



**HAL**  
open science

# An explicit hybrid high-order method for structural dynamics

Morgane Steins

► **To cite this version:**

Morgane Steins. An explicit hybrid high-order method for structural dynamics. Dynamical Systems [math.DS]. École des Ponts ParisTech, 2023. English. NNT : 2023ENPC0039 . tel-04421265v2

**HAL Id: tel-04421265**

**<https://hal.science/tel-04421265v2>**

Submitted on 19 Mar 2024

**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.

# An explicit hybrid high-order method for structural dynamics

École doctorale N°532, Mathématiques et Sciences et Technologies de  
l'Information et de la Communication (MSTIC)

Mathématiques Appliquées

Thèse préparée au Centre d'Enseignement et de Recherche en  
Mathématiques et Calcul Scientifique (CERMICS) de l'École des Ponts  
ParisTech, au Laboratoire d'Études de Dynamique du CEA (SEMT/DYN) et au  
sein de l'équipe-projet SERENA (INRIA)

---

Thèse soutenue le 5 décembre 2023

**Morgane STEINS**

---

Composition du jury :

Paola ANTONIETTI Professeure, Politecnico di Milano	<i>Examinatrice</i>
Hélène BARUCQ Directrice de Recherche, INRIA	<i>Rapportrice</i>
Erik BURMAN Professeur, University College of London	<i>Examineur</i>
Alexandre ERN Professeur, École des Ponts ParisTech	<i>Directeur de thèse</i>
Sébastien IMPERIALE Chargé de Recherche, INRIA	<i>Rapporteur</i>
Olivier JAMOND Docteur, CEA Saclay	<i>Encadrant industriel</i>
Patrick LE TALLEC Professeur, École Polytechnique	<i>Président du jury</i>
Nicolas PIGNET Docteur, EDF R&D	<i>Examineur</i>



# Résumé

Dans cette thèse, nous nous intéressons au développement de la méthode HHO (hybride d'ordre élevé) pour la dynamique des structures. La méthode HHO est formulée en termes d'inconnues de faces définies sur le squelette du maillage, et d'inconnues de cellules définies dans le volume. La méthode HHO présente de nombreux avantages dans le cadre de la mécanique des structures : elle s'écrit sous forme primale, elle est robuste au verrouillage volumique, elle permet l'utilisation naturelle de maillages polyédriques quelconques, et elle est efficace en termes de coût de calcul. Plus particulièrement, dans cette thèse, nous nous penchons sur la possibilité de coupler la méthode HHO pour la discrétisation spatiale d'équations d'ondes linéaires et non-linéaires avec une intégration explicite en temps. Nous prouvons la convergence optimale de l'approximation dans le cas linéaire avec le schéma des différences centrées en temps. Puis, nous proposons une méthode itérative pour résoudre le couplage statique existant entre les inconnues de cellule et de face à chaque pas de temps. Cette méthode repose sur un *splitting* d'opérateurs et converge si la stabilisation est multipliée par un facteur suffisamment grand. Une preuve de convergence est donnée dans le cas linéaire. L'extension de cette méthode à l'équation des ondes non-linéaires, à l'équation des ondes élastiques et à la dynamique des structures en grandes transformations avec plasticité est ensuite présentée. De nombreux exemples numériques permettent de vérifier la précision de la méthode proposée, d'étudier l'influence des paramètres du *splitting* et de comparer la méthode aux éléments finis en termes de précision et de coût de calcul.

**Mots-clés :** Méthode HHO, schéma explicite en temps, équation des ondes, ondes non-linéaires, méthode itérative, ondes élastiques, plasticité en grandes transformations, absence de verrouillage volumique.

# Abstract

In this Thesis, we focus on the development of the Hybrid High-Order (HHO) method for structural dynamics. The HHO method is formulated in terms of face unknowns defined on the mesh skeleton, and cell unknowns defined in the volume. The HHO method has many advantages in the context of structural mechanics: it can be written in primal form, it is robust to volume locking, it accommodates polyhedral meshes, and it is computationally efficient. More specifically, in this Thesis, we investigate the possibility of coupling the HHO method for the spatial discretization of linear and nonlinear wave equations with explicit time-integration schemes. We prove the optimal convergence of the approximation in the linear case with the central finite difference scheme in time. We then propose an iterative method for solving the static coupling between the cell and face unknowns at each time step. The method is based on an operator splitting and converges if the stabilisation is multiplied by a sufficiently large factor. A convergence proof is given in the linear case. We extend this method to the nonlinear wave equation, to the elastic wave equation and to structural dynamics with finite-strain plasticity. Numerous numerical experiments are conducted to verify the accuracy of the proposed method, examine the influence of splitting parameters and compare the method to a standard finite element approach in terms of accuracy and computational cost.

**Keywords:** HHO method, explicit time scheme, wave equation, nonlinear waves, iterative method, elastic waves, finite-strain plasticity, locking-free.



# Remerciements

Ce manuscrit représente l'aboutissement de trois années de travail intense, marquant une étape cruciale dans ma vie. J'exprime ma profonde gratitude envers toutes les personnes qui ont contribué à mon parcours jusqu'ici.

En premier lieu, je tiens à remercier sincèrement mon directeur de thèse, Alexandre. Nos innombrables sessions de tableau blanc, ses relectures minutieuses de papiers et de ce manuscrit, ainsi que ses précieux conseils ont été indispensables. Un immense merci également à mes encadrants au CEA, Olivier et Florence, pour leur soutien quotidien, leur expertise dans la prise en main de MANTA, leurs interprétations judicieuses de mes nombreux bugs, et leur assistance face aux méandres administratifs.

Je souhaite exprimer ma reconnaissance envers les membres du jury qui ont accepté de participer à ma soutenance et ont manifesté un intérêt significatif pour mon travail. Mes remerciements particuliers vont à Hélène Barucq et Sébastien Impériale pour avoir assumé le rôle de rapporteurs et pour leurs retours éclairés.

Ma thèse a été co-encadrée par trois organismes, j'ai donc rencontré beaucoup de collègues. Je tiens en particulier à remercier mes co-doctorants du CEA pour le temps passé ensemble, les repas et les promenades sur le centre : Antoine, Filippo, Matthieu, Pau, Léopold, Stan. Un merