and 4.17). A possible choice for the scaling parameter  $\gamma > 0$  in (A.6) is  $\gamma = \underline{\sigma}$ . For all  $\underline{v}_h \in \underline{U}_{h,D}^k$ , it holds (the proof follows from [74, Lemma 4]):

$$C_{\text{eq}}^{-1} \|\underline{v}_h\|_{\epsilon,h}^2 \leq \sum_{T \in \mathcal{T}_h} \left( \|\mathbf{G}_{s,T}^k \underline{v}_T\|_T^2 + \sum_{F \in \mathcal{F}_T} \frac{\gamma}{h_F} \|\Delta_{TF}^k \underline{v}_T\|_F^2 \right) \leq C_{\text{eq}} \|\underline{v}_h\|_{\epsilon,h}^2,$$

where  $C_{\text{eq}} > 0$  is independent of  $h$ . By (A.2b), this implies the coercivity of  $a_h$ .

### A.3.2 Darcy operator

The discretization of the Darcy operator is based on the Symmetric Weighted Interior Penalty method of [77], cf. also [76, Sec. 4.5]. At each time step, the discrete pore pressure is sought in the broken polynomial space  $P_h^k$  defined in (2.17). For all  $F \in \mathcal{F}_h^i$  and all  $q_h \in \mathbb{P}_d^k(\mathcal{T}_h; \mathbb{R})$ , we define the jump and average operators such that, denoting by  $q_T$  and  $\kappa_T$  the restrictions of  $q_h$  and  $\kappa$  to an element  $T \in \mathcal{T}_h$ ,

$$[q_h]_F := q_{T_{F,1}} - q_{T_{F,2}}, \quad \{q_h\}_F := \frac{\kappa_{T_{F,2}}}{\kappa_{T_{F,1}} + \kappa_{T_{F,2}}} q_{T_{F,1}} + \frac{\kappa_{T_{F,1}}}{\kappa_{T_{F,1}} + \kappa_{T_{F,2}}} q_{T_{F,2}},$$

where  $T_{F,1}, T_{F,2} \in \mathcal{T}_h$  are such that  $F \subset T_{F,1} \cap T_{F,2}$ . The bilinear form  $c_h$  on  $P_h^k \times P_h^k$  is defined such that, for all  $q_h, r_h \in P_h^k$ ,

$$\begin{aligned} c_h(r_h, q_h) := & \int_{\Omega} \boldsymbol{\kappa} \nabla_h r_h \cdot \nabla_h q_h + \sum_{F \in \mathcal{F}_h^i} \frac{2\varsigma \kappa_{T_{F,1}} \kappa_{T_{F,2}}}{h_F (\kappa_{T_{F,1}} + \kappa_{T_{F,2}})} \int_F [r_h]_F [q_h]_F \\ & - \sum_{F \in \mathcal{F}_h^i} \int_F ([r_h]_F \{\boldsymbol{\kappa} \nabla_h q_h\}_F + [q_h]_F \{\boldsymbol{\kappa} \nabla_h r_h\}_F) \cdot \mathbf{n}_{T_{F,1}F}, \end{aligned} \quad (\text{A.7})$$

where  $\nabla_h$  denotes the broken gradient and  $\varsigma > 0$  is a user-defined penalty parameter chosen large enough to ensure the coercivity of  $c_h$  (cf. [76, Lemma 4.51]). In the numerical tests of Sec. A.5, we took  $\varsigma = (N_{\partial} + 0.1)k^2$ , with  $N_{\partial}$  equal to the maximum number of faces between the elements in  $\mathcal{T}_h$ . The fact that the boundary terms only appear on internal faces in (A.7) reflects the Neumann boundary condition (A.1d).

### A.3.3 Hydro-mechanical coupling

The hydro-mechanical coupling is realized by means of the bilinear form  $b_h$  on  $\underline{U}_{h,D}^k \times \mathbb{P}_d^k(\mathcal{T}_h; \mathbb{R})$  such that, for all  $\underline{v}_h \in \underline{U}_{h,D}^k$  and all  $q_h \in \mathbb{P}_d^k(\mathcal{T}_h; \mathbb{R})$ ,

$$b_h(\underline{v}_h, q_h) := - \sum_{T \in \mathcal{T}_h} \int_T \mathbf{G}_{s,T}^k \underline{v}_T : q_h|_T \mathbf{I}_d,$$

where  $\mathbf{I}_d \in \mathbb{R}_{\text{sym}}^{d \times d}$  is the identity matrix. A simple verification shows that there exists  $C_{\text{bd}} > 0$  independent of  $h$  such that  $b_h(\underline{\mathbf{v}}_h, q_h) \leq C_{\text{bd}} \|\underline{\mathbf{v}}_h\|_{\epsilon, h} \|q_h\|_{\Omega}$ . Additionally, using definition (A.5) of  $\mathbf{G}_{s, T}^k$ , it can be proved that, for all  $\underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_{h, D}^k$ ,  $b_h(\underline{\mathbf{v}}_h, 1) = 0$ . The following inf-sup condition expresses the stability of the coupling:

**Proposition A.1.** *There is a real  $\beta$  independent of  $h$  such that, for all  $q_h \in \mathbb{P}_{d0}^k(\mathcal{T}_h; \mathbb{R})$ ,*

$$\|q_h\|_{\Omega} \leq \beta \sup_{\underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_{h, D}^k \setminus \{\mathbf{0}\}} \frac{b_h(\underline{\mathbf{v}}_h, q_h)}{\|\underline{\mathbf{v}}_h\|_{\epsilon, h}}. \quad (\text{A.8})$$

*Proof.* Let  $q_h \in \mathbb{P}_{d0}^k(\mathcal{T}_h; \mathbb{R})$ . There is  $\mathbf{v}_{q_h} \in H_0^1(\Omega; \mathbb{R}^d)$  such that  $\nabla \cdot \mathbf{v}_{q_h} = q_h$  and  $\|\mathbf{v}_{q_h}\|_{H^1(\Omega; \mathbb{R}^d)} \leq C_{\text{sj}} \|q_h\|_{\Omega}$ , with  $C_{\text{sj}} > 0$  independent of  $h$ . Owing to (A.4) we get

$$\|\underline{\mathbf{I}}_h^k \mathbf{v}_{q_h}\|_{\epsilon, h} \leq C_{\text{st}} \|\mathbf{v}_{q_h}\|_{H^1(\Omega; \mathbb{R}^d)} \leq C_{\text{st}} C_{\text{sj}} \|q_h\|_{\Omega}.$$

Therefore, using the commuting property (4.19), denoting by  $\mathbf{S}$  the supremum in (A.8), and using the previous inequality, it is inferred that

$$\|q_h\|_{\Omega}^2 = \sum_{T \in \mathcal{T}_h} \int_T (\nabla_s \mathbf{v}_{q_h} : q_h \mathbf{I}_d)|_T = -b_h(\underline{\mathbf{I}}_h^k \mathbf{v}_{q_h}, q_h) \leq \mathbf{S} \|\underline{\mathbf{I}}_h^k \mathbf{v}_{q_h}\|_{\epsilon, h} \leq C_{\text{st}} C_{\text{sj}} \mathbf{S} \|q_h\|_{\Omega}. \quad \square$$

If a HHO discretization were used also for the Darcy operator, only cell DOFs would be controlled by the inf-sup condition.

## A.4 Formulation of the method

For the time discretization, we consider a uniform mesh of the time interval  $(0, t_F)$  of step  $\tau := t_F/N$  with  $N \in \mathbb{N}^*$ , and introduce the discrete times  $t^n := n\tau$  for all  $0 \leq n \leq N$ . For any  $\varphi \in C^1([0, t_F]; V)$ , we set  $\varphi^n := \varphi(t^n) \in V$  and let, for all  $1 \leq n \leq N$ ,  $\delta_t \varphi^n := (\varphi^n - \varphi^{n-1})/\tau \in V$ . For all  $1 \leq n \leq N$ , the discrete solution  $(\underline{\mathbf{u}}_h^n, p_h^n) \in \underline{\mathbf{U}}_{h, D}^k \times P_h^k$  at time  $t^n$  is such that, for all  $(\underline{\mathbf{v}}_h, q_h) \in \underline{\mathbf{U}}_{h, D}^k \times \mathbb{P}_d^k(\mathcal{T}_h; \mathbb{R})$ ,

$$\begin{aligned} a_h(\underline{\mathbf{u}}_h^n, \underline{\mathbf{v}}_h) + b_h(\underline{\mathbf{v}}_h, p_h^n) &= \sum_{T \in \mathcal{T}_h} \int_T \mathbf{f}^n \cdot \mathbf{v}_T, \\ c_0 \int_{\Omega} (\delta_t p_h^n) q_h - b_h(\delta_t \underline{\mathbf{u}}_h^n, q_h) + c_h(p_h^n, q_h) &= \int_{\Omega} g^n q_h. \end{aligned} \quad (\text{A.9})$$

If  $c_0 = 0$ , we set the initial discrete displacement as  $\underline{\mathbf{u}}_h^0 = \underline{\mathbf{I}}_h^k \mathbf{u}^0$ . If  $c_0 > 0$ , the usual way to enforce the initial condition is to compute a displacement from the given initial pressure  $p^0$ . We let  $p_h^0 := \pi_h^k p^0$  and set  $\underline{\mathbf{u}}_h^0 \in \underline{\mathbf{U}}_{h, D}^k$  as the solution of

$$a_h(\underline{\mathbf{u}}_h^0, \underline{\mathbf{v}}_h) = \sum_{T \in \mathcal{T}_h} \int_T \mathbf{f}^0 \cdot \mathbf{v}_T - b_h(\underline{\mathbf{v}}_h, p_h^0) \quad \forall \underline{\mathbf{v}}_h \in \underline{\mathbf{U}}_{h, D}^k.$$

At each time step  $n$  the discrete nonlinear equations (A.9) are solved by the Newton's method using as initial guess the solution at step  $n - 1$ . At each Newton's iteration the Jacobian matrix is computed analytically and in the linearized system the displacement element unknowns can be statically condensed (cf. Section 2.5).

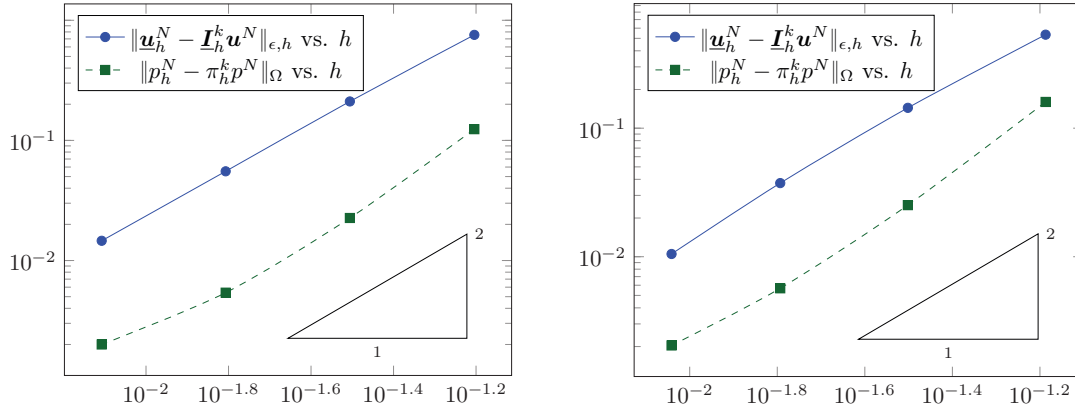


Figure A.1: Convergence tests on a Cartesian mesh family (left) and on a Voronoi mesh family (right).

## A.5 Numerical results

We consider a regular exact solution in order to assess the convergence of the method for polynomial degree  $k = 1$ . Specifically, we solve problem (A.1) in the square domain  $\Omega = (0, 1)^2$  with  $t_F = 1$  and physical parameters  $c_0 = 0$  and  $\kappa = 1$ . As nonlinear constitutive law we take the Hencky–Mises relation given by

$$\sigma(\nabla_s \mathbf{u}) = (2 \exp^{-\text{dev}(\nabla_s \mathbf{u})} - 1) \text{tr}(\nabla_s \mathbf{u}) \mathbf{I}_d + (4 - 2 \exp^{-\text{dev}(\nabla_s \mathbf{u})}) \nabla_s \mathbf{u},$$

where  $\text{tr}(\boldsymbol{\tau}) := \boldsymbol{\tau} : \mathbf{I}_d$  and  $\text{dev}(\boldsymbol{\tau}) = \text{tr}(\boldsymbol{\tau}^2) - \frac{1}{d} \text{tr}(\boldsymbol{\tau})^2$  are the trace and deviatoric operators. It can be checked that the previous stress-strain relation satisfies (A.2). The exact displacement  $\mathbf{u}$  and exact pressure  $p$  are given by

$$\begin{aligned} \mathbf{u}(\mathbf{x}, t) &= t^2 (\sin(\pi x_1) \sin(\pi x_2), \sin(\pi x_1) \sin(\pi x_2)), \\ p(\mathbf{x}, t) &= -\pi^{-1} t (\sin(\pi x_1) \cos(\pi x_2) + \cos(\pi x_1) \sin(\pi x_2)). \end{aligned}$$

The volumetric load  $\mathbf{f}$ , the source term  $g$ , and the boundary conditions are inferred from the exact solutions. The time step  $\tau$  on the coarsest mesh is 0.2 and it decreases with the mesh size  $h$  according to  $\tau_1/\tau_2 = 2h_1/h_2$ . In Fig. A.1 we display the convergence results obtained on two mesh families. The method exhibits second order convergence with respect to the mesh size  $h$  for both the energy norm of the displacement and the  $L^2$ -norm of the pressure at final time  $t_F$ .

# Appendix B

---

## A posteriori error estimation for nonlinear elasticity

---

This appendix has been published in the following peer-reviewed journal (see [36])

**Computational Methods in Applied Mathematics,**  
*Published online: 20 June 2018, DOI: 10.1515/cmam-2018-0012.*

### Contents

---

<b>B.1 Introduction</b>	<b>147</b>
<b>B.2 Setting</b>	<b>149</b>
<b>B.3 Equilibrated stress reconstruction</b>	<b>150</b>
B.3.1 Patchwise construction in the Arnold–Falk–Winther spaces	150
B.3.2 Discretization and linearization error stress reconstructions	152
<b>B.4 A posteriori error estimate and adaptive algorithm</b>	<b>154</b>
B.4.1 Guaranteed upper bound	154
B.4.2 Adaptive algorithm	156
B.4.3 Local efficiency	157
<b>B.5 Numerical results</b>	<b>158</b>
B.5.1 L-shaped domain	158
B.5.2 Notched specimen plate	161

---

### B.1 Introduction

In this appendix we develop equilibrated  $H(\text{div})$ -conforming stress tensor reconstructions for a class of (linear and) nonlinear elasticity problems in the small deformation regime. Based on these reconstructions, we can derive an a posteriori error estimate distinguishing the discretization and linearization errors, as proposed in [95] for nonlinear diffusion problems. Thanks to this distinction we can, at each iteration, compare these two error contributions

and stop the linearization algorithm once its contribution is negligible compared to the discretization error.

Let  $\Omega \in \mathbb{R}^d$ ,  $d \in \{2, 3\}$ , be a bounded, simply connected polyhedron, which is occupied by a body subjected to a volumetric force field  $\mathbf{f} \in [L^2(\Omega)]^d$ . For the sake of simplicity, we assume that the body is fixed on its boundary  $\partial\Omega$ . The nonlinear elasticity problem consists in finding a vector-valued displacement field  $\mathbf{u} : \Omega \rightarrow \mathbb{R}^d$  solving

$$-\nabla \cdot \boldsymbol{\sigma}(\nabla_s \mathbf{u}) = \mathbf{f} \quad \text{in } \Omega, \quad (\text{B.1a})$$

$$\mathbf{u} = \mathbf{0} \quad \text{on } \partial\Omega, \quad (\text{B.1b})$$

where  $\nabla_s \mathbf{u} = \frac{1}{2}((\nabla \mathbf{u})^T + \nabla \mathbf{u})$  denotes the symmetric gradient and expresses the strain tensor associated to  $\mathbf{u}$ . The stress-strain law  $\boldsymbol{\sigma} : \Omega \times \mathbb{R}_{\text{sym}}^{d \times d} \rightarrow \mathbb{R}_{\text{sym}}^{d \times d}$  is assumed to satisfy regularity requirements inspired by [33, 69]. Problem (B.1) describes the mechanical behavior of soft materials [191] and metal alloys [169]. In these applications, the solution is often approximated using  $H^1$ -conforming finite elements. For nonlinear mechanical behavior laws, the resulting discrete nonlinear equation can then be solved using an iterative linearization algorithm yielding at each iteration a linear algebraic system to be solved, until the residual of the nonlinear equation lies under a predefined threshold.

The a posteriori error estimate is based on equilibrated stress reconstructions. It is well known that, in contrast to the analytical solution, the discrete stress tensor resulting from the conforming finite element method does not have continuous normal components across mesh interfaces, and that its divergence is not locally in equilibrium with the source term  $\mathbf{f}$  on mesh elements. We consider here the stress reconstruction proposed in [175] for linear elasticity to restore these two properties. This reconstruction uses the Arnold–Falk–Winther mixed finite element spaces [7], leading to weakly symmetric tensors. In Section B.3 we apply this reconstruction to the nonlinear case by constructing two stress tensors: one playing the role of the discrete stress and one expressing the linearization error. They are obtained by summing up the solutions of minimization problems on cell patches around each mesh vertex, so that they are  $H(\text{div})$ -conforming and the sum of the two reconstructions verifies locally the mechanical equilibrium (B.1a). The patch-wise equilibration technique was introduced in [38, 67] for the Poisson problem using the Raviart–Thomas finite element spaces. In [83] it is extended to linear elasticity without any symmetry constraint by using linewise Raviart–Thomas reconstructions. Elementwise reconstructions from local Neumann problems requiring some pre-computations to determine the fluxes to obtain an equilibrated stress tensor can be found in [4, 51, 134, 161], whereas in [158] the direct prescription of the degrees of freedom in the Arnold–Winther finite element space is considered.

Based on the equilibrated stress reconstructions, we develop the a posteriori error estimate in Section B.4 and prove that it is efficient, meaning that, up to a generic constant, it is also a local lower bound for the error. The idea goes back to [170] and was advanced amongst others by [2, 132, 133, 172] for the upper bound. Local lower error bounds are derived in [38, 67, 94, 97, 144]. Using equilibrated fluxes for a posteriori error estimation offers several advantages. The first one is, as mentioned above, the possible distinction of error components by expressing them in terms of fluxes. Secondly, the error upper bound is obtained with computable constants. In our case these constants depend only on the parameters of the

stress-strain relation. Thirdly, since the estimate is based on the discrete stress (and not the displacement), it does not depend on the mechanical behavior law (except for the constant in the upper bound). In addition, equilibrated error estimates were proven to be polynomial-degree robust for several linear problems in 2D, as the Poisson problem in [37, 97], linear elasticity in [83] and the related Stokes problem in [49] and recently in 3D in [96].

The material is organized as follows. In Section B.2 we introduce the weak and the discrete formulations of problem (B.1), along with some useful notation. In Section B.3 we present the equilibrated stress tensor reconstructions, first assuming that we solve the nonlinear discrete problem exactly, and then distinguish its discrete and its linearization error part at each iteration of a linearization solver. In Section B.4 we derive the a posteriori error estimate distinguishing the different error components. We then propose an algorithm equilibrating the error sources using adaptive stopping criteria for the linearization and adaptive remeshing. We finally show the efficiency of the error estimate. In Section B.5 we evaluate the performance of the estimates. We perform a linear test with analytical solution in order to compare the error estimator and the analytical error and to show that effectivity indices are close to one. We then apply the proposed adaptive algorithm to application-oriented tests with nonlinear stress-strain relations.

## B.2 Setting

In this section we write the weak and the considered discrete formulation of problem (B.1). In the following, we will adopt the notation that bold greek letters are tensor-valued functions and bold latin letters are vector-valued functions. For  $X \subset \overline{\Omega}$ , we respectively denote by  $(\cdot, \cdot)_X$  and  $\|\cdot\|_X$  the standard inner product and norm in  $L^2(X)$ , with the convention that the subscript is omitted whenever  $X = \Omega$ . The same notation is used in the vector- and tensor-valued cases.  $[H^1(\Omega)]^d$  and  $\mathbb{H}(\text{div}, \Omega)$  stand for the Sobolev spaces composed of vector-valued  $[L^2(\Omega)]^d$  functions with weak gradient in  $[L^2(\Omega)]^{d \times d}$ , and tensor-valued  $[L^2(\Omega)]^{d \times d}$  functions with weak divergence in  $[L^2(\Omega)]^d$ , respectively. In what follows, we consider stress-strain functions  $\boldsymbol{\sigma}$  satisfying the following assumption.

**Assumption B.1** (Stress-strain relation). The symmetric stress tensor  $\boldsymbol{\sigma} : \mathbb{R}_{\text{sym}}^{d \times d} \rightarrow \mathbb{R}_{\text{sym}}^{d \times d}$  is continuous on  $\mathbb{R}_{\text{sym}}^{d \times d}$ , satisfies the growth and strong monotonicity assumptions (3.2c) and (3.40b), and  $\boldsymbol{\sigma}(\mathbf{0}) = \mathbf{0}$ .

Multiplying (B.1a) by a test function  $\boldsymbol{v} \in [H_0^1(\Omega)]^d$  and integrating by parts yields

$$(\boldsymbol{\sigma}(\nabla_s \boldsymbol{u}), \nabla_s \boldsymbol{v}) = (\boldsymbol{f}, \boldsymbol{v}). \quad (\text{B.2})$$

Owing to (3.2c), for all  $\boldsymbol{v}, \boldsymbol{w} \in [H^1(\Omega)]^d$ , the form  $a(\boldsymbol{v}, \boldsymbol{w}) := (\boldsymbol{\sigma}(\nabla_s \boldsymbol{v}), \nabla_s \boldsymbol{w})$  is well defined and, from equation (B.2), we can derive the following weak formulation of (B.1):

$$\text{Given } \boldsymbol{f} \in [L^2(\Omega)]^d, \text{ find } \boldsymbol{u} \in [H_0^1(\Omega)]^d \text{ s.t., } \forall \boldsymbol{v} \in [H_0^1(\Omega)]^d, a(\boldsymbol{u}, \boldsymbol{v}) = (\boldsymbol{f}, \boldsymbol{v}). \quad (\text{B.3})$$

From (B.3) it is clear that the analytical stress tensor  $\boldsymbol{\sigma}(\nabla_s \boldsymbol{u})$  lies in the space  $\mathbb{H}_s(\text{div}, \Omega) := \{\boldsymbol{\tau} \in [L^2(\Omega)]^{d \times d} \mid \nabla \cdot \boldsymbol{\tau} \in [L^2(\Omega)]^d \text{ and } \boldsymbol{\tau} \text{ is symmetric}\}$ .

The discretization of (B.3) is based on a conforming triangulation  $\mathcal{T}_h$  of  $\Omega$ , i.e. a set of closed triangles or tetrahedra with union  $\overline{\Omega}$  and such that, for any distinct  $T_1, T_2 \in \mathcal{T}_h$ , the set



$T_1 \cap T_2$  is either a common edge, a vertex, the empty set or, if  $d = 3$ , a common face. The set of vertices of the mesh is denoted by  $\mathcal{V}_h$ ; it is decomposed into interior vertices  $\mathcal{V}_h^{\text{int}}$  and boundary vertices  $\mathcal{V}_h^{\text{ext}}$ . For all  $a \in \mathcal{V}_h$ ,  $\mathcal{T}_a$  is the patch of elements sharing the vertex  $a$ ,  $\omega_a$  the corresponding open subdomain in  $\Omega$  and  $\mathcal{V}_a$  the set of vertices in  $\omega_a$ . For all  $T \in \mathcal{T}_h$ ,  $\mathcal{V}_T$  denotes the set of vertices of  $T$ ,  $h_T$  its diameter and  $\mathbf{n}_T$  its unit outward normal vector.

For all  $p \in \mathbb{N}$  and all  $T \in \mathcal{T}_h$ , we denote by  $\mathbb{P}^p(T)$  the space of  $d$ -variate polynomials in  $T$  of total degree at most  $p$  and by  $\mathbb{P}^p(\mathcal{T}_h) = \{\varphi \in L^2(\Omega) \mid \varphi|_T \in \mathbb{P}^p(T) \forall T \in \mathcal{T}_h\}$  the corresponding broken space over  $\mathcal{T}_h$ . In what follows we will focus on conforming discretizations of problem (B.2) of polynomial degree  $p \geq 2$  to avoid numerical locking, cf [196]. The discrete formulation reads: find  $\mathbf{u}_h \in [H_0^1(\Omega)]^d \cap [\mathbb{P}^p(\mathcal{T}_h)]^d$  such that

$$\forall \mathbf{v}_h \in [H_0^1(\Omega)]^d \cap [\mathbb{P}^p(\mathcal{T}_h)]^d, \quad a(\mathbf{u}_h, \mathbf{v}_h) = (\mathbf{f}, \mathbf{v}_h). \quad (\text{B.4})$$

This problem is usually solved using some iterative linearization algorithm defining at each iteration  $k \geq 1$  a linear approximation  $\sigma_{\text{lin}}^{k-1}$  of  $\sigma$ . Then the linearized formulation reads: find  $\mathbf{u}_h^k \in [H_0^1(\Omega)]^d \cap [\mathbb{P}^p(\mathcal{T}_h)]^d$  such that

$$\forall \mathbf{v}_h \in [H_0^1(\Omega)]^d \cap [\mathbb{P}^p(\mathcal{T}_h)]^d, \quad (\sigma_{\text{lin}}^{k-1}(\nabla_s \mathbf{u}_h^k), \nabla_s \mathbf{v}_h) = (\mathbf{f}, \mathbf{v}_h). \quad (\text{B.5})$$

For the Newton algorithm the linearized stress  $\sigma_{\text{lin}}^{k-1}$  is defined, for all  $\boldsymbol{\eta} \in \mathbb{R}_{\text{sym}}^{d \times d}$ , such that

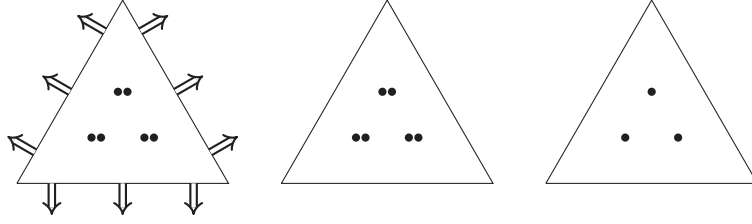
$$\sigma_{\text{lin}}^{k-1}(\boldsymbol{\eta}) := \frac{\partial \sigma(\boldsymbol{\tau})}{\partial \boldsymbol{\tau}} \Big|_{\boldsymbol{\tau} = \nabla_s \mathbf{u}_h^{k-1}} (\boldsymbol{\eta} - \nabla_s \mathbf{u}_h^{k-1}) + \sigma(\nabla_s \mathbf{u}_h^{k-1}).$$

## B.3 Equilibrated stress reconstruction

In general, the discrete stress tensor  $\sigma(\nabla_s \mathbf{u}_h)$  resulting from (B.4) does not lie in  $\mathbb{H}(\text{div}, \Omega)$  and thus cannot verify the equilibrium equation (B.1a). In this section we will reconstruct from  $\sigma(\nabla_s \mathbf{u}_h)$  a discrete stress tensor  $\sigma_h$  satisfying these properties. Based on this reconstruction, we then devise two equilibrated stress tensors representing the discrete stress and the linearization error respectively, which will be useful for the distinction of error components in the a posteriori error estimate of Section B.4.

### B.3.1 Patchwise construction in the Arnold–Falk–Winther spaces

Let us for now suppose that  $\mathbf{u}_h$  solves (B.4) exactly, before considering iterative linearization methods such as (B.5) in Section B.3.2. For the stress reconstruction we will use mixed finite element formulations on patches around mesh vertices in the spirit of [174, 175]. The goal is to obtain a stress tensor  $\sigma_h$  in a suitable (i.e.  $H(\text{div})$ -conforming) finite element space by summing up these local solutions. The local problems are posed such that this global stress tensor is close to the discrete stress tensor  $\sigma(\nabla_s \mathbf{u}_h)$  obtained from (B.4), and that it satisfies the mechanical equilibrium on each element. In [174] the stress tensor is reconstructed in the Arnold–Winther finite element space [8], directly providing symmetric tensors, but requiring high computational effort. As done in [175], we weaken the symmetry constraint and impose

Figure B.1: Element diagrams for  $(\Sigma_T, V_T, \Lambda_T)$  in the case  $d = q = 2$ 

it weakly, as proposed in [7]: for each element  $T \in \mathcal{T}_h$ , the local Arnold–Falk–Winther mixed finite element spaces of degree  $q \geq 1$  hinge on the Brezzi–Douglas–Marini mixed finite element spaces [41] for each line of the stress tensor and are defined by

$$\Sigma_T := [\mathbb{P}^q(T)]^{d \times d}, \quad V_T := [\mathbb{P}^{q-1}(T)]^d, \quad \Lambda_T := \{\boldsymbol{\mu} \in [\mathbb{P}^{q-1}(T)]^{d \times d} \mid \boldsymbol{\mu} = -\boldsymbol{\mu}^T\}. \quad (\text{B.6})$$

For  $q = 2$ , the degrees of freedom are displayed in Figure B.1. On a patch  $\omega_a$  the global space  $\Sigma_h(\omega_a)$  is the subspace of  $\mathbb{H}(\text{div}, \omega_a)$  composed of functions belonging piecewise to  $\Sigma_T$ . The spaces  $V_h(\omega_a)$  and  $\Lambda_h(\omega_a)$  consist of functions lying piecewise in  $V_T$  and  $\Lambda_T$  respectively, with no continuity conditions between two elements.

Let now  $q = p$ . On each patch we consider subspaces where a zero normal component is enforced on the stress tensor on the boundary of the patch, so that the sum of the local solutions will have continuous normal component across any mesh interface. Since the boundary condition (B.1b) prescribes the displacement and not the normal stress, we distinguish the case whether  $a$  is an interior vertex or a boundary vertex. If  $a \in \mathcal{V}_h^{\text{int}}$  we set

$$\begin{aligned} \Sigma_h^a &:= \{\boldsymbol{\tau}_h \in \Sigma_h(\omega_a) \mid \boldsymbol{\tau}_h \mathbf{n}_{\omega_a} = \mathbf{0} \text{ on } \partial\omega_a\}, \\ V_h^a &:= \{\mathbf{v}_h \in V_h(\omega_a) \mid (\mathbf{v}_h, \mathbf{z})_{\omega_a} = 0 \forall \mathbf{z} \in \mathbf{RM}^d\}, \\ \Lambda_h^a &:= \Lambda_h(\omega_a), \end{aligned}$$

where  $\mathbf{RM}^2 := \{\mathbf{b} + c(x_2, -x_1)^T \mid \mathbf{b} \in \mathbb{R}^2, c \in \mathbb{R}\}$  and  $\mathbf{RM}^3 := \{\mathbf{b} + \mathbf{a} \times \mathbf{x} \mid \mathbf{b} \in \mathbb{R}^3, \mathbf{a} \in \mathbb{R}^3\}$  are the spaces of rigid-body motions respectively for  $d = 2$  and  $d = 3$ . If  $a \in \mathcal{V}_h^{\text{ext}}$  we set

$$\begin{aligned} \Sigma_h^a &:= \{\boldsymbol{\tau}_h \in \Sigma_h(\omega_a) \mid \boldsymbol{\tau}_h \mathbf{n}_{\omega_a} = \mathbf{0} \text{ on } \partial\omega_a \setminus \partial\Omega\}, \\ V_h^a &:= V_h(\omega_a), \\ \Lambda_h^a &:= \Lambda_h(\omega_a). \end{aligned}$$

For each vertex  $a \in \mathcal{V}_h$  we define its hat function  $\psi_a \in \mathbb{P}^1(\mathcal{T}_h)$  as the piecewise linear function taking value one at the vertex  $a$  and zero on all other mesh vertices.

**Construction B.2** (Stress tensor reconstruction). Let  $\mathbf{u}_h$  solve (B.4). For each  $a \in \mathcal{V}_h$  find  $(\boldsymbol{\sigma}_h^a, \mathbf{r}_h^a, \boldsymbol{\lambda}_h^a) \in \Sigma_h^a \times V_h^a \times \Lambda_h^a$  such that for all  $(\boldsymbol{\tau}_h, \mathbf{v}_h, \boldsymbol{\mu}_h) \in \Sigma_h^a \times V_h^a \times \Lambda_h^a$ ,

$$(\boldsymbol{\sigma}_h^a, \boldsymbol{\tau}_h)_{\omega_a} + (\mathbf{r}_h^a, \nabla \cdot \boldsymbol{\tau}_h)_{\omega_a} + (\boldsymbol{\lambda}_h^a, \boldsymbol{\tau}_h)_{\omega_a} = (\psi_a \boldsymbol{\sigma}(\nabla_s \mathbf{u}_h), \boldsymbol{\tau}_h)_{\omega_a}, \quad (\text{B.7a})$$

$$(\nabla \cdot \boldsymbol{\sigma}_h^a, \mathbf{v}_h)_{\omega_a} = (-\psi_a \mathbf{f} + \boldsymbol{\sigma}(\nabla_s \mathbf{u}_h) \nabla \psi_a, \mathbf{v}_h)_{\omega_a}, \quad (\text{B.7b})$$

$$(\boldsymbol{\sigma}_h^a, \boldsymbol{\mu}_h)_{\omega_a} = 0. \quad (\text{B.7c})$$

Then, extending  $\boldsymbol{\sigma}_h^a$  by zero outside  $\omega_a$ , set  $\boldsymbol{\sigma}_h := \sum_{a \in \mathcal{V}_h} \boldsymbol{\sigma}_h^a$ .

For interior vertices, the source term in (B.7b) has to verify the Neumann compatibility condition

$$(-\psi_a \mathbf{f} + \boldsymbol{\sigma}(\nabla_s \mathbf{u}_h) \nabla \psi_a, \mathbf{z})_{\omega_a} = 0 \quad \forall \mathbf{z} \in \mathbf{RM}^d. \quad (\text{B.8})$$

Taking  $\psi_a \mathbf{z}$  as a test function in (B.4), we see that (B.8) holds and we obtain the following result.

**Lemma B.3** (Properties of  $\boldsymbol{\sigma}_h$ ). *Let  $\boldsymbol{\sigma}_h$  be prescribed by Construction B.2. Then  $\boldsymbol{\sigma}_h \in \mathbb{H}(\text{div}, \Omega)$ , and for all  $T \in \mathcal{T}_h$ , the following holds:*

$$(\mathbf{f} + \nabla \cdot \boldsymbol{\sigma}_h, \mathbf{v})_T = 0 \quad \forall \mathbf{v} \in \mathbf{V}_T \quad \forall T \in \mathcal{T}_h. \quad (\text{B.9})$$

*Proof.* All the fields  $\boldsymbol{\sigma}_h^a$  are in  $\mathbb{H}(\text{div}, \omega_a)$  and satisfy appropriate zero normal conditions so that their zero-extension to  $\Omega$  is in  $\mathbb{H}(\text{div}, \Omega)$ . Hence,  $\boldsymbol{\sigma}_h \in \mathbb{H}(\text{div}, \Omega)$ . Let us prove (B.9). Since (B.8) holds for all  $a \in \mathcal{V}_h^{\text{int}}$ , we infer that (B.7b) is actually true for all  $\mathbf{v}_h \in \mathbf{V}_h(\omega_a)$ . The same holds if  $a \in \mathcal{V}_h^{\text{ext}}$  by definition of  $\mathbf{V}_h^a$ . Since  $\mathbf{V}_h(\omega_a)$  is composed of piecewise polynomials that can be chosen independently in each cell  $T \in \mathcal{T}_a$ , and using  $\boldsymbol{\sigma}_h|_T = \sum_{a \in \mathcal{V}_T} \boldsymbol{\sigma}_h^a|_T$  and the partition of unity  $\sum_{a \in \mathcal{V}_T} \psi_a = 1$ , we infer that  $(\mathbf{f} + \nabla \cdot \boldsymbol{\sigma}_h, \mathbf{v})_T = 0$  for all  $\mathbf{v} \in \mathbf{V}_T$  and all  $T \in \mathcal{T}_h$ .  $\square$

*Remark B.4* (Choice of  $q$ ). In contrast to flux reconstructions in the Raviart–Thomas space (cf. [38] or [83]), the choice  $q = p-1$  in (B.6) does not lead to optimal convergence rates of the resulting error estimate in Section B.4. This is due to the fact that the Brezzi–Douglas–Marini space is smaller than the Raviart–Thomas space of the same degree.

### B.3.2 Discretization and linearization error stress reconstructions

Let now, for  $k \geq 1$ ,  $\mathbf{u}_h^k$  solve (B.5). We will construct two different equilibrated  $H(\text{div})$ -conforming stress tensors. The first one,  $\boldsymbol{\sigma}_{h,\text{disc}}^k$ , represents as above the discrete stress tensor  $\boldsymbol{\sigma}(\nabla_s \mathbf{u}_h^k)$ , for which we will have to modify Construction B.2, because the Neumann compatibility condition (B.8) is not satisfied anymore. The second stress tensor  $\boldsymbol{\sigma}_{h,\text{lin}}^k$  will be a measure for the linearization error and approximate  $\boldsymbol{\sigma}_{\text{lin}}^{k-1}(\nabla_s \mathbf{u}_h^k) - \boldsymbol{\sigma}(\nabla_s \mathbf{u}_h^k)$ . The matrix resulting from the left side of (B.7) will stay unchanged and we will only modify the source terms.

We denote by  $\overline{\boldsymbol{\sigma}}(\nabla_s \mathbf{u}_h^k)$  the  $L^2$ -orthogonal projection of  $\boldsymbol{\sigma}(\nabla_s \mathbf{u}_h^k)$  onto  $[\mathbb{P}^{p-1}(\mathcal{T}_h)]^{d \times d}$  such that  $(\boldsymbol{\sigma}(\nabla_s \mathbf{u}_h^k) - \overline{\boldsymbol{\sigma}}(\nabla_s \mathbf{u}_h^k), \boldsymbol{\tau}_h) = 0$  for any  $\boldsymbol{\tau}_h \in [\mathbb{P}^{p-1}(\mathcal{T}_h)]^{d \times d}$ .

**Construction B.5** (Discrete stress reconstruction). For each  $a \in \mathcal{V}_h$  solve (B.7) with  $\mathbf{u}_h^k$  instead of  $\mathbf{u}_h$ ,  $\overline{\boldsymbol{\sigma}}(\nabla_s \mathbf{u}_h^k)$  instead of  $\boldsymbol{\sigma}(\nabla_s \mathbf{u}_h^k)$  and the source term in (B.7b) replaced by

$$-\psi_a \mathbf{f} + \overline{\boldsymbol{\sigma}}(\nabla_s \mathbf{u}_h^k) \nabla \psi_a - \mathbf{y}_{\text{disc}}^k,$$

where  $\mathbf{y}_{\text{disc}}^k \in \mathbf{RM}^d$  is the unique solution of

$$(\mathbf{y}_{\text{disc}}^k, \mathbf{z})_{\omega_a} = -(\mathbf{f}, \psi_a \mathbf{z})_{\omega_a} + (\overline{\boldsymbol{\sigma}}(\nabla_s \mathbf{u}_h^k), \nabla_s(\psi_a \mathbf{z}))_{\omega_a} \quad \forall \mathbf{z} \in \mathbf{RM}^d. \quad (\text{B.10})$$

The so obtained problem reads: find  $(\boldsymbol{\sigma}_h^a, \mathbf{r}_h^a, \boldsymbol{\lambda}_h^a) \in \boldsymbol{\Sigma}_h^a \times \mathbf{V}_h^a \times \boldsymbol{\Lambda}_h^a$  such that for all  $(\boldsymbol{\tau}_h, \mathbf{v}_h, \boldsymbol{\mu}_h) \in \boldsymbol{\Sigma}_h^a \times \mathbf{V}_h^a \times \boldsymbol{\Lambda}_h^a$ ,

$$\begin{aligned} (\boldsymbol{\sigma}_h^a, \boldsymbol{\tau}_h)_{\omega_a} + (\mathbf{r}_h^a, \nabla \cdot \boldsymbol{\tau}_h)_{\omega_a} + (\boldsymbol{\lambda}_h^a, \boldsymbol{\tau}_h)_{\omega_a} &= (\psi_a \bar{\boldsymbol{\sigma}}(\nabla_s \mathbf{u}_h^k), \boldsymbol{\tau}_h)_{\omega_a}, \\ (\nabla \cdot \boldsymbol{\sigma}_h^a, \mathbf{v}_h)_{\omega_a} &= (-\psi_a \mathbf{f} + \bar{\boldsymbol{\sigma}}(\nabla_s \mathbf{u}_h^k) \nabla \psi_a - \mathbf{y}_{\text{disc}}^k, \mathbf{v}_h)_{\omega_a}, \\ (\boldsymbol{\sigma}_h^a, \boldsymbol{\mu}_h)_{\omega_a} &= 0. \end{aligned}$$

Then set  $\boldsymbol{\sigma}_{h,\text{disc}}^k := \sum_{a \in \mathcal{V}_h} \boldsymbol{\sigma}_h^a$ .

**Construction B.6** (Linearization error stress reconstruction). For each  $a \in \mathcal{V}_h$  solve (B.7) with  $\mathbf{u}_h^k$  instead of  $\mathbf{u}_h$ , the source term in (B.7a) replaced by

$$\psi_a (\boldsymbol{\sigma}_{\text{lin}}^{k-1}(\nabla_s \mathbf{u}_h^k) - \bar{\boldsymbol{\sigma}}(\nabla_s \mathbf{u}_h^k)),$$

and the source term in (B.7b) replaced by

$$(\boldsymbol{\sigma}_{\text{lin}}^{k-1}(\nabla_s \mathbf{u}_h^k) - \bar{\boldsymbol{\sigma}}(\nabla_s \mathbf{u}_h^k)) \nabla \psi_a + \mathbf{y}_{\text{disc}}^k,$$

where  $\mathbf{y}_{\text{disc}}^k \in \mathbf{RM}^d$  is defined by (B.10). The corresponding local problem is to find  $(\boldsymbol{\sigma}_h^a, \mathbf{r}_h^a, \boldsymbol{\lambda}_h^a) \in \boldsymbol{\Sigma}_h^a \times \mathbf{V}_h^a \times \boldsymbol{\Lambda}_h^a$  such that for all  $(\boldsymbol{\tau}_h, \mathbf{v}_h, \boldsymbol{\mu}_h) \in \boldsymbol{\Sigma}_h^a \times \mathbf{V}_h^a \times \boldsymbol{\Lambda}_h^a$ ,

$$\begin{aligned} (\boldsymbol{\sigma}_h^a, \boldsymbol{\tau}_h)_{\omega_a} + (\mathbf{r}_h^a, \nabla \cdot \boldsymbol{\tau}_h)_{\omega_a} + (\boldsymbol{\lambda}_h^a, \boldsymbol{\tau}_h)_{\omega_a} &= (\psi_a (\boldsymbol{\sigma}_{\text{lin}}^{k-1}(\nabla_s \mathbf{u}_h^k) - \bar{\boldsymbol{\sigma}}(\nabla_s \mathbf{u}_h^k)), \boldsymbol{\tau}_h)_{\omega_a}, \\ (\nabla \cdot \boldsymbol{\sigma}_h^a, \mathbf{v}_h)_{\omega_a} &= ((\boldsymbol{\sigma}_{\text{lin}}^{k-1}(\nabla_s \mathbf{u}_h^k) - \bar{\boldsymbol{\sigma}}(\nabla_s \mathbf{u}_h^k)) \nabla \psi_a + \mathbf{y}_{\text{disc}}^k, \mathbf{v}_h)_{\omega_a}, \\ (\boldsymbol{\sigma}_h^a, \boldsymbol{\mu}_h)_{\omega_a} &= 0. \end{aligned}$$

Then set  $\boldsymbol{\sigma}_{h,\text{lin}}^k := \sum_{a \in \mathcal{V}_h} \boldsymbol{\sigma}_h^a$ .

Notice that the role of  $\mathbf{y}_{\text{disc}}^k$  is to guarantee that, for interior vertices, the source terms in Constructions B.5 and B.6 satisfy the Neumann compatibility conditions

$$\begin{aligned} (-\psi_a \mathbf{f} + \bar{\boldsymbol{\sigma}}(\nabla_s \mathbf{u}_h^k) \nabla \psi_a - \mathbf{y}_{\text{disc}}^k, \mathbf{z})_{\omega_a} &= 0 \quad \forall \mathbf{z} \in \mathbf{RM}^d, \\ ((\boldsymbol{\sigma}_{\text{lin}}^{k-1}(\nabla_s \mathbf{u}_h^k) - \bar{\boldsymbol{\sigma}}(\nabla_s \mathbf{u}_h^k)) \nabla \psi_a + \mathbf{y}_{\text{disc}}^k, \mathbf{z})_{\omega_a} &= 0 \quad \forall \mathbf{z} \in \mathbf{RM}^d. \end{aligned}$$

**Lemma B.7** (Properties of the discretization and linearization error stress reconstructions). Let  $\boldsymbol{\sigma}_{h,\text{disc}}^k$  and  $\boldsymbol{\sigma}_{h,\text{lin}}^k$  be prescribed by Constructions B.5 and B.6. Then it holds

1.  $\boldsymbol{\sigma}_{h,\text{disc}}^k, \boldsymbol{\sigma}_{h,\text{lin}}^k \in \mathbb{H}(\text{div}, \Omega)$ ,
2.  $(\mathbf{f} + \nabla \cdot (\boldsymbol{\sigma}_{h,\text{disc}}^k + \boldsymbol{\sigma}_{h,\text{lin}}^k), \mathbf{v})_T = 0 \quad \forall \mathbf{v} \in \mathbf{V}_T \quad \forall T \in \mathcal{T}_h$ ,
3. As the Newton solver converges,  $\boldsymbol{\sigma}_{h,\text{lin}}^k \rightarrow \mathbf{0}$ .

*Proof.* The proof is similar to the proof of Lemma B.3. The first property is again satisfied due to the definition of  $\boldsymbol{\Sigma}_h^a$ . In order to show that the second property holds, we add the two equations (B.7b) obtained for each of the constructions. The right hand side of this sum then reads  $(-\psi_a \mathbf{f} + \boldsymbol{\sigma}_{\text{lin}}^{k-1}(\nabla_s \mathbf{u}_h^k) \nabla \psi_a, \mathbf{v}_h)_{\omega_a}$ . Once again we can, for any  $\mathbf{z} \in \mathbf{RM}^d$ , take  $\psi_a \mathbf{z}$  as a test function in (B.5) to show that this term is zero if  $\mathbf{v}_h \in \mathbf{RM}^d$ , and so the equation holds for all  $\mathbf{v}_h \in \mathbf{V}_h(\omega_a)$ . Then we proceed as in the proof of Lemma B.3.  $\square$

## B.4 A posteriori error estimate and adaptive algorithm

In this section we first derive an upper bound on the error between the analytical solution of (B.3) and the solution  $\mathbf{u}_h$  of (B.4), in which we then identify and distinguish the discretization and linearization error components at each Newton iteration for the solution  $\mathbf{u}_h^k$  of (B.5). Based on this distinction, we present an adaptive algorithm stopping the Newton iterations once the linearization error estimate is dominated by the estimate of the discretization error. Finally, we show the efficiency of the error estimate.

### B.4.1 Guaranteed upper bound

We measure the error in the energy norm  $\|\mathbf{v}\|_{\text{en}}^2 := a(\mathbf{v}, \mathbf{v}) = (\boldsymbol{\sigma}(\nabla_s \mathbf{v}), \nabla_s \mathbf{v})$ , for which we obtain the properties

$$\sigma_*^2 C_K^{-2} \|\nabla \mathbf{v}\|^2 \leq \|\mathbf{v}\|_{\text{en}}^2 \leq \bar{\sigma} \|\nabla_s \mathbf{v}\|^2, \quad (\text{B.13})$$

by applying (3.40b) and the Korn inequality for the left inequality, and the Cauchy–Schwarz inequality and (3.2c) for the right one. In our case it holds  $C_K = \sqrt{2}$ , owing to (B.1b).

**Theorem B.8** (Basic a posteriori error estimate). *Let  $\mathbf{u}$  be the solution of (B.3) and  $\mathbf{u}_h$  the solution of (B.4). Let  $\boldsymbol{\sigma}_h$  be the stress tensor defined in Construction B.2. Then,*

$$\|\mathbf{u} - \mathbf{u}_h\|_{\text{en}} \leq \sqrt{2} \bar{\sigma} \sigma_*^{-3} \left( \sum_{T \in \mathcal{T}_h} \left( \frac{h_T}{\pi} \|\mathbf{f} + \nabla \cdot \boldsymbol{\sigma}_h\|_T + \|\boldsymbol{\sigma}_h - \boldsymbol{\sigma}(\nabla_s \mathbf{u}_h)\|_T \right)^2 \right)^{1/2}. \quad (\text{B.14})$$

*Remark B.9* (Constants  $\bar{\sigma}$  and  $\sigma_*$ ). For the estimate to be computable, the constants  $\bar{\sigma}$  and  $\sigma_*$  have to be specified. For the linear elasticity model (3.3) we set  $\bar{\sigma} := 2\mu + d\lambda$  and  $\sigma_* := \sqrt{2\mu}$ , whereas for the Hencky–Mises model (3.4) we set  $\bar{\sigma} := 2\tilde{\mu}(0) + d\tilde{\lambda}(0)$  and  $\sigma_* := \sqrt{2\tilde{\mu}(0)}$ . Following [175], we obtain a sharper bound in the case of linear elasticity, with  $\mu^{-1/2}$  instead of  $\sqrt{2}\bar{\sigma}\sigma_*^{-3}$  in (B.14).

*Proof of Theorem B.8.* We start by bounding the energy norm of the error by the dual norm of the residual of the weak formulation (B.3). Using (B.13), (3.40b), the linearity of  $a$  in its second argument, and (B.3) we obtain

$$\begin{aligned} \|\mathbf{u} - \mathbf{u}_h\|_{\text{en}}^2 &\leq \bar{\sigma} \|\nabla_s(\mathbf{u} - \mathbf{u}_h)\|^2 \leq \bar{\sigma} \sigma_*^{-2} |a(\mathbf{u}, \mathbf{u} - \mathbf{u}_h) - a(\mathbf{u}_h, \mathbf{u} - \mathbf{u}_h)| \\ &= \bar{\sigma} \sigma_*^{-2} \|\nabla(\mathbf{u} - \mathbf{u}_h)\| \left| a\left(\mathbf{u}, \frac{\mathbf{u} - \mathbf{u}_h}{\|\nabla(\mathbf{u} - \mathbf{u}_h)\|}\right) - a\left(\mathbf{u}_h, \frac{\mathbf{u} - \mathbf{u}_h}{\|\nabla(\mathbf{u} - \mathbf{u}_h)\|}\right) \right| \\ &\leq \bar{\sigma} \sigma_*^{-3} C_K \|\mathbf{u} - \mathbf{u}_h\|_{\text{en}} \sup_{\mathbf{v} \in [H_0^1(\Omega)]^d, \|\nabla \mathbf{v}\|=1} \{a(\mathbf{u}, \mathbf{v}) - a(\mathbf{u}_h, \mathbf{v})\} \\ &= \bar{\sigma} \sigma_*^{-3} C_K \|\mathbf{u} - \mathbf{u}_h\|_{\text{en}} \sup_{\mathbf{v} \in [H_0^1(\Omega)]^d, \|\nabla \mathbf{v}\|=1} \{(\mathbf{f}, \mathbf{v}) - (\boldsymbol{\sigma}(\nabla_s \mathbf{u}_h), \nabla_s \mathbf{v})\}. \end{aligned}$$

and thus

$$\|\mathbf{u} - \mathbf{u}_h\|_{\text{en}} \leq \bar{\sigma} \sigma_*^{-3} C_K \sup_{\mathbf{v} \in [H_0^1(\Omega)]^d, \|\nabla \mathbf{v}\|=1} \{(\mathbf{f}, \mathbf{v}) - (\boldsymbol{\sigma}(\nabla_s \mathbf{u}_h), \nabla_s \mathbf{v})\}. \quad (\text{B.15})$$

Note that, due to the symmetry of  $\sigma$  we can replace  $\nabla_s \mathbf{v}$  by  $\nabla \mathbf{v}$  in the second term inside the supremum. Now fix  $\mathbf{v} \in [H_0^1(\Omega)]^d$ , such that  $\|\nabla \mathbf{v}\| = 1$ . Since  $\sigma_h \in \mathbb{H}(\text{div}, \Omega)$ , we can insert  $(\nabla \cdot \sigma_h, \mathbf{v}) + (\sigma_h, \nabla \mathbf{v}) = 0$  into the term inside the supremum and obtain

$$(\mathbf{f}, \mathbf{v}) - (\sigma(\nabla_s \mathbf{u}_h), \nabla \mathbf{v}) = (\mathbf{f} + \nabla \cdot \sigma_h, \mathbf{v}) + (\sigma_h - \sigma(\nabla_s \mathbf{u}_h), \nabla \mathbf{v}). \quad (\text{B.16})$$

For the first term of the right hand side of (B.16) we obtain, using (B.9) on each  $T \in \mathcal{T}_h$  to insert  $\Pi_T^0 \mathbf{v}$ , which denotes the  $L^2$ -projection of  $\mathbf{v}$  onto  $[\mathbb{P}^0(T)]^d$ , the Cauchy–Schwarz inequality and the Poincaré inequality on simplices,

$$|(\mathbf{f} + \nabla \cdot \sigma_h, \mathbf{v})| \leq \left| \sum_{T \in \mathcal{T}_h} (\mathbf{f} + \nabla \cdot \sigma_h, \mathbf{v} - \Pi_T^0 \mathbf{v})_T \right| \leq \sum_{T \in \mathcal{T}_h} \frac{h_T}{\pi} \|\mathbf{f} + \nabla \cdot \sigma_h\|_T \|\nabla \mathbf{v}\|_T, \quad (\text{B.17})$$

whereas the Cauchy–Schwarz inequality applied to the second term directly yields

$$|(\sigma_h - \sigma(\nabla_s \mathbf{u}_h), \nabla \mathbf{v})| \leq \sum_{T \in \mathcal{T}_h} \|\sigma_h - \sigma(\nabla_s \mathbf{u}_h)\|_T \|\nabla \mathbf{v}\|_T.$$

Inserting in (B.15) and applying again the Cauchy–Schwarz inequality yields the result.  $\square$

The goal of the following result is to elaborate the error estimate (B.14) so as to distinguish different error components using the equilibrated stress tensors of Constructions B.5 and B.6. This distinction is essential for the development of Algorithm B.12, where the mesh and the stopping criteria for the iterative solver are chosen adaptively. Notice that in Theorem B.8 we don't necessarily need  $\sigma_h$  to be the stress tensor obtained in Construction B.2. We only need it to satisfy two properties: First, equation (B.16) requires  $\sigma_h$  to lie in  $\mathbb{H}(\text{div}, \Omega)$ . Second, in order to be able to apply the Poincaré inequality in (B.17),  $\sigma_h$  has to satisfy the local equilibrium relation

$$(\mathbf{f} - \nabla \cdot \sigma_h, \mathbf{v})_T = 0 \quad \forall \mathbf{v} \in [\mathbb{P}^0(T)]^d \quad \forall T \in \mathcal{T}_h.$$

Thus, the theorem also holds for  $\sigma_h := \sigma_{h,\text{disc}}^k + \sigma_{h,\text{lin}}^k$ , where  $\sigma_{h,\text{disc}}^k$  and  $\sigma_{h,\text{lin}}^k$  are defined in Constructions B.5.

**Theorem B.10** (A posteriori error estimate distinguishing different error sources). *Let  $\mathbf{u}$  be the solution of (B.3),  $\mathbf{u}_h^k$  the solution of (B.5), and  $\sigma_h := \sigma_{h,\text{disc}}^k + \sigma_{h,\text{lin}}^k$ . Then,*

$$\|\mathbf{u} - \mathbf{u}_h^k\|_{\text{en}} \leq \sqrt{2\bar{\sigma}} \sigma_*^{-3} (\eta_{\text{disc}}^k + \eta_{\text{lin}}^k + \eta_{\text{quad}}^k + \eta_{\text{osc}}^k),$$

where the local discretization, linearization, quadrature and oscillation error estimators on each  $T \in \mathcal{T}_h$  are defined as

$$\begin{aligned} \eta_{\text{disc},T}^k &:= \|\sigma_{h,\text{disc}}^k - \bar{\sigma}(\nabla_s \mathbf{u}_h^k)\|_T, \\ \eta_{\text{lin},T}^k &:= \|\sigma_{h,\text{lin}}^k\|_T, \\ \eta_{\text{quad},T}^k &:= \|\bar{\sigma}(\nabla_s \mathbf{u}_h^k) - \sigma(\nabla_s \mathbf{u}_h^k)\|_T, \\ \eta_{\text{osc},T}^k &:= h_T \pi^{-1} \|\mathbf{f} - \Pi_T^{p-1} \mathbf{f}\|_T, \end{aligned}$$

with  $\Pi_T^{p-1}$  denoting the  $L^2$ -projection onto  $[\mathbb{P}^{p-1}(T)]^d$ , and for each error source the global estimator is given by  $\eta^k := 2 \left( \sum_{T \in \mathcal{T}_h} (\eta_{\cdot,T}^k)^2 \right)^{1/2}$ .

*Remark B.11* (Quadrature error). In practice, the projection  $\bar{\sigma}(\nabla_s \mathbf{u}_h^k)$  of  $\sigma(\nabla_s \mathbf{u}_h^k)$  onto  $[\mathbb{P}^{p-1}(\mathcal{T}_h)]^{d \times d}$  for a general nonlinear stress-strain relation cannot be computed exactly. The quadrature error estimator  $\eta_{\text{quad},T}^k$  measures the quality of this projection and is computed by using a considerably higher number of Gauss points.

## B.4.2 Adaptive algorithm

Based on the error estimate of Theorem B.10, we propose an adaptive algorithm where the mesh size is locally adapted, and a dynamic stopping criterion is used for the linearization iterations. The idea is to compare the estimators for the different error sources with each other in order to concentrate the computational effort on reducing the dominant one. For this purpose, let  $\gamma_{\text{lin}}, \gamma_{\text{quad}} \in (0, 1)$  be user-given weights representing the relative size of the quadrature and linearization errors with respect to the total error.

**Algorithm B.12** (Mesh adaptation algorithm).

1. Choose an initial function  $\mathbf{u}_h^0 \in [H_0^1(\Omega)]^d \cap [\mathbb{P}^p(\mathcal{T}_h)]^d$  and set  $k := 1$
2. Set the initial quadrature precision  $\nu := 2p$  (exactness for polynomials up to degree  $\nu$ )
3. Linearization iterations
  - (a) Calculate  $\sigma_{\text{lin}}^{k-1}(\nabla_s \mathbf{u}_h^k)$ ,  $\mathbf{u}_h^k$ , and  $\sigma(\nabla_s \mathbf{u}_h^k)$
  - (b) Calculate  $\bar{\sigma}(\nabla_s \mathbf{u}_h^k)$ , the stress reconstructions  $\sigma_{h,\text{disc}}^k$  and  $\sigma_{h,\text{lin}}^k$ , and the error estimators  $\eta_{\text{disc}}^k$ ,  $\eta_{\text{lin}}^k$ ,  $\eta_{\text{osc}}^k$  and  $\eta_{\text{quad}}^k$
  - (c) Improve the quadrature rule (setting  $\nu := \nu + 1$ ) and go back to step 3(b) until
$$\eta_{\text{quad}}^k \leq \gamma_{\text{quad}}(\eta_{\text{disc}}^k + \eta_{\text{lin}}^k + \eta_{\text{osc}}^k) \quad (\text{B.18a})$$
  - (d) Set  $k := k + 1$  and go back to step 3(a) until
$$\eta_{\text{lin}}^k \leq \gamma_{\text{lin}}(\eta_{\text{disc}}^k + \eta_{\text{osc}}^k) \quad (\text{B.18b})$$
4. Refine or coarsen the mesh  $\mathcal{T}_h$  such that the local discretization error estimators  $\eta_{\text{disc},T}^k$  are distributed evenly

Instead of using the global stopping criteria (B.18a) and (B.18b), which are evaluated over all mesh elements, we can also define the local criteria

$$\begin{aligned} \eta_{\text{quad},T}^k &\leq \gamma_{\text{quad}}(\eta_{\text{disc},T}^k + \eta_{\text{lin},T}^k + \eta_{\text{osc},T}^k) \quad \forall T \in \mathcal{T}_h, \\ \eta_{\text{lin},T}^k &\leq \gamma_{\text{lin}}(\eta_{\text{disc},T}^k + \eta_{\text{osc},T}^k) \quad \forall T \in \mathcal{T}_h, \end{aligned} \quad (\text{B.19})$$

where it is also possible to define local weights  $\gamma_{\text{lin},T}$  and  $\gamma_{\text{quad},T}$  for each element. These local stopping criteria are necessary to establish the local efficiency of the error estimators in the following section, whereas the global criteria are only sufficient to prove global efficiency.

### B.4.3 Local efficiency

Let us start by introducing some additional notation used in this section. For a given element  $T \in \mathcal{T}_h$ , the set  $\mathcal{T}_T$  collects the elements sharing at least a vertex with  $T$ . The set  $\mathcal{F}_h$  contains all faces (if  $d = 2$  we will, for simplicity, refer to the edges as faces) of the mesh and is decomposed into boundary faces  $\mathcal{F}_h^{\text{ext}}$  and interfaces  $\mathcal{F}_h^{\text{int}}$ . For any  $T \in \mathcal{T}_h$  the set  $\mathcal{F}_T$  contains the faces of  $T$ , whereas  $\mathcal{F}_{\mathcal{T}_T}$  collects all faces sharing at least a vertex with  $T$  and we denote  $\mathcal{F}_{\mathcal{T}_T}^{\text{int}} = \mathcal{F}_{\mathcal{T}_T} \cap \mathcal{F}_h^{\text{int}}$ . We also assume that  $\mathcal{T}_h$  verifies the minimum angle condition, i.e., there exists  $\alpha_{\min} > 0$  uniform with respect to all considered meshes such that the minimum angle  $\alpha_T$  of each  $T \in \mathcal{T}_h$  satisfies  $\alpha_T \geq \alpha_{\min}$ . In what follows we let  $a \lesssim b$  stand for  $a \leq Cb$  with a generic constant  $C$ , which is independent of the mesh size, the domain  $\Omega$  and the stress-strain relation, but that can depend on  $\alpha_{\min}$  and on the polynomial degree  $p$ .

To prove efficiency, we will use a posteriori error estimators of residual type. Following [194, 195] we define for  $X \subseteq \Omega$  the functional  $\mathcal{R}_X : [H^1(X)]^d \rightarrow [H^{-1}(X)]^d$  such that, for all  $\mathbf{v} \in [H^1(X)]^d, \mathbf{w} \in [H_0^1(X)]^d$ ,

$$\langle \mathcal{R}_X(\mathbf{v}), \mathbf{w} \rangle_X := (\boldsymbol{\sigma}(\nabla_s \mathbf{v}), \nabla_s \mathbf{w})_X - (\mathbf{f}, \mathbf{w})_X.$$

Define, for each  $T \in \mathcal{T}_h$ ,

$$\begin{aligned} (\eta_{\sharp, T}^k)^2 &:= \sum_{T' \in \mathcal{T}_T} h_{T'}^2 \|\nabla \cdot \bar{\boldsymbol{\sigma}}(\nabla_s \mathbf{u}_h^k) + \mathbf{\Pi}_T^p \mathbf{f}\|_{T'}^2 + \sum_{F \in \mathcal{F}_{\mathcal{T}_T}^{\text{int}}} h_F \|\llbracket \bar{\boldsymbol{\sigma}}(\nabla_s \mathbf{u}_h^k) \mathbf{n}_F \rrbracket_F\|_F^2, \\ (\eta_{b, T}^k)^2 &:= \sum_{T' \in \mathcal{T}_T} h_{T'}^2 \|\nabla \cdot (\boldsymbol{\sigma}(\nabla_s \mathbf{u}_h^k) - \bar{\boldsymbol{\sigma}}(\nabla_s \mathbf{u}_h^k))\|_{T'}^2 + \sum_{F \in \mathcal{F}_{\mathcal{T}_T}^{\text{int}}} h_F \|\llbracket (\boldsymbol{\sigma}(\nabla_s \mathbf{u}_h^k) - \bar{\boldsymbol{\sigma}}(\nabla_s \mathbf{u}_h^k)) \mathbf{n}_F \rrbracket_F\|_F^2, \end{aligned}$$

with  $[\cdot]_F$  denoting the jump operators defined in (2.18). The quantity  $\eta_{b, T}^k$  measures the quality of the approximation of  $\boldsymbol{\sigma}(\nabla_s \mathbf{u}_h^k)$  by  $\bar{\boldsymbol{\sigma}}(\nabla_s \mathbf{u}_h^k)$  and can be estimated explicitly. The following result is shown in [194, Section 3.3]. Denoting by  $\eta_{\text{osc}, \mathcal{T}_T}^k := \{2 \sum_{T' \in \mathcal{T}_T} (\eta_{\text{osc}, T'}^k)^2\}^{1/2}$ , it holds

$$\eta_{\sharp, T}^k \lesssim \|\mathcal{R}_{\mathcal{T}_T}(\mathbf{u}_h^k)\|_{[H^{-1}(\mathcal{T}_T)]^d} + \eta_{b, T}^k + \eta_{\text{osc}, \mathcal{T}_T}^k. \quad (\text{B.20})$$

In order to bound the dual norm of the residual, we need to additionally assume that the stress-strain relation  $\boldsymbol{\sigma}$  is Lipschitz continuous, cf. Assumption 3.14. We recall that the three stress-strain relations presented in Examples 3.2, 3.3, and 3.4 satisfy the previous Lipschitz continuity assumptions. Owing to the definition of the functional  $\mathcal{R}_{\mathcal{T}_T}$  and to the fact that  $-\nabla \cdot \boldsymbol{\sigma}(\nabla_s \mathbf{u}) = \mathbf{f} \in L^2(\mathcal{T}_T)^d$ , using the Cauchy–Schwarz inequality and the Lipschitz continuity (3.40a) of  $\boldsymbol{\sigma}$ , it is inferred that

$$\begin{aligned} \|\mathcal{R}_{\mathcal{T}_T}(\mathbf{u}_h^k)\|_{[H^{-1}(\mathcal{T}_T)]^d} &:= \sup_{\mathbf{w} \in [H_0^1(\mathcal{T}_T)]^d, \|\mathbf{w}\|_{[H_0^1(\mathcal{T}_T)]^d} \leq 1} (\boldsymbol{\sigma}(\nabla_s \mathbf{u}_h^k), \nabla_s \mathbf{w})_{\mathcal{T}_T} - (\mathbf{f}, \mathbf{w})_{\mathcal{T}_T} \\ &= \sup_{\mathbf{w} \in [H_0^1(\mathcal{T}_T)]^d, \|\mathbf{w}\|_{[H_0^1(\mathcal{T}_T)]^d} \leq 1} (\boldsymbol{\sigma}(\nabla_s \mathbf{u}_h^k) - \boldsymbol{\sigma}(\nabla_s \mathbf{u}), \nabla_s \mathbf{w})_{\mathcal{T}_T} \\ &\leq \sigma^* \|\nabla_s(\mathbf{u} - \mathbf{u}_h^k)\|_{\mathcal{T}_T}. \end{aligned}$$



Thus, by (B.20), the previous bound, and the strong monotonicity (3.40b) it holds

$$\eta_{\sharp,T}^k \lesssim \sigma^* \sigma_*^{-1} \|\mathbf{u} - \mathbf{u}_h^k\|_{\text{en},\mathcal{T}_T} + \eta_{b,T}^k + \eta_{\text{osc},\mathcal{T}_T}^k. \quad (\text{B.21})$$

It is well known that there exist nonconforming finite element methods which are equivalent to mixed finite element methods using the Brezzi–Douglas–Marini spaces (see e.g. [5]). Following the ideas of [91, 95, 118] and references therein, we use these spaces to prove that  $\eta_{\text{disc},T}^k \lesssim \eta_{\sharp,T}^k$ . As a result of this bound together with (B.21), we obtain the following Theorem (see [36, Section 4.4] for the detailed proof).

**Theorem B.13** (Local efficiency). *Let  $\mathbf{u} \in [H_0^1(\Omega)]^d$  be the solution of (B.3),  $\mathbf{u}_h^k \in [H_0^1(\Omega)]^d \cap [\mathbb{P}^p(\mathcal{T}_h)]^d$  be arbitrary and  $\sigma_{h,\text{disc}}^k$  and  $\sigma_{h,\text{lin}}^k$  defined by Constructions B.5 and B.6. Let the local stopping criteria (B.19) be verified. Then it holds for all  $T \in \mathcal{T}_h$ ,*

$$\eta_{\text{disc},T}^k + \eta_{\text{lin},T}^k + \eta_{\text{quad},T}^k + \eta_{\text{osc},T}^k \lesssim \sigma^* \sigma_*^{-1} \|\mathbf{u} - \mathbf{u}_h^k\|_{\text{en},\mathcal{T}_T} + \eta_{b,T}^k + \eta_{\text{osc},\mathcal{T}_T}^k.$$

## B.5 Numerical results

In this section we illustrate numerically our results on two test cases, both performed with the Code\_Aster<sup>1</sup> software, which uses conforming finite elements of degree  $p = 2$ . Our intention is, first, to show the relevance of the discretization error estimators used as mesh refinement indicators, and second, to propose a stopping criterion for the Newton iterations based on the linearization error estimator. All the triangulations are conforming, since in the remeshing progress hanging nodes are removed by bisecting the neighboring element.

### B.5.1 L-shaped domain

Following [4, 129, 158], we consider the L-shaped domain  $\Omega = (-1, 1)^2 \setminus ([0, 1] \times [-1, 0])$ , where for the linear elasticity case an analytical solution is given by

$$\mathbf{u}(r, \theta) = \frac{1}{2\mu} r^\alpha \begin{pmatrix} \cos(\alpha\theta) - \cos((\alpha - 2)\theta) \\ A \sin(\alpha\theta) + \sin((\alpha - 2)\theta) \end{pmatrix},$$

with the parameters

$$\mu = 1.0, \quad \lambda = 5.0, \quad \alpha = 0.6, \quad A = 31/9.$$

This solution is imposed as Dirichlet boundary condition on  $\partial\Omega$ , together with  $\mathbf{f} = \mathbf{0}$  in  $\Omega$ . We perform this test for two different stress-strain relations. First on the linear elasticity model (3.3), where we can compare the error estimate (B.14) to the analytical error  $\|\mathbf{u} - \mathbf{u}_h\|_{\text{en}}$ . The second relation is the nonlinear Hencky–Mises model (3.4), for which we distinguish the discretization and linearization error components and use the adaptive algorithm from Section B.4.2.

We compute the analytical error and its estimate on two series of unstructured meshes. Starting with the same initial mesh, we use uniform mesh refinement for the first one and

<sup>1</sup><http://web-code-aster.org>

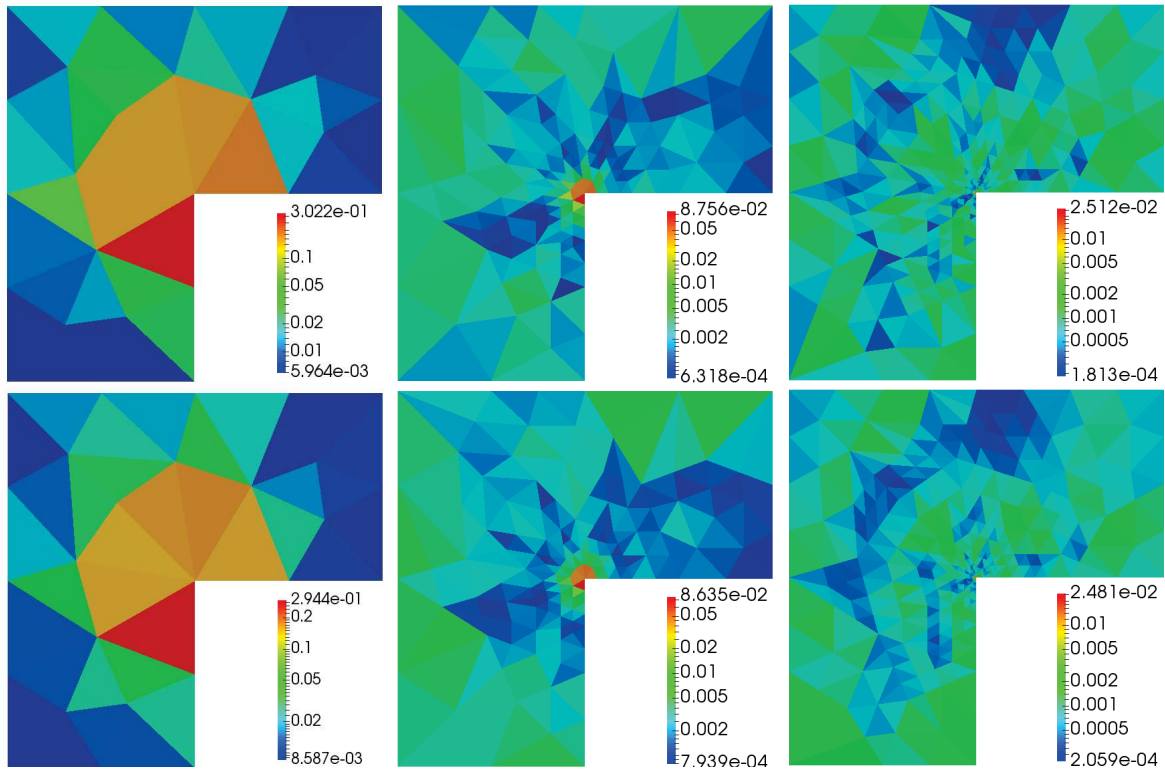


Figure B.2: L-shaped domain with linear elasticity model. Distribution of the error estimators (top) and the analytical error (bottom) for the initial mesh (left) and after three (middle) and six (right) adaptive mesh refinements.

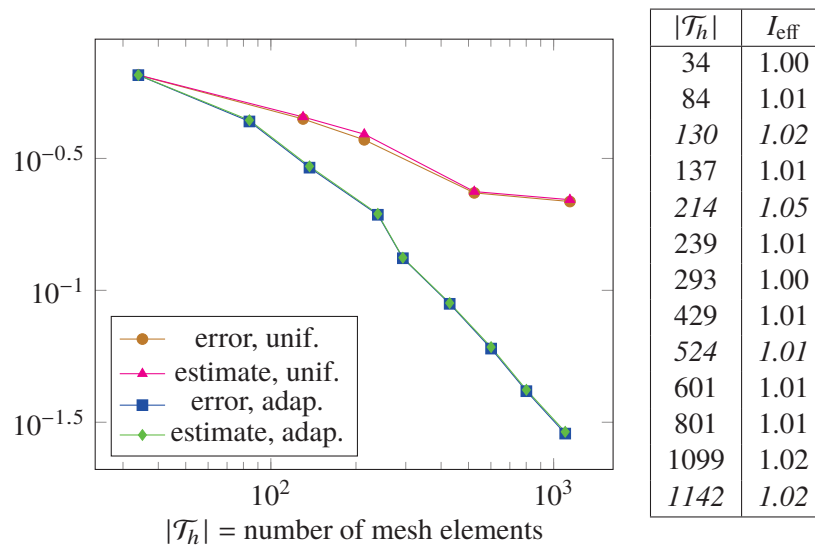


Figure B.3: L-shaped domain with linear elasticity model. *Left*: Comparison of the error estimate (B.14) and  $\|\mathbf{u} - \mathbf{u}_h\|_{\text{en}}$  on two series of meshes, obtained by uniform and adaptive remeshing. *Right*: Effectivity indices of the estimate for each mesh, with the meshes stemming from uniform refinement in italic.

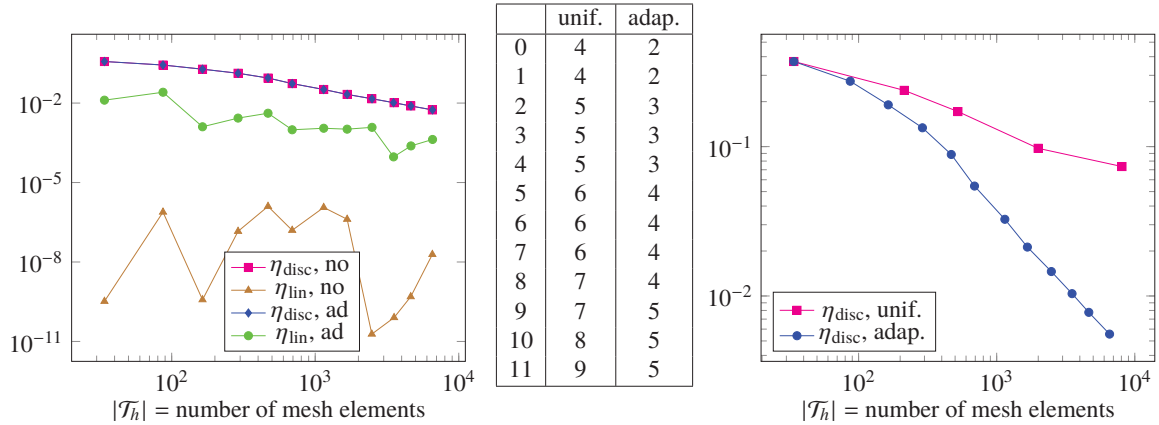


Figure B.4: L-shaped domain with Hencky–Mises model. *Left*: Comparison of the discretization and linearization error estimators on a series of meshes, without and with adaptive stopping criterion for the Newton algorithm. *Middle*: Number of Newton iterations without and with adaptive stopping criterion for each mesh. *Right*: Discretization error estimate for uniform and adaptive remeshing.

adaptive refinement based on the error estimate for the second series. Figure B.2 compares the distribution of the error and the estimators on the initial and two adaptively refined meshes. The error estimators reflect the distribution of the analytical error, which makes them a good indicator for adaptive remeshing. Figure B.3 shows the global estimates and errors for each mesh, as well as their effectivity index corresponding to the ratio of the estimate to the error. We obtain effectivity indices close to one, showing that the estimated error value lies close to the actual one, what we can also observe in the graphics on the left. As expected, the adaptively refined mesh series has a higher convergence rate, with corresponding error an order of magnitude lower for  $10^3$  elements.

For the Hencky–Mises model we choose the Lamé functions

$$\tilde{\mu}(\rho) := a + b(1 + \rho^2)^{-1/2}, \quad \tilde{\lambda}(\rho) := \kappa - \frac{3}{2}\tilde{\mu}(\rho),$$

corresponding to the Carreau law for elastoplastic materials (see, e.g. [104, 143, 177]), and we set  $a = 1/20$ ,  $b = 1/2$ , and  $\kappa = 17/3$  so that the shear modulus reduces progressively to approximately 10% of its initial value. This model allows us to soften the singularity observed in the linear case and to validate our results on more homogeneous error distributions. We apply Algorithm B.12 with  $\gamma_{\text{lin}} = 0.1$  and compare the obtained results to those without the adaptive stopping criterion for the Newton solver. In both cases, we use adaptive remeshing based on the spatial error estimators.

The results are shown in Figure B.4. In the left graphic we observe that the linearization error estimate in the adaptive case is much higher than in the one without adaptive stopping criterion. We see that this does not affect the discretization error estimator. The table shows the number of performed Newton iterations for both cases. The algorithm using the adaptive stopping criterion is more efficient. In the right graphic we compare the discretization error estimate on two series of meshes, one refined uniformly and the other one adaptively, based on the local discretization error estimators. Again the convergence rate is higher for the adaptively refined mesh series.

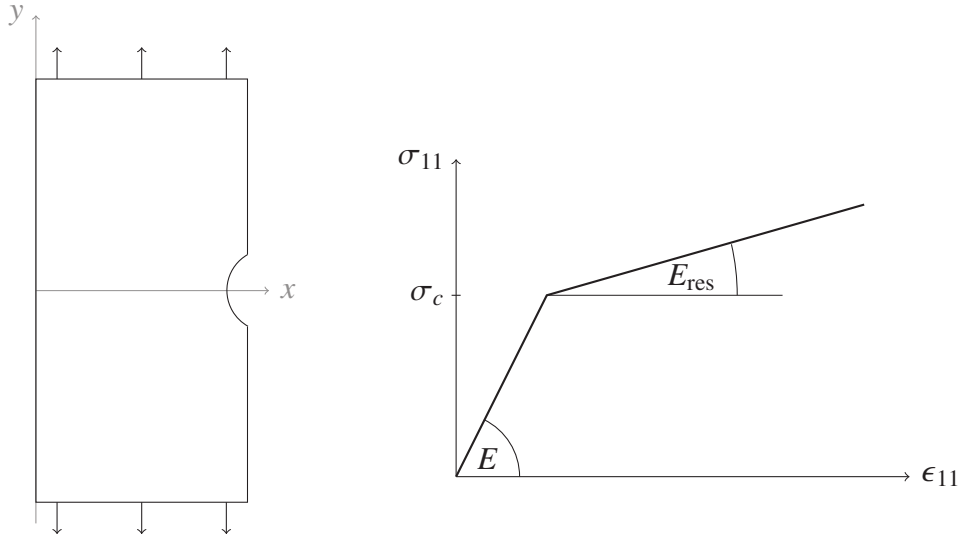


Figure B.5: *Left*: Notched specimen plate. *Right*: Uniaxial traction curve

### B.5.2 Notched specimen plate

In our second test we use the two nonlinear models of Examples 3.3 and 3.4 on a more application-oriented test. The idea is to set a special sample geometry yielding to a model discrimination test, namely different physical results for different models. We simulate the uniform traction of a notched specimen under plain strain assumption (cf. Figure B.5). The notch is meant to favor strain localization phenomenon. We consider a domain  $\Omega = (0, 10\text{m}) \times (-10\text{m}, 10\text{m}) \setminus \{\mathbf{x} \in \mathbb{R}^2 \mid \|\mathbf{x}\text{m} - (0, 11\text{m})^T\| \leq 2\text{m}\}$ , we take  $\mathbf{f} = \mathbf{0}$ , and we prescribe a displacement on the boundary leading to the following Dirichlet conditions:

$$u_x = 0\text{m} \text{ if } x = 0\text{m}, \quad u_y = -1.1 \cdot 10^{-3}\text{m} \text{ if } y = -10\text{m}, \quad u_y = 1.1 \cdot 10^{-3}\text{m} \text{ if } y = 10\text{m}.$$

In many applications, the information about the material properties are obtained in uniaxial experiments, yielding a relation between  $\sigma_{ii}$  and  $\epsilon_{ii}$  for a space direction  $x_i$ . Since we only consider isotropic materials, we can choose  $i = 1$ . From this curve one can compute the nonlinear Lamé functions of (3.4) and the damage function in (3.5). Although the uniaxial relation is the same, the resulting stress-strain relations will be different. In our test, we use the  $\sigma_{11} - \epsilon_{11}$ -relation indicated in the right of Figure B.5 with

$$\sigma_c = 3 \cdot 10^4\text{Pa}, \quad E = \frac{\mu(3\lambda + 2\mu)}{\lambda + \mu} = 3 \cdot 10^8\text{Pa}, \quad E_{\text{res}} = 3 \cdot 10^7\text{Pa},$$

corresponding to the Lamé parameters  $\mu = \frac{3}{26} \cdot 10^9\text{Pa}$  and  $\lambda = \frac{9}{52} \cdot 10^9\text{Pa}$ . For both stress-strain relations we apply Algorithm B.12 with  $\gamma_{\text{lin}} = 0.1$ . We first compare the results to a computation on a very fine mesh to evaluate the remeshing based on the discretization error estimators. Secondly, we perform adaptive remeshing based on these estimators but without applying the adaptive stopping of the Newton iterations and compare the two series of meshes. As in Section B.5.1, we verify if the reduced number of iterations impacts the discretization error.

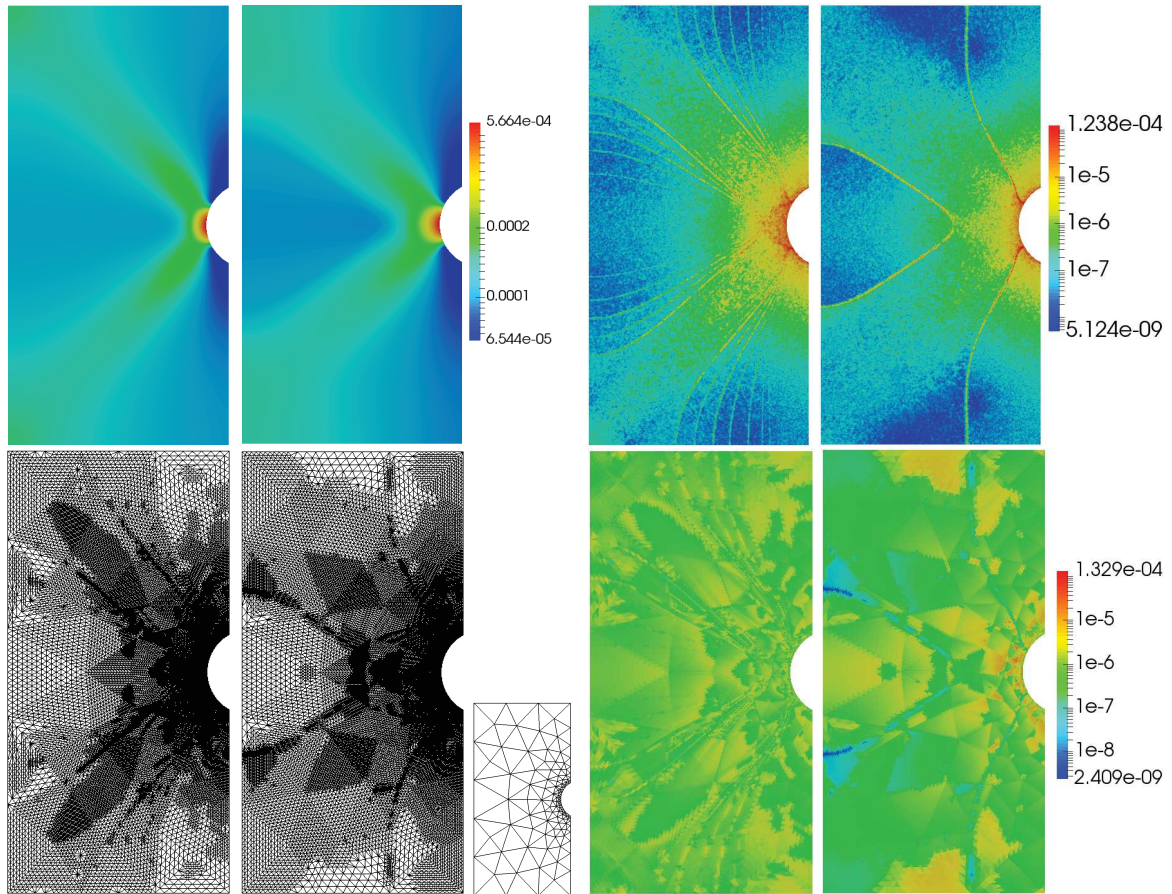


Figure B.6: Notched specimen plate, comparison between Hencky–Mises (left in each picture) and damage model (right). *Top left*:  $\text{tr}(\nabla_s \mathbf{u}_h)$ . *Top right*:  $\eta_{\text{disc}}$  on a fine mesh (no adaptive refinement). *Bottom left*: meshes after six adaptive refinements. *Bottom middle*: initial mesh. *Bottom right*:  $\eta_{\text{disc}}$  on the adaptively refined meshes.

Figure B.6 shows the result of the first part of the test. In each of the four images the left specimen corresponds to the Hencky–Mises and the right to the isotropic damage model. To illustrate the difference of the two models, the top left picture shows the trace of the strain tensor. This scalar value is a good indicator for both models, representing locally the relative volume increase which could correspond to either a damage or shear band localization zone. In the top right picture we see the distribution of the discretization error estimators in the reference computation (209,375 elements), whereas the distribution of the estimators on the sixth adaptively refined mesh is shown in the bottom right picture (60,618 elements for Hencky–Mises, 55,718 elements for the damage model). The corresponding meshes and the initial mesh for the adaptive algorithm are displayed in the bottom left of the figure. To ensure a good discretization of the notch after repeated mesh refinement, the initial mesh cannot be too coarse in this curved area. We observe that the adaptively refined meshes match the distribution of the discretization error estimators on the uniform mesh, and that the estimators are more evenly distributed on these meshes.

The results of the second part of the test are illustrated in Figure B.7. As for the L-shape

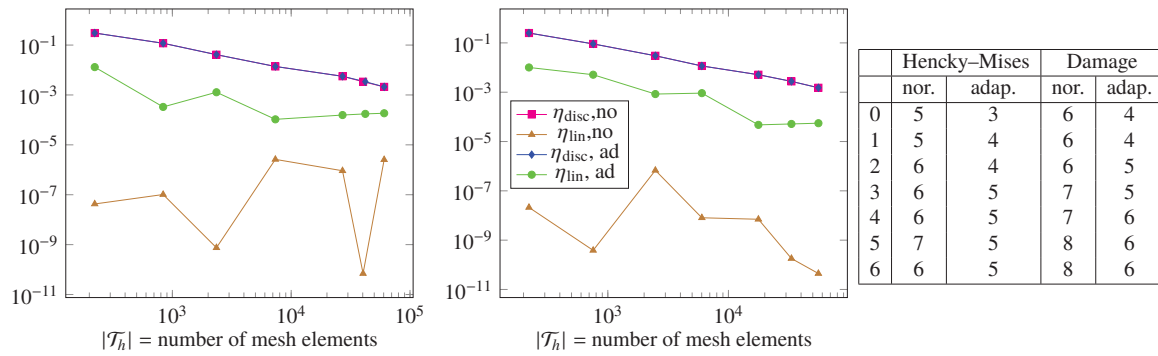


Figure B.7: Notched specimen plate. Comparison of the global discretization and linearization error estimators without and with adaptive stopping criterion for the Hencky–Mises model (left) and the damage model (middle), and comparison of the number of performed Newton iterations (right).

test, we observe that the reduced number of Newton iterations does not affect the discretization error estimate, nor the overall error estimate which is dominated by the discretization error estimate if the Newton algorithm is stopped.



---

## List of Figures

---

1.1	Geologic storage of carbon dioxide in saline aquifers . . . . .	2
1.2	Example of polyhedral mesh with degenerate elements in reservoir applications	11
2.1	Triangular, hexagonal-dominant, Voronoi, and nonmatching quadrangular meshes for the numerical tests. The triangular and nonmatching quadrangular meshes were originally proposed for the FVCA5 benchmark [120], whereas the hexagonal-dominant mesh is the same used in [80, Section 4.2.3]. . . .	40
2.2	Errors vs. $h$ . . . . .	41
2.3	Time convergence rate with BDF2, hexagonal mesh . . . . .	42
2.4	Pressure field on the deformed domain at different times for the finest Cartesian mesh containing 4,192 elements . . . . .	43
2.5	Pressure profiles along the diagonal $(0, 0)$ – $(1, 1)$ of the domain for different normalized times $\hat{t}$ and meshes ( $k = 1$ ). The time step is here $\tau = (2\pi/\beta) \cdot 10^{-2}$ .	43
2.6	Pressure profiles along the diagonal $(0, 0)$ – $(1, 1)$ of the domain for $\kappa = 1 \cdot 10^{-6}$ and time step $\tau = 1 \cdot 10^{-4}$ . Small oscillations are present on the Cartesian mesh (left), whereas no sign of oscillations is present on the hexagonal mesh (right). . . . .	44
3.1	Local DOFs for $k = 1$ (left) and $k = 2$ (right). Shaded DOFs can be locally eliminated by static condensation when solving linearized versions of problem (3.21). . . . .	56
3.2	Triangular, hexagonal-dominant, Voronoi, and nonmatching quadrangular meshes for the numerical tests. The triangular and nonmatching quadrangular meshes were originally proposed for the FVCA5 benchmark [120]. The (predominantly) hexagonal was used in [80]. The Voronoi mesh family was obtained using the PolyMesher algorithm of [188]. . . . .	70
3.3	Tensile test description and resulting stress components for the linear case. Values in $10^5$ Pa . . . . .	71
3.4	Shear test description and resulting stress components for the linear case. Values in $10^5$ Pa . . . . .	71
3.5	Tensile test case: Stress norm on the deformed domain. Values in $10^5$ Pa . .	73
3.6	Shear test case: Stress norm on the deformed domain. Values in $10^4$ Pa . . .	73
3.7	Energy vs $h$ , tensile and shear test cases . . . . .	75



- 5.1 Left: Scatter plot of 5000 realizations  $(\mu, \lambda)$  for  $K \sim \mathcal{U}([1, 20])$  GPa and  $\nu \sim \mathcal{U}([0.1, 0.4])$ . Right: Scatter plot of 5000 realizations  $(K, \nu)$  for  $\mu \sim \mathcal{U}([1, 10])$  GPa and  $\lambda \sim \mathcal{U}([2, 20])$  GPa. . . . . 118
- 5.2 Distribution of the bulk modulus  $K$  (left) and the Poisson ratio  $\nu$  (right) corresponding to  $10^6$  realizations for  $\mu \sim \mathcal{U}([1, 10])$  GPa and  $\lambda \sim \mathcal{U}([2, 20])$  GPa. 119
- 5.3 Distribution of  $10^5$  realizations of  $c_0$  obtained by solving (5.14) with uniformly distributed random input coefficients  $\lambda, \mu$ , and  $\alpha$ . . . . . 120
- 5.4 Examples of 1-D nested sequences including 1, 3, 5, and 9 quadrature points. 128
- 5.5 Examples of tensor-product quadrature grids for  $N = 2$  and  $l = (3, 3)$ . In the full tensorization (top) the colored points correspond to the colored contributions in  $\mathcal{Q}_l^{\text{FT}} = Q_0 \otimes Q_3 + \Delta_1^Q \otimes Q_3 + \Delta_2^Q \otimes Q_3 + \Delta_3^Q \otimes Q_3$ . In Smolyak's tensorization (bottom) the colored points correspond to the colored contributions in  $\mathcal{Q}_L = Q_0 \otimes Q_3 + \Delta_1^Q \otimes Q_2 + \Delta_2^Q \otimes Q_1 + \Delta_3^Q \otimes Q_0$ . . . . . 129
- 5.6 Comparison of polynomials multi-index sets  $\mathcal{K}(\mathcal{L})$  (below the red line) and  $\mathcal{K}^*(\mathcal{L})$  (below the blue line) using a quadrature rule  $\mathcal{Q}_L$  in  $d = 2$  dimensions. Points are multi-indices  $(k_1, k_2)$  such  $F \in P^{k_1}([0, 1]) \times P^{k_2}([0, 1])$  is exactly integrated by  $\mathcal{Q}_L$ . . . . . 131
- 5.7 Pressure plotted on the deformed domain at  $t = t_F = 1$  (data are in KPa). Mean (Left) and Variance (Right) of the PSP method with level  $l = 5$  and  $N_p(l) = 2561$ . . . . . 133
- 5.8 Pressure plotted on the deformed domain at  $t = 0.2$ . Mean (left) and Variance (right) computed using the PSP method with level  $l = 5$  and  $N_p(l) = 2561$ . 134
- 5.9 Errors  $\|\text{MSE}(\mathbf{u} - \mathbf{u}_{\mathcal{K}})\|_{L^2(D)}$  and  $\|\text{MSE}(p - p_{\mathcal{K}})\|_{L^2(D)}$  vs. level  $l$  of the Sparse Grid. Point injection test (left) and footing test (right) with model coefficients  $\mu, \lambda, \alpha, \kappa$  distributed according to (5.36). . . . . 135
- 5.10 Displacement  $\text{MSE}(u_1 - u_{\mathcal{K},1})$  field (left) and pressure  $\text{MSE}(p - p_{\mathcal{K}})$  field (right) of the injection test case. Results are computed using the PSP method with  $l = 1$  (top),  $l = 3$  (middle), and  $l = 5$  (bottom). . . . . 136
- 5.11 Displacement  $\text{MSE}(u_2 - u_{\mathcal{K},2})$  field (left) and pressure  $\text{MSE}(p - p_{\mathcal{K}})$  field (right) corresponding to the footing test case. Results computed using the PSP method with  $l = 1$  (top),  $l = 3$  (middle), and  $l = 5$  (bottom). . . . . 137
- 5.12 Point injection test: Comparison between the total variance  $\text{Var}(u_1(\mathbf{x}, \xi))$  and the sum of the first-order conditional variances  $\sum_{i=1}^4 \text{Var}(u_1(\mathbf{x}, \xi)|\xi_i)$  estimated using the PSP method with  $l = 5$ . . . . . 138
- 5.13 Point injection test: Conditional variances  $\text{Var}(u_1(\mathbf{x}, \xi)|\xi_1)$  (top left)  $\text{Var}(u_1(\mathbf{x}, \xi)|\xi_2)$  (top right)  $\text{Var}(u_1(\mathbf{x}, \xi)|\xi_3)$  (bottom left)  $\text{Var}(u_1(\mathbf{x}, \xi)|\xi_4)$  (bottom right) estimated using the PSP method with  $l = 5$ . . . . . 138
- 5.14 Poroelastic footing test: Comparison between the total variance  $\text{Var}(p(\mathbf{x}, \xi))$  and the sum of the first-order conditional variances  $\sum_{i=1}^4 \text{Var}(p(\mathbf{x}, \xi)|\xi_i)$  estimated using the PSP method with  $l = 5$ . . . . . 139
- 5.15 Poroelastic footing test: Conditional variances  $\text{Var}(p(\mathbf{x}, \xi)|\xi_1)$  (top left)  $\text{Var}(p(\mathbf{x}, \xi)|\xi_2)$  (top right)  $\text{Var}(p(\mathbf{x}, \xi)|\xi_3)$  (bottom left)  $\text{Var}(p(\mathbf{x}, \xi)|\xi_4)$  (bottom right) estimated using the PSP method with  $l = 5$ . . . . . 139

A.1	Convergence tests on a Cartesian mesh family (left) and on a Voronoi mesh family (right). . . . .	146
B.1	Element diagrams for $(\Sigma_T, V_T, \Lambda_T)$ in the case $d = q = 2$ . . . . .	151
B.2	L-shaped domain with linear elasticity model. Distribution of the error estimators (top) and the analytical error (bottom) for the initial mesh (left) and after three (middle) and six (right) adaptive mesh refinements. . . . .	159
B.3	L-shaped domain with linear elasticity model. <i>Left:</i> Comparison of the error estimate (B.14) and $\ \mathbf{u} - \mathbf{u}_h\ _{\text{en}}$ on two series of meshes, obtained by uniform and adaptive remeshing. <i>Right:</i> Effectivity indices of the estimate for each mesh, with the meshes stemming from uniform refinement in italic. . . . .	159
B.4	L-shaped domain with Hencky–Mises model. <i>Left:</i> Comparison of the discretization and linearization error estimators on a series of meshes, without and with adaptive stopping criterion for the Newton algorithm. <i>Middle:</i> Number of Newton iterations without and with adaptive stopping criterion for each mesh. <i>Right:</i> Discretization error estimate for uniform and adaptive remeshing. . . . .	160
B.5	<i>Left:</i> Notched specimen plate. <i>Right:</i> Uniaxial traction curve . . . . .	161
B.6	Notched specimen plate, comparison between Hencky–Mises (left in each picture) and damage model (right). <i>Top left:</i> $\text{tr}(\nabla_s \mathbf{u}_h)$ . <i>Top right:</i> $\eta_{\text{disc}}$ on a fine mesh (no adaptive refinement). <i>Bottom left:</i> meshes after six adaptive refinements. <i>Bottom middle:</i> initial mesh. <i>Bottom right:</i> $\eta_{\text{disc}}$ on the adaptively refined meshes. . . . .	162
B.7	Notched specimen plate. Comparison of the global discretization and linearization error estimators without and with adaptive stopping criterion for the Hencky–Mises model (left) and the damage model (middle), and comparison of the number of performed Newton iterations (right). . . . .	163



---

## Bibliography

---

- [1] M. Abbas, A. Ern, and N. Pignet. *Hybrid High-Order methods for finite deformations of hyperelastic materials*. Submitted. 2018.
- [2] M. Ainsworth and J.T. Oden. “A posteriori error estimation in finite element analysis”. In: *Pure and Applied Mathematics, Wiley-Interscience [John Wiley & Sons], New York* (2000).
- [3] M. Ainsworth and J.T. Oden. “A posteriori error estimators for second order elliptic systems Part 2. An optimal order process for calculating self-equilibrating fluxes”. In: *Computers & Mathematics with Applications* 26.9 (1993), pp. 75–87. DOI: [10.1016/0898-1221\(93\)90007-1](https://doi.org/10.1016/0898-1221(93)90007-1).
- [4] M. Ainsworth and Rankin R. “Guaranteed computable error bounds for conforming and nonconforming finite element analysis in planar elasticity”. In: *Internat. J. Numer. Methods Engrg* 82 (2010), pp. 1114–1157.
- [5] T. Arbogast and Z. Chen. “On the implementation of mixed methods as nonconforming methods for second-order elliptic problems”. In: *Math. Comp.* 64.211 (1995), pp. 943–972.
- [6] C.G. Armstrong, W.M. Lai, and V.C. Mow. “An analysis of the unconfined compression of articular cartilage”. In: *J. Biomech. Eng.* 106.2 (1984), pp. 165–173.
- [7] D. N. Arnold, R. S. Falk, and R. Winther. “Mixed finite element methods for linear elasticity with weakly imposed symmetry”. In: *Math. Comp.* 76 (2007), pp. 1699–1723.
- [8] D. N. Arnold and Winther R. “Mixed finite elements for elasticity”. In: *Numer. Math.* 92 (2002), pp. 401–419.
- [9] J.-L. Auriault and E. Sanchez-Palencia. “Étude du comportement macroscopique d’un milieu poreux saturé déformable”. In: *J. Méc.* 16.4 (1977), pp. 575–603.
- [10] I. Babuška and P. Chatzipantelidis. “On solving elliptic stochastic partial differential equations”. In: *Comput. Methods Appl. Mech. Engrg.* 191.37 (2002), pp. 4093–4122. ISSN: 0045-7825. DOI: [10.1016/S0045-7825\(02\)00354-7](https://doi.org/10.1016/S0045-7825(02)00354-7).
- [11] I. Babuška, F. Nobile, and R. Tempone. “A stochastic collocation method for elliptic partial differential equations with random input data”. In: *SIAM J. Numer. Anal.* 45.3 (2007), pp. 1005–1034. DOI: [10.1137/050645142](https://doi.org/10.1137/050645142).

- [12] I. Babuška, R. Tempone, and G.E. Zouraris. “Galerkin Finite Element Approximations of Stochastic Elliptic Partial Differential Equations”. In: *SIAM J. Numer. Anal.* 42.2 (2004), pp. 800–825. ISSN: 0036-1429. DOI: [10.1137/S0036142902418680](https://doi.org/10.1137/S0036142902418680).
- [13] A. M. Barrientos, N. G. Gatica, and P. E. Stephan. “A mixed finite element method for nonlinear elasticity: two-fold saddle point approach and a-posteriori error estimate”. In: *Numer. Math.* 91.2 (2002), pp. 197–222. DOI: [10.1007/s002110100337](https://doi.org/10.1007/s002110100337).
- [14] S. Barry and G. Mercer. “Exact solution for two-dimensional time dependent flow and deformation within a poroelastic medium”. In: *J. Appl. Mech.* 66.2 (1999), pp. 536–540.
- [15] S.I. Barry and G.K. Aldis. “Comparison of models for flow induced deformation of soft biological tissue”. In: *J. Biomech.* 23.7 (1990), pp. 647–654.
- [16] F. Bassi, L. Botti, A. Colombo, D. A. Di Pietro, and P. Tesini. “On the flexibility of agglomeration based physical space discontinuous Galerkin discretizations”. In: *J. Comput. Phys.* 231.1 (2012), pp. 45–65.
- [17] F. Bassi, S. Rebay, G. Mariotti, S. Pedinotti, and M. Savini. “A high-order accurate discontinuous finite element method for inviscid and viscous turbomachinery flows”. In: *Proceedings of the 2nd European Conference on Turbomachinery Fluid Dynamics and Thermodynamics* (1997), pp. 99–109.
- [18] J. Bear. *Dynamics of Fluids In Porous Media*. Vol. 1. American Elsevier Publishing Company, New York, 1972.
- [19] L. Beirão Da Veiga. “A mimetic discretization method for linear elasticity”. In: *M2AN Math. Model. Numer. Anal.* 44.2 (2010), pp. 231–250.
- [20] L. Beirão da Veiga, F. Brezzi, and L. D. Marini. “Virtual elements for linear elasticity problems”. In: *SIAM J. Numer. Anal.* 2.51 (2013), pp. 794–812.
- [21] L. Beirão da Veiga, C. Lovadina, and D. Mora. “A Virtual Element Method for elastic and inelastic problems on polytope meshes”. In: *Comput. Methods Appl. Mech. Engrg.* 295 (2015), pp. 327–346.
- [22] E. Bemmer, M. Boutéca, O. Vincké, N. Hoteit, and O. Ozanam. “Poromechanics: From linear to nonlinear poroelasticity and poroviscoelasticity”. In: *Oil & Gas Science and Technologies—Rev. IFP* 56.6 (2001), pp. 531–544. DOI: [10.2516/OGST:2001043](https://doi.org/10.2516/OGST:2001043).
- [23] L. Berger, R. Bordas, D. Kay, and S. Tavener. “A stabilized finite element method for finite-strain three-field poroelasticity”. In: *Comput. Mech.* 60 (2017), pp. 51–68. DOI: [10.1007/s00466-017-1381-8](https://doi.org/10.1007/s00466-017-1381-8).
- [24] A. Bespalov, C.E. Powell, and D. Silvester. “A priori error analysis of stochastic Galerkin mixed approximations of elliptic PDEs with random data”. In: *SIAM J. Numer. Anal.* 50.4 (2012), pp. 2039–2063. DOI: [10.1137/110854898](https://doi.org/10.1137/110854898).

- [25] C. Bi and Y. Lin. “Discontinuous Galerkin method for monotone nonlinear elliptic problems”. In: *Int. J. Numer. Anal. Model* 9 (2012), pp. 999–1024.
- [26] S.O.R. Biabanaki, A.R. Khoei, and P. Wriggers. “Polygonal finite element methods for contact-impact problems on non-conformal meshes”. In: *Comput. Meth. Appl. Mech. Engrg.* 269 (2014), pp. 198–221. DOI: [10.1016/j.cma.2013.10.025](https://doi.org/10.1016/j.cma.2013.10.025).
- [27] M. A. Biot. “General theory of threedimensional consolidation”. In: *J. Appl. Phys.* 12.2 (1941), pp. 155–164.
- [28] M. A. Biot. “Nonlinear and semilinear rheology of porous solids”. In: *J. Geoph. Res.* 78.23 (1973), pp. 4924–4937.
- [29] M. A. Biot. “Theory of elasticity and consolidation for a porous anisotropic solid”. In: *J. Appl. Phys.* 26.2 (1955), pp. 182–185.
- [30] D. Boffi, M. Botti, and D. A. Di Pietro. “A nonconforming high-order method for the Biot problem on general meshes”. In: *SIAM J. Sci. Comput.* 38.3 (2016), A1508–A1537. DOI: [10.1137/15M1025505](https://doi.org/10.1137/15M1025505).
- [31] D. Boffi, F. Brezzi, and M. Fortin. *Mixed finite element methods and applications*. Vol. 44. Springer Series in Computational Mathematics. Berlin Heidelberg: Springer, 2013. DOI: [10.1007/978-3-642-36519-5](https://doi.org/10.1007/978-3-642-36519-5).
- [32] L. Botti, D. A. Di Pietro, and J. Droniou. *A Hybrid High-Order discretisation of the Brinkman problem robust in the Darcy and Stokes limits*. Submitted. Mar. 2018. URL: <http://hal.archives-ouvertes.fr/hal-01746367>.
- [33] M. Botti, D. A. Di Pietro, and P. Sochala. “A Hybrid High-Order method for nonlinear elasticity”. In: *SIAM J. Numer. Anal.* 55.6 (2017), pp. 2687–2717. DOI: [10.1137/16M1105943](https://doi.org/10.1137/16M1105943).
- [34] M. Botti, D. A. Di Pietro, and P. Sochala. “A nonconforming high-order method for nonlinear poroelasticity”. In: *Finite Volumes for Complex Applications VIII – Hyperbolic, Elliptic and Parabolic Problems*. 2017, pp. 537–545.
- [35] M. Botti, D. A. Di Pietro, and P. Sochala. *Analysis of a Hybrid High-Order-discontinuous Galerkin discretization method for nonlinear poroelasticity*. Submitted. May 2018. URL: <http://hal.archives-ouvertes.fr/hal-01785810>.
- [36] M. Botti and R. Riedlbeck. “Equilibrated stress tensor reconstruction and a posteriori error estimation for nonlinear elasticity”. In: *Comput. Methods Appl. Math.* (June 2018). DOI: [10.1515/cmam-2018-0012](https://doi.org/10.1515/cmam-2018-0012).
- [37] D. Braess, V. Pillwein, and J. Schöberl. “Equilibrated residual error estimates are  $p$ -robust”. In: *Comput. Methods Appl. Mech. Engrg.* 198 (2009), pp. 1189–1197.
- [38] D. Braess and J. Schöberl. “Equilibrated residual error estimator for edge elements”. In: *Math. Comp.* 77.262 (2008), pp. 651–672.
- [39] S. C. Brenner. “Korn’s inequalities for piecewise  $H^1$  vector fields”. In: *Math. Comp.* 73.247 (2004), pp. 1067–1087. ISSN: 0025-5718.

- [40] H. Brezis. *Functional Analysis, Sobolev Spaces and Partial Differential Equations*. Universitext. Springer New York, 2010. ISBN: 9780387709130.
- [41] F. Brezzi, J. Douglas, and L. D. Marini. “Recent results on mixed finite element methods for second order elliptic problems”. In: *Vistas in applied mathematics. Numerical analysis, atmospheric sciences, immunology*. Ed. by Dorodnitsyn Balakrishnan and Lions Eds. Optimization Software Inc., Publications Division, New York, 1986, pp. 25–43.
- [42] F. Brezzi, G. Manzini, L. D. Marini, P. Pietra, and A. Russo. “Discontinuous Galerkin approximations for elliptic problems”. In: *Numer. Methods Partial Differential Equations* 16.4 (2000), pp. 365–378. ISSN: 0749-159X.
- [43] M. K. Brun, I. Berre, J. M. Nordbotten, and F. A. Radu. “Upscaling of the Coupling of Hydromechanical and Thermal Processes in a Quasi-static Poroelastic Medium”. In: *Transport in Porous Media* 124.1 (2018), pp. 137–158. ISSN: 1573-1634. DOI: [10.1007/s11242-018-1056-8](https://doi.org/10.1007/s11242-018-1056-8).
- [44] H. Bryne and L. Preziosi. “Modeling tumour growth using the theory of mixtures”. In: *Mathematical Medicine and Biology* 20.4 (2003), pp. 341–366.
- [45] R. Burridge and Keller J. B. “Poroelasticity equations derived from microstructure”. In: *J. Acoust. Soc. Am.* 70.4 (1981), pp. 1140–1146. DOI: [10.1121/1.386945](https://doi.org/10.1121/1.386945).
- [46] R. H. Cameron and W. T. Martin. “The orthogonal development of nonlinear functionals in series of Fourier–Hermite functionals”. In: *Ann. Math.* 48 (1947), pp. 385–392.
- [47] A. Cangiani, E. H. Georgoulis, and P. Houston. “*hp*-version discontinuous Galerkin methods on polygonal and polyhedral meshes”. In: *Math. Models Methods Appl. Sci.* 24.10 (2014), pp. 2009–2041. ISSN: 0218-2025. DOI: [10.1142/S0218202514500146](https://doi.org/10.1142/S0218202514500146).
- [48] Y. Cao, S. Chen, and A. J. Mier. “Analysis and numerical approximations of equations of nonlinear poroelasticity”. In: *Discrete & Continuous Dynamical Systems - B* 18 (2013), pp. 1253–1273. DOI: [10.3924/dcdsb.2013.18.1253](https://doi.org/10.3924/dcdsb.2013.18.1253).
- [49] M. Čermák, F. Hecht, Z. Tang, and M. Vohralík. “Adaptive inexact iterative algorithms based on polynomial-degree-robust a posteriori estimates for the Stokes problem”. In: *Numer. Math.* 138.4 (2018), pp. 1027–1065. DOI: [10.1007/s00211-017-0925-3](https://doi.org/10.1007/s00211-017-0925-3).
- [50] M. Cervera, M. Chiumenti, and R. Codina. “Mixed stabilized finite element methods in nonlinear solid mechanics: Part II: Strain localization”. In: *Comput. Methods in Appl. Mech. and Engrg.* 199.37–40 (2010), pp. 2571–2589. DOI: [10.1016/j.cma.2010.04.005](https://doi.org/10.1016/j.cma.2010.04.005).
- [51] L. Chamoin, P. Ladevèze, and F. Pled. “On the techniques for constructing admissible stress fields in model verification: Performances on engineering examples”. In: *Internat. J. Numer. Methods Engrg.* 88.5 (2011), pp. 409–441. DOI: [10.1002/nme.3180](https://doi.org/10.1002/nme.3180).

- [52] C.S. Chang. “Uncertainty in one-dimensional consolidation analysis”. In: *J. Geotech. Eng.* 111.12 (1985), pp. 1411–1424.
- [53] J. Charrier. “Strong and weak error estimates for elliptic partial differential equations with random coefficients”. In: *SIAM J. Numer. Anal.* 50.1 (2012), pp. 216–246. DOI: [10.1137/100800531](https://doi.org/10.1137/100800531).
- [54] H. Chi, L. Beirão da Veiga, and G.H. Paulino. “Some basic formulations of the virtual element method (VEM) for finite deformations”. In: *Computer Methods in Applied Mechanics and Engineering* 318 (2017), pp. 148–192. ISSN: 0045-7825. DOI: [10.1016/j.cma.2016.12.020](https://doi.org/10.1016/j.cma.2016.12.020).
- [55] H. Chi, C. Talischi, O. Lopez-Pamies, and G. H. Paulino. “Polygonal finite elements for finite elasticity”. In: *Int. J. Numer. Methods Eng.* 101.4 (2015), pp. 305–328. ISSN: 1097-0207. DOI: [10.1002/nme.4802](https://doi.org/10.1002/nme.4802).
- [56] P. G. Ciarlet and P. Ciarlet. “Another approach to linearized elasticity and a new proof of Korn’s inequality”. In: *Math. Models Methods Appl. Sci.* 15.02 (2005), pp. 259–271. DOI: [10.1142/S0218202505000352](https://doi.org/10.1142/S0218202505000352).
- [57] B. Cockburn, D. A. Di Pietro, and A. Ern. “Bridging the Hybrid High-Order and Hybridizable Discontinuous Galerkin Methods”. In: *M2AN Math. Model. Numer. Anal.* (2015). Published online. DOI [10.1051/m2an/2015051](https://doi.org/10.1051/m2an/2015051).
- [58] P. R. Conrad and Y. M. Marzouk. “Adaptive Smolyak Pseudospectral Approximations”. In: *SIAM J. Sci. Comp.* 35.6 (2013), A2643–A2670. DOI: [10.1137/120890715](https://doi.org/10.1137/120890715).
- [59] P. G. Constantine, M. S. Eldred, and E. T. Phipps. “Sparse pseudospectral approximation method”. In: *Comput. Methods Appl. Mech. Engrg.* 229 (2012), pp. 1–12.
- [60] P. Cosenza, M. Ghoreychi, G. De Marsily, G. Vasseur, and S. Violette. “Theoretical prediction of poroelastic properties of argillaceous rocks from in situ specific storage coefficient”. In: *Water Resour. Res.* 38.10 (2002), p. 1207. DOI: [10.1029/2001WR001201](https://doi.org/10.1029/2001WR001201).
- [61] O. Coussy. *Poromechanics*. J. Wiley and Sons, Ltd., 2004. DOI: [10.1002/0470092718](https://doi.org/10.1002/0470092718).
- [62] A.A. Darrag and M.A. Tawil. “The consolidation of soils under stochastic initial excess pore pressure”. In: *Applied Mathematical Modelling* 17.11 (1993), pp. 609–612. DOI: [10.1016/0307-904X\(93\)90069-S](https://doi.org/10.1016/0307-904X(93)90069-S).
- [63] R. Davies, G. Foulger, A. Bindley, and P. Styles. “Induced seismicity and hydraulic fracturing for the recovery of hydrocarbons”. In: *Mar. Petrol. Geol.* 45.1 (2013), pp. 171–185. DOI: [10.1016/j.marpetgeo.2013.03.016](https://doi.org/10.1016/j.marpetgeo.2013.03.016).
- [64] K. Deimling. “Nonlinear functional analysis”. In: *Springer-Verlag, Berlin* (1985). DOI: [10.1007/978-3-662-00547-7](https://doi.org/10.1007/978-3-662-00547-7).



- [65] P. Delgado and V. Kumar. “A stochastic Galerkin approach to uncertainty quantification in poroelastic media”. In: *Applied Mathematics and Computation* 266.1 (2015), pp. 328–338. DOI: [10.1016/j.amc.2015.04.127](https://doi.org/10.1016/j.amc.2015.04.127).
- [66] M. Destrade and R. W. Ogden. “On the third- and fourth-order constants of incompressible isotropic elasticity”. In: *The Journal of the Acoustical Society of America* 128.6 (2010), pp. 3334–3343. DOI: [10.1121/1.3505102](https://doi.org/10.1121/1.3505102).
- [67] P. Destuynder and B. Métivet. “Explicit error bounds in a conforming finite element method”. In: *Math. Comp.* 68.228 (1999), pp. 1379–1396.
- [68] E. Detournay and A. H. D. Cheng. “Fundamentals of poroelasticity”. In: *Comprehensive rock engineering. Vol. 2*. Ed. by J.A. Hudson. Pergamon, 1993, pp. 113–171.
- [69] D. A. Di Pietro and J. Droniou. “A Hybrid High-Order method for Leray–Lions elliptic equations on general meshes”. In: *Math. Comp.* 86.307 (2017), pp. 2159–2191. DOI: [10.1142/S0218202517500191](https://doi.org/10.1142/S0218202517500191).
- [70] D. A. Di Pietro and J. Droniou. “ $W^{s,p}$ -approximation properties of elliptic projectors on polynomial spaces, with application to the error analysis of a Hybrid High-Order discretization of Leray–Lions problems”. In: *Math. Models Methods Appl. Sci.* 27.5 (2017), pp. 879–908. DOI: [10.1142/S0218202517500191](https://doi.org/10.1142/S0218202517500191).
- [71] D. A. Di Pietro, J. Droniou, and A. Ern. “A Discontinuous-Skeletal Method for Advection-Diffusion-Reaction on General Meshes”. In: *SIAM J. Numer. Anal.* 53.5 (2015), pp. 2135–2157. DOI: [10.1137/140993971](https://doi.org/10.1137/140993971).
- [72] D. A. Di Pietro, J. Droniou, and G. Manzini. “Discontinuous Skeletal Gradient Discretisation methods on polytopal meshes”. In: *J. Comput. Phys.* 355 (2018), pp. 397–425. DOI: [10.1016/j.jcp.2017.11.018](https://doi.org/10.1016/j.jcp.2017.11.018).
- [73] D. A. Di Pietro and A. Ern. “A family of arbitrary order mixed methods for heterogeneous anisotropic diffusion on general meshes”. In: *IMA J. Numer. Anal.* (2016). Published online. DOI: [10.1093/imanum/drw003](https://doi.org/10.1093/imanum/drw003).
- [74] D. A. Di Pietro and A. Ern. “A hybrid high-order locking-free method for linear elasticity on general meshes”. In: *Comput. Meth. Appl. Mech. Engrg.* 283 (2015), pp. 1–21. DOI: [10.1016/j.cma.2014.09.009](https://doi.org/10.1016/j.cma.2014.09.009).
- [75] D. A. Di Pietro and A. Ern. “Discrete functional analysis tools for discontinuous Galerkin methods with application to the incompressible Navier–Stokes equations”. In: *Math. Comp.* 79 (2010), pp. 1303–1330. DOI: [10.1090/S0025-5718-10-02333-1](https://doi.org/10.1090/S0025-5718-10-02333-1).
- [76] D. A. Di Pietro and A. Ern. *Mathematical aspects of discontinuous Galerkin methods*. Vol. 69. Mathématiques & Applications. Berlin: Springer-Verlag, 2012. DOI: [10.1007/978-3-642-22980-0](https://doi.org/10.1007/978-3-642-22980-0).
- [77] D. A. Di Pietro, A. Ern, and J.-L. Guermond. “Discontinuous Galerkin methods for anisotropic semi-definite diffusion with advection”. In: *SIAM J. Numer. Anal.* 46.2 (2008), pp. 805–831. DOI: [10.1137/060676106](https://doi.org/10.1137/060676106).

- [78] D. A. Di Pietro, A. Ern, and S. Lemaire. “An arbitrary-order and compact-stencil discretization of diffusion on general meshes based on local reconstruction operators”. In: *Comput. Methods Appl. Math.* 14.4 (2014), pp. 461–472.
- [79] D. A. Di Pietro, B. Kapidani, R. Specogna, and F. Trevisan. “An Arbitrary-Order Discontinuous Skeletal Method for Solving Electrostatics on General Polyhedral Meshes”. In: *IEEE Transactions on Magnetics* 53.6 (June 2017), pp. 1–4. ISSN: 0018-9464. DOI: [10.1109/TMAG.2017.2666546](https://doi.org/10.1109/TMAG.2017.2666546).
- [80] D. A. Di Pietro and S. Lemaire. “An extension of the Crouzeix–Raviart space to general meshes with application to quasi-incompressible linear elasticity and Stokes flow”. In: *Math. Comp.* 84.291 (2015), pp. 1–31.
- [81] D. A. Di Pietro and S. Nicaise. “A locking-free discontinuous Galerkin method for linear elasticity in locally nearly incompressible heterogeneous media”. In: *Appl. Numer. Math.* 63 (2013), pp. 105–116. ISSN: 0168-9274. DOI: [10.1016/j.apnum.2012.09.009](https://doi.org/10.1016/j.apnum.2012.09.009).
- [82] D. A. Di Pietro and R. Tittarelli. “Numerical Methods for PDEs. State of the Art Techniques”. In: ed. by L. Formaggia D. A. Di Pietro A. Ern. SEMA-SIMAI 15. Springer, 2018. Chap. An introduction to Hybrid High-Order methods. URL: <http://arxiv.org/abs/1703.05136>.
- [83] P. Dörsek and J. Melenk. “Symmetry-free,  $p$ -robust equilibrated error indication for the  $hp$ -version of the FEM in nearly incompressible linear elasticity”. In: *Comput. Methods Appl. Math.* 13 (2013), pp. 291–304.
- [84] J. Droniou. “Finite volume schemes for fully non-linear elliptic equations in divergence form”. In: *ESAIM: Mathematical Modelling and Numerical Analysis - Modélisation Mathématique et Analyse Numérique* 40.6 (2006), pp. 1069–1100.
- [85] J. Droniou, R. Eymard, T. Gallouët, C. Guichard, and R. Herbin. *The gradient discretisation method*. Preprint [hal-01366646](https://hal.archives-ouvertes.fr/hal-01366646). Nov. 2016.
- [86] J. Droniou, R. Eymard, T. Gallouët, C. Guichard, and R. Herbin. *The gradient discretisation method: A framework for the discretisation and numerical analysis of linear and nonlinear elliptic and parabolic problems*. Maths & Applications. To appear. Springer, 2017, 509p. URL: <https://hal.archives-ouvertes.fr/hal-01382358>.
- [87] J. Droniou and B. P. Lamichhane. “Gradient Schemes for Linear and Non-linear Elasticity Equations”. In: *Numer. Math.* 129.2 (2015), pp. 251–277. ISSN: 0029-599X. DOI: [10.1007/s00211-014-0636-y](https://doi.org/10.1007/s00211-014-0636-y).
- [88] T. Dupont and R. Scott. “Polynomial approximation of functions in Sobolev spaces”. In: *Math. Comp.* 34.150 (1980), pp. 441–463.
- [89] G. Duvaut and Lions J.L. *Inequalities in Mechanics and Physics*. 1st ed. Grundlehren der mathematischen Wissenschaften 219. Springer-Verlag Berlin Heidelberg, 1976. ISBN: 978-3-642-66167-9,978-3-642-66165-5.

- [90] W. Ehlers and G. Eipper. “Finite Elastic Deformations in Liquid-Saturated and Empty Porous Solids”. In: *Transport in Porous Media* 34.1 (1999), pp. 179–191. ISSN: 1573-1634. DOI: [10.1023/A:1006565509095](https://doi.org/10.1023/A:1006565509095).
- [91] L. El Alaoui, A. Ern, and M. Vohralík. “Guaranteed and robust a posteriori error estimates and balancing discretization and linearization errors for monotone nonlinear problems”. In: *Comput. Methods Appl. Mech. Engrg.* 200 (2011), pp. 2782–2795.
- [92] A. Ern and J.-L. Guermond. *Theory and Practice of Finite Elements*. Vol. 159. Applied Mathematical Sciences. New York, NY: Springer-Verlag, 2004.
- [93] A. Ern, A. F. Stephansen, and P. Zunino. “A discontinuous Galerkin method with weighted averages for advection-diffusion equations with locally small and anisotropic diffusivity”. In: *IMA J. Numer. Anal.* 29.2 (2009), pp. 235–256. ISSN: 0272-4979. DOI: [10.1093/imanum/drm050](https://doi.org/10.1093/imanum/drm050).
- [94] A. Ern and M. Vohralík. “A posteriori error estimation based on potential and flux reconstruction for the heat equation”. In: *SIAM J. Numer. Anal.* 48.1 (2010), pp. 198–223.
- [95] A. Ern and M. Vohralík. “Adaptive inexact Newton methods with a posteriori stopping criteria for nonlinear diffusion PDEs”. In: *SIAM J. Sci. Comput.* 35.4 (2013), A1761–A1791.
- [96] A. Ern and M. Vohralík. *Broken stable  $H^1$  and  $H(\text{div})$  polynomial extensions for polynomial-degree-robust potential and flux reconstruction in three space dimensions*. preprint hal-01422204. 2017.
- [97] A. Ern and M. Vohralík. “Polynomial-degree-robust a posteriori estimates in a unified setting for conforming, nonconforming, discontinuous Galerkin, and mixed discretizations”. In: *SIAM J. Numer. Anal.* 53.2 (2015), pp. 1058–1081.
- [98] O. G. Ernst, A. Mugler, H.-J. Starkloff, and E. Ullmann. “On the convergence of generalized polynomial chaos expansions”. In: *ESAIM: Mathematical Modelling and Numerical Analysis* 46 (02 Mar. 2012), pp. 317–339. ISSN: 1290-3841. DOI: [10.1051/m2an/2011045](https://doi.org/10.1051/m2an/2011045).
- [99] L. C. Evans and R. F. Gariepy. *Measure theory and fine properties of functions*. Studies in advanced mathematics. Boca Raton (Fla.): CRC Press, 1992. URL: <http://opac.inria.fr/record=b1089059>.
- [100] D.G. Frias, M.A. Murad, and F. Pereira. “Stochastic computational modelling of highly heterogeneous poroelastic media with long-range correlations”. In: *Int. J. Numer. Anal. Meth. Geomech.* 28.1 (2004), pp. 1–32. DOI: [10.1002/nag.323](https://doi.org/10.1002/nag.323).
- [101] G. Fu. *A high-order HDG method for the Biot’s consolidation model*. Submitted. 2018.

- [102] F.J. Gaspar, F.J. Lisbona, C.W. Oosterlee, and P.N. Vabishchevich. “An efficient multigrid solver for a reformulated version of the poroelasticity system”. In: *Comput. Methods Appl. Mech. and Engrg.* 196 (2007), pp. 1447–1457. DOI: [10.1016/j.cma.2006.03.020](https://doi.org/10.1016/j.cma.2006.03.020).
- [103] F. Gassmann. “On Elasticity of Porous Media”. In: *Vierteljahresheft der Naturforschenden Gessellschaft* 96 (1951), pp. 1–23.
- [104] G. N. Gatica, A. Márquez, and W. Rudolph. “A priori and a posteriori error analyses of augmented twofold saddle point formulations for nonlinear elasticity problems”. In: *Comput. Meth. Appl. Mech. Engrg.* 264 (2013), pp. 23–48. DOI: [10.1016/j.cma.2013.05.010](https://doi.org/10.1016/j.cma.2013.05.010).
- [105] G. N. Gatica and E. P. Stephan. “A mixed-FEM formulation for nonlinear incompressible elasticity in the plane”. In: *Numerical Methods for Partial Differential Equations* 18.1 (2002), pp. 105–128. ISSN: 1098-2426. DOI: [10.1002/num.1046](https://doi.org/10.1002/num.1046).
- [106] G. N. Gatica and W. L. Wendland. “Coupling of Mixed Finite Elements and Boundary Elements for A Hyperelastic Interface Problem”. In: *SIAM J. on Numer. Anal.* 34.6 (1997), pp. 2335–2356. DOI: [10.1137/S0036142995291317](https://doi.org/10.1137/S0036142995291317).
- [107] B. Gatmiri and P. Delage. “A formulation of fully coupled thermal-hydraulic-mechanical behavior of saturated porous media—numerical approach”. In: *Int. J. Numerical and Anal. Meth. Geomech.* 21.3 (1997), pp. 199–225. DOI: [10.1002/\(SICI\)1096-9853\(199703\)21:3<199::AID-NAG865>3.0.CO;2-M](https://doi.org/10.1002/(SICI)1096-9853(199703)21:3<199::AID-NAG865>3.0.CO;2-M).
- [108] T. Gerstner and M. Griebel. “Dimension-Adaptive Tensor-Product Quadrature”. In: *Computing* 71.1 (2003), pp. 65–87.
- [109] T. Gerstner and M. Griebel. “Numerical Integration Using Sparse Grids”. In: *Numerical Algorithms* (1998), pp. 209–232.
- [110] R. G. Ghanem. “Probabilistic characterization of transport in heterogeneous media”. In: *Comput. Meth. Appl. Mech. Eng.* 158.3 (1998), pp. 199–220. ISSN: 0045-7825. DOI: [10.1016/S0045-7825\(97\)00250-8](https://doi.org/10.1016/S0045-7825(97)00250-8).
- [111] R. G. Ghanem and S. Dham. “Stochastic Finite Element Analysis for Multiphase Flow in Heterogeneous Porous Media”. In: *Transport in Porous Media* 32.3 (1998), pp. 239–262. ISSN: 1573-1634. DOI: [10.1023/A:1006514109327](https://doi.org/10.1023/A:1006514109327).
- [112] R. G. Ghanem and S. D. Spanos. *Stochastic Finite Elements: a Spectral Approach*. Springer Verlag, 1991.
- [113] S.M. Gorelick and M. D. Zoback. “Earthquake triggering and large-scale geologic storage of carbon dioxide”. In: *Proc. Natl. Acad. Sci. U.S.A.* 109.26 (2012), pp. 10164–10168. DOI: [10.1073/pnas.1202473109](https://doi.org/10.1073/pnas.1202473109).
- [114] D. H. Green and H. F. Wang. “Specific storage as a poroelastic coefficient”. In: *Water Resour. Res.* 26.7 (1990), pp. 1631–1637. DOI: [10.1029/WR026i007p01631](https://doi.org/10.1029/WR026i007p01631).

- [115] P. Grisvard. *Elliptic problems in nonsmooth domains*. Pitman Advanced Pub. Program. London, 1985.
- [116] G. Guennebaud and B. Jacob. *Eigen v3*. <http://eigen.tuxfamily.org>. 2010.
- [117] J. B. Haga, H. Osnes, and H. P. Langtangen. “On the causes of pressure oscillations in low-permeable and low-compressible porous media”. In: *Int. J. Numer. Anal. Methods Geomech.* 36.12 (2012), pp. 1507–1522. DOI: [10.1002/nag.1062](https://doi.org/10.1002/nag.1062).
- [118] A. Hannukainen, R. Stenberg, and M. Vohralík. “A unified framework for a posteriori error estimation for the Stokes problem”. In: *Numer. Math.* 122 (2012), pp. 725–769.
- [119] T. D. Hein and M. Kleiber. “Stochastic Finite Element modeling in linear transient heat transfer”. In: *Comput. Meth. Appl. Mech. Eng.* 144 (1997), pp. 111–124.
- [120] R. Herbin and F. Hubert. “Benchmark on discretization schemes for anisotropic diffusion problems on general grids”. In: *Finite Volumes for Complex Applications V*. Ed. by R. Eymard and J.-M. Hérard. John Wiley & Sons, 2008, pp. 659–692.
- [121] J. G. Heywood and R. Rannacher. “Finite-element approximation of the nonstationary Navier–Stokes problem. Part IV: error analysis for second-order time discretization”. In: *SIAM J. Numer. Anal.* 27.2 (1990), pp. 353–384.
- [122] H. P. Hong. “One-dimensional consolidation with uncertain properties”. In: *Canadian Geotechnical journal* 29 (1991), pp. 161–165.
- [123] L. Hu, P. H. Winterfield, P. Fakcharoenphol, and Wu Y. S. “A novel fully-coupled flow and geomechanics model in enhanced geothermal reservoirs”. In: *J. Pet. Sci. Eng.* 107 (2013), pp. 1–11. DOI: [10.1016/j.petrol.2013.04.005](https://doi.org/10.1016/j.petrol.2013.04.005).
- [124] D. S. Hughes and J. L. Kelly. “Second-Order Elastic Deformation of Solids”. In: *Phys. Rev.* 92 (5 1953), pp. 1145–1149.
- [125] M. Iskandarani, S. Wang, A. Srinivasan, W. C. Thacker., J. Winokur, and O. Knio. “An overview of uncertainty quantification techniques with application to oceanic and oil-spill simulations”. In: *J. Geophys. Res.: Oceans* 121.4 (2016), pp. 2789–2808. ISSN: 2169-9291. DOI: [10.1002/2015JC011366](https://doi.org/10.1002/2015JC011366).
- [126] B. Jah and R. Juanes. “Coupled multiphase flow and poromechanics: A computational model of pore pressure effects on fault slip and earthquake triggering”. In: *Water Resour. Res.* 50.5 (2014), pp. 3776–3808. DOI: [10.1002/2013WR015175](https://doi.org/10.1002/2013WR015175).
- [127] L. Jin and M. D. Zoback. “Fully coupled nonlinear fluid flow and poroelasticity in arbitrarily fractured porous media: A hybrid-dimensional computational model”. In: *Journal Geophysical Research: Solid Earth* 122 (2017), pp. 7626–7658. DOI: [10.1002/2017JB014892](https://doi.org/10.1002/2017JB014892).
- [128] A. Khan, C. E. Powell, and D. J. Silvester. *Robust preconditioning for stochastic Galerkin formulations of parameter-dependent linear elasticity equations*. Submitted. Mar. 2018.

- [129] Kwang-Yeon Kim. “Guaranteed a posteriori error estimator for mixed finite element methods of linear elasticity with weak stress symmetry”. In: *SIAM J. Numer. Anal.* 48.6 (2011), pp. 2364–2385.
- [130] A.E. Kolesov, P.N. Vabishchevich, and M.V. Vasilyeva. “Splitting scheme for poroelasticity and thermoelasticity problems”. In: *Comput. and Math. with Appl.* 67 (2014), pp. 2185–2198. DOI: [10.1016/j.camwa.2014.02.005](https://doi.org/10.1016/j.camwa.2014.02.005).
- [131] J. Korsawe, G. Starke, W. Wang, and O. Kolditz. “Finite element analysis of poro-elastic consolidation in porous media: Standard and mixed approaches”. In: *Comput. Methods Appl. Mech. Engrg.* 195 (2006), pp. 1096–1115. DOI: [10.1016/j.cma.2005.04.011](https://doi.org/10.1016/j.cma.2005.04.011).
- [132] P. Ladevèze. “Comparaison de modèles de milieux continus”. PhD thesis. Université Pierre et Marie Curie (Paris 6), 1975.
- [133] P. Ladevèze and D. Leguillon. “Error estimate procedure in the finite element method and applications”. In: *SIAM J. Numer. Anal.* 20 (1983), pp. 485–509.
- [134] P. Ladevèze, J. P. Pelle, and P. Rougeot. “Error estimation and mesh optimization for classical finite elements”. In: *Engrg. Comp.* 8.1 (1991), pp. 69–80.
- [135] L. D. Landau and E. M. Lifshitz. *Theory of elasticity*. Pergamon London, 1959, 134p.
- [136] O. P. Le Maître and O. M. Knio. *Spectral Methods for Uncertainty Quantification*. Scientific Computation. Springer, 2010.
- [137] O. Le Maître, M. T. Reagan, H. N. Najm, R. G. Ghanem, and O. M. Knio. “A Stochastic Projection Method for Fluid Flow II.: Random Process”. In: *J. Comput. Phys.* 181.1 (2002), pp. 9–44. ISSN: 0021-9991. DOI: [10.1006/jcph.2002.7104](https://doi.org/10.1006/jcph.2002.7104).
- [138] S. Lemaire. “Discretisations non-conformes d’un modèle poromécanique sur maillages généraux”. PhD thesis. Université Paris-Est, Dec. 2013.
- [139] S. E. Leon, D. W. Spring, and G. H. Paulino. “Reduction in mesh bias for dynamic fracture using adaptive splitting of polygonal finite elements”. In: *Int. J. Numer. Methods Eng.* 100.8 (2014), pp. 555–576. ISSN: 1097-0207. DOI: [10.1002/rme.4744](https://doi.org/10.1002/rme.4744).
- [140] R. W. Lewis and B. A. Schrefler. *The finite element method in the static and dynamic deformation and consolidation of porous media*. Numerical methods in engineering. John Wiley, 1998. ISBN: 9780471928096.
- [141] H. Liu, B. Hu, and Z. W. Yu. “Stochastic finite element method for random temperature in concrete structures”. In: *Int. J. Solids Struct.* 38.1 (2001), pp. 6965–6983.
- [142] M. Loève. *Probability theory*. Vol. 2. Springer-Verlang, New York, 1977.
- [143] A. F. D. Loula and J. N. C. Guerreiro. “Finite element analysis of nonlinear creeping flows”. In: *Comput. Meth. Appl. Mech. Engrg.* 79.1 (1990), pp. 87–109.

- [144] R. Luce and B. I. Wohlmuth. “A local a posteriori error estimator based on equilibrated fluxes”. In: *SIAM J. Numer. Anal.* 42 (2004), pp. 1394–1414.
- [145] P. Luo, C. Rodrigo, F. J. Gaspar, and C. W. Oosterlee. “Multigrid method for nonlinear poroelasticity equations”. In: *Comput. Visual. Sci.* 17 (2015), pp. 255–265. DOI: [10.1007/s00791-016-0260-8](https://doi.org/10.1007/s00791-016-0260-8).
- [146] Bowen R. M. “Compressible porous media models by use of the theory of mixtures”. In: *Int. J. Eng. Sci.* 20.6 (1982), pp. 697–735. DOI: [10.1016/0020-7225\(82\)90082-9](https://doi.org/10.1016/0020-7225(82)90082-9).
- [147] Bowen R. M. “Incompressible porous media models by use of the theory of mixtures”. In: *Int. J. Eng. Sci.* 18.9 (1980), pp. 1129–1148. DOI: [10.1016/0020-7225\(80\)90114-7](https://doi.org/10.1016/0020-7225(80)90114-7).
- [148] H. G. Matthies and A. Keese. “Galerkin methods for linear and nonlinear elliptic stochastic partial differential equations”. In: *Comput. Methods Appl. Mech. Engrg.* 194.12-16 (2005), pp. 1295–1331.
- [149] V. Maz’ya and J. Rossmann. *Elliptic equations in polyhedral domains*. Vol. 162. AMS: Mathematical Surveys and Monographs, 2010.
- [150] A. Mehrabian and Y. N. Abousleiman. “Gassmann equation and the constitutive relations for multiple-porosity and multiple-permeability poroelasticity with applications to oil and gas shale”. In: *Int. J. Numer. Anal. Meth. Geomech.* 39 (2015), pp. 1547–1569. DOI: [10.1002/nag.2399](https://doi.org/10.1002/nag.2399).
- [151] S. E. Minkoff, C. M. Stone, S. Bryant, M. Peszynsak, and M. F. Wheeler. “Coupled fluid flow and geomechanical deformation modeling”. In: *J. Pet. Sci. Eng.* 38 (2003), pp. 37–56. DOI: [10.1016/S0920-4105\(03\)00021-4](https://doi.org/10.1016/S0920-4105(03)00021-4).
- [152] G. J. Minty. “On a “monotonicity” method for the solution of nonlinear equations in Banach spaces”. In: *Proceedings of the National Academy of Sciences of the United States of America* 50.6 (1963), pp. 1038–1041.
- [153] D. Moos, P. Peska, T. Finkbeiner, and M.D. Zoback. “Comprehensive Wellbore Stability Analysis Utilizing Quantitative Risk Assessment”. In: *J. Pet. Sci. Eng.* 38.3 (2003), pp. 97–109. DOI: [10.1016/S0920-4105\(03\)00024-X](https://doi.org/10.1016/S0920-4105(03)00024-X).
- [154] M. A. Murad and F. D. Loula. “Improved accuracy in finite element analysis of Biot’s consolidation problem”. In: *Comput. Methods Appl. Mech. Engrg.* 93.3 (1992), pp. 359–382.
- [155] M. A. Murad and F. D. Loula. “On stability and convergence of finite element approximations of Biot’s consolidation problem”. In: *Interat. J. Numer. Methods Engrg.* 37.4 (1994). DOI: [10.1002/nme.1620370407](https://doi.org/10.1002/nme.1620370407).
- [156] A. Naumovich. “On finite volume discretization of the three-dimensional Biot poroelasticity system in multilayer domains”. In: *Comput. Meth. App. Math.* 6.3 (2006), pp. 306–325.

- [157] J. Nečas. *Introduction to the theory of nonlinear elliptic equations*. Chichester: A Wiley-Interscience Publication. John Wiley & Sons Ltd., 1986. Reprint of the 1983 edition.
- [158] S. Nicaise, K. Witowski, and B. I. Wohlmuth. “An a posteriori error estimator for the Lamé equation based on H(div)-conforming stress approximations”. In: *IMA J. Numer. Anal.* 28.2 (2008), pp. 331–353. DOI: [10.1093/imanum/drm008](https://doi.org/10.1093/imanum/drm008).
- [159] S. Nishimura, K. Shimada, and H. Fujii. “Consolidation inverse analysis considering spacial variability and non-linearity of soil parameters”. In: *Soils and Foundations* 42.3 (2002), pp. 45–61. DOI: [10.3208/sandf.42.3\\_45](https://doi.org/10.3208/sandf.42.3_45).
- [160] J. M. Nordbotten. “Cell-centered finite volume discretizations for deformable porous media”. In: *Int. J. Numer. Meth. Eng.* 100.6 (2014), pp. 399–418. DOI: [10.1002/nme.4734](https://doi.org/10.1002/nme.4734).
- [161] S. Ohnibus, E. Stein, and E. Walhorn. “Local error estimates of FEM for displacements and stresses in linear elasticity by solving local Neumann problems”. In: *Int. J. Numer. Meth. Engng.* 52 (2001), pp. 727–746.
- [162] C.M. Oldenburg. “The risk of induced seismicity: Is cap-rock on shaly ground ?” In: *Greenh. Gases: Sci. Technol.* 2.4 (2012), pp. 217–218. DOI: [10.1002/ghg.1299](https://doi.org/10.1002/ghg.1299).
- [163] C. Ortner and E. Süli. “Discontinuous Galerkin Finite Element Approximation of Nonlinear Second-Order Elliptic and Hyperbolic Systems”. In: *SIAM J. on Numer. Anal.* 45.4 (2007), pp. 1370–1397. DOI: [10.1137/06067119X](https://doi.org/10.1137/06067119X).
- [164] P. J. Phillips. “Finite element methods in linear poroelasticity: Theoretical and computational results”. PhD thesis. University of Texas at Austin, Dec. 2005.
- [165] P. J. Phillips and M. F. Wheeler. “A coupling of mixed and continuous Galerkin finite element methods for poroelasticity I: the continuous in time case”. In: *Comput. Geosci.* 11 (2007), pp. 131–144.
- [166] P. J. Phillips and M. F. Wheeler. “A coupling of mixed and continuous Galerkin finite element methods for poroelasticity II: the discrete-in-time case”. In: *Comput. Geosci.* 11 (2007), pp. 145–158.
- [167] P. J. Phillips and M. F. Wheeler. “A coupling of mixed and discontinuous Galerkin finite-element methods for poroelasticity”. In: *Comput. Geosci.* 12 (2008), pp. 417–435.
- [168] P. J. Phillips and M. F. Wheeler. “Overcoming the problem of locking in linear elasticity and poroelasticity: An heuristic approach”. In: *Comput. Geosci.* 13 (2009), pp. 5–12. DOI: [10.1007/s10596-008-9114-x](https://doi.org/10.1007/s10596-008-9114-x).
- [169] M. Pitteri and G. Zanutto. *Continuum models for phase transitions and twinning in crystals*. Chapman & Hall/CRC, 2003.
- [170] W. Prager and J. L. Synge. “Approximations in elasticity based on the concept of function space”. In: *Quart. Appl. Math.* 5 (1947), pp. 241–269.



- [171] M. T. Reagan, H. N. Najm, R. G. Ghanem, and O. M. Knio. “Uncertainty quantification in reacting-flow simulations through non-intrusive spectral projection”. In: *Combustion and Flame* 132.3 (2003), pp. 545–555.
- [172] S. I. Repin. *A posteriori estimates for partial differential equations*. Vol. 4. Radon Series on Computational and Applied Mathematics. Walter de Gruyter GmbH & Co. KG, Berlin, 2008.
- [173] J. R. Rice and M. P. Cleary. “Some basic stress diffusion solutions for fluid-saturated elastic porous media with compressible constituents”. In: *Rev. Geophys.* 14.2 (1976), pp. 227–241. DOI: [10.1029/RG014i002p00227](https://doi.org/10.1029/RG014i002p00227).
- [174] R. Riedlbeck, D. A. Di Pietro, A. Ern, S. Granet, and K. Kazymyrenko. “Stress and flux reconstruction in Biot’s poro-elasticity problem with application to a posteriori analysis”. In: *Comput. and Math. with Appl.* 73.7 (2017), pp. 1593–1610.
- [175] R. Riedlbeck, D. A. Di Pietro, and E. Ern. “Equilibrated stress reconstructions for linear elasticity problems with application to a posteriori error analysis”. In: *Finite Volumes for Complex Applications VIII – Methods and Theoretical Aspects*. 2017, pp. 293–301.
- [176] C. Rodrigo, F.J. Gaspar, X. Hu, and L.T. Zikatanov. “Stability and monotonicity for some discretizations of the Biot’s consolidation model”. In: *Comput. Methods Appl. Mech. and Engrg.* 298 (2016), pp. 183–204. DOI: [10.1016/j.cma.2015.09.019](https://doi.org/10.1016/j.cma.2015.09.019).
- [177] D. Sandri. “Sur l’approximation numérique des écoulements quasi-Newtoniens dont la viscosité suit la loi puissance ou la loi de Carreau”. In: *Math. Modelling and Num. Anal.* 27.2 (1993), pp. 131–155.
- [178] O.D. Schirra. “New Korn-type inequalities and regularity of solutions to linear elliptic systems and anisotropic variational problems involving the trace-free part of the symmetric gradient”. In: *Calc. Var.* 43 (2012), pp. 147–172. DOI: [10.1007/s00526-011-0406-y](https://doi.org/10.1007/s00526-011-0406-y).
- [179] D. Schötzau and Schwab C. “Time Discretization of Parabolic Problems by the HP-Version of the Discontinuous Galerkin Finite Element Method”. In: *SIAM J. Numer. Anal.* 38.3 (2000), pp. 837–875. DOI: [10.1137/S0036142999352394](https://doi.org/10.1137/S0036142999352394).
- [180] C. Schwab and R. A. Todor. “Karhunen-Loève Approximation of Random Fields by Generalized Fast Multipole Methods”. In: *J. Comput. Phys.* 217.1 (2006), pp. 100–122. ISSN: 0021-9991. DOI: [10.1016/j.jcp.2006.01.048](https://doi.org/10.1016/j.jcp.2006.01.048).
- [181] R. E. Showalter. “Diffusion in poro-elastic media”. In: *J. Math. Anal. Appl.* 251 (2000), pp. 310–340. DOI: [10.1006/jmaa.2000.7048](https://doi.org/10.1006/jmaa.2000.7048).
- [182] R. Shukla, P. Ranjith, S. Choi, and A. Haque. “Study of Caprock Integrity in Geosequestration of Carbon Dioxide”. In: *Int. J. Geomech.* 11.4 (2010), pp. 294–301. DOI: [10.1061/\(ASCE\)GM.1943-5622.0000015](https://doi.org/10.1061/(ASCE)GM.1943-5622.0000015).
- [183] I. Smears. “Robust and efficient preconditioners for the discontinuous Galerkin time-stepping method”. In: *IMA J. Numer. Anal.* 37.4 (2017), pp. 1961–1985. DOI: [10.1093/imanum/drw050](https://doi.org/10.1093/imanum/drw050).

- [184] S.A. Smolyak. “Quadrature and interpolation formulas for tensor products of certain classes of functions”. In: *Dokl. Akad. Nauk SSSR* 4.240-243 (1963), p. 123.
- [185] S.-C. Soon, B. Cockburn, and H. K. Stolarski. “A hybridizable discontinuous Galerkin method for linear elasticity”. In: *Internat. J. Numer. Methods Engrg.* 80.8 (2009), pp. 1058–1092.
- [186] D. W. Spring, S. E. Leon, and G. H. Paulino. “Unstructured polygonal meshes with adaptive refinement for the numerical simulation of dynamic cohesive fracture”. In: *Int. J. Fracture* 189.1 (2014), pp. 33–57. ISSN: 1573-2673. DOI: [10.1007/s10704-014-9961-5](https://doi.org/10.1007/s10704-014-9961-5).
- [187] C.C. Swan, R.S. Lakes, R.A. Brand, and K.J. Stewart. “Micromechanically based poroelastic modeling of fluid flow in haversian bone”. In: *J. Biomech. Eng.* 125.1 (2003), pp. 25–37.
- [188] C. Talischi, G. H. Paulino, A. Pereira, and I. F. M. Menezes. “PolyMesher: a general-purpose mesh generator for polygonal elements written in Matlab”. In: *Structural and Multidisciplinary Optimization* 45.3 (2012), pp. 309–328. DOI: [10.1007/s00158-011-0706-z](https://doi.org/10.1007/s00158-011-0706-z).
- [189] K. Terzaghi. “Die berechnung der durchlässigkeitszier des tones aus dem verlauf der hydrodynamischen spannungserscheinungen”. In: *Sitz. Akad. Wissen.* 132.2a (1923), pp. 105–124.
- [190] K. Terzaghi. *Theoretical soil mechanics*. New York: Wiley, 1943.
- [191] L. R. G. Treolar. *The Physics of Rubber Elasticity*. Oxford: Clarendon Press, 1975.
- [192] J. Troyen, O. Le Maître, M. Ndjinga, and A. Ern. “Intrusive projection methods with upwinding for uncertain nonlinear hyperbolic systems”. In: *J. Comput. Phys.* 229.18 (2010), pp. 6485–6511. DOI: [10.1016/j.cam.2010.05.043](https://doi.org/10.1016/j.cam.2010.05.043).
- [193] B. Tully and Y. Ventikos. “Cerebral water transport using multiple-network poroelastic theory: application to normal pressure hydrocephalus”. In: *J. Fluid Mech.* 667 (2011), pp. 188–215. DOI: [10.1017/S0022112010004428](https://doi.org/10.1017/S0022112010004428).
- [194] R. Verfürth. *A review of a posteriori error estimation and adaptive mesh-refinement techniques*. Wiley-Teubner, Chichester, 1996.
- [195] R. Verfürth. “A review of a posteriori error estimation techniques for elasticity problems”. In: *Comput. Meth. Appl. Mech. Engrg.* 176 (1999), pp. 419–440.
- [196] M. Vogelius. “An analysis of the p-version of the finite element method for nearly incompressible materials. Uniformly valid, optimal error estimates”. In: *Numer. Math.* 41 (1983), pp. 39–53.
- [197] C. Wang, J. Wang, R. Wang, and R. Zhang. “A Locking-free Weak Galerkin Finite Element Method for Elasticity Problems in the Primal Formulation”. In: *J. Comput. Appl. Math.* 307.C (2016), pp. 346–366. ISSN: 0377-0427. DOI: [10.1016/j.cam.2015.12.015](https://doi.org/10.1016/j.cam.2015.12.015).

- [198] H. F. Wang. *Theory of linear poroelasticity with applications to geomechanics and hydrogeology*. Princeton University Press, 2000.
- [199] J. Wang and X. Ye. “A weak Galerkin mixed finite element method for second order elliptic problems”. In: *Math. Comp.* 83.289 (2014), pp. 2101–2126.
- [200] M. F. Wheeler, G. Xue, and I. Yotov. “Coupling multipoint flux mixed finite element methods with continuous Galerkin methods for poroelasticity”. In: *Comput. Geosci.* 18 (2014), pp. 57–75.
- [201] N. Wiener. “The Homogeneous Chaos”. In: *Amer. J. Math.* 60 (1938), pp. 897–936.
- [202] P. Wriggers, W. T. Rust, and B. D. Reddy. “A virtual element method for contact”. In: *Computational Mechanics* 58.6 (2016), pp. 1039–1050. ISSN: 1432-0924. DOI: [10.1007/s00466-016-1331-x](https://doi.org/10.1007/s00466-016-1331-x).
- [203] D. Xiu and G.E. Karniadakis. “The Wiener-Askey Polynomial Chaos for Stochastic Differential Equations”. In: *SIAM J. Sci. Comp.* 24.2 (2002), pp. 619–644.
- [204] S.-Y. Yi. “Convergence analysis of a new mixed finite element method for Biot’s consolidation model”. In: *Numer. Methods Partial Differential Equations* 30.4 (2014), pp. 1189–1210. DOI: [10.1002/num.21865](https://doi.org/10.1002/num.21865).
- [205] A. Ženíšek. “The existence and uniqueness theorem in Biot’s consolidation theory”. In: *Aplikace Matematiky* 29 (1984), pp. 194–211.
- [206] R. W. Zimmerman. “Coupling in poroelasticity and thermoelasticity”. In: *Int. J. Rock Mech. Min. Sci.* 37.1–2 (Jan. 2000), pp. 79–87. ISSN: 0148-9062.



---

## Résumé

---

Dans cette thèse, on s'intéresse à de nouveaux schémas de discrétisation afin de résoudre les équations couplées de la poroélasticité et nous présentons des résultats analytiques et numériques concernant des problèmes issus de la poromécanique. Nous proposons de résoudre ces problèmes en utilisant les méthodes Hybrid High-Order (HHO), une nouvelle classe de méthodes de discrétisation polyédriques d'ordre arbitraire. Cette thèse a été conjointement financé par le Bureau de Recherches Géologiques et Minières (BRGM) et le LabEx NUMEV. Le couplage entre l'écoulement souterrain et la déformation géomécanique est un sujet de recherche crucial pour les deux institutions de cofinancement.

**Mots clés:** méthodes Hybrides d'Ordre Élevé, méthodes Galerkin discontinues, maillages polyédriques, ordre d'approximation arbitraire, poromécanique, problème de Biot, élasticité non-linéaire, écoulement de Darcy, couplage hydraulique, formulation primale, analyse fonctionnelle discrète, quantification d'incertitude.

---

## Abstract

---

In this manuscript we focus on novel discretization schemes for solving the coupled equations of poroelasticity and we present analytical and numerical results for poromechanics problems relevant to geoscience applications. We propose to solve these problems using Hybrid High-Order (HHO) methods, a new class of nonconforming high-order methods supporting general polyhedral meshes. This Ph.D. thesis was jointly funded by the Bureau de recherches géologiques et minières (BRGM) and LabEx NUMEV. The coupling between subsurface flow and geomechanical deformation is a crucial research topic for both cofunding institutions.

**Keywords:** Hybrid High-Order methods, discontinuous Galerkin methods, arbitrary order, polyhedral meshes, poromechanics, Biot problem, nonlinear elasticity, Darcy flow, hydro-mechanical coupling, primal formulation, discrete functional analysis, uncertainty quantification.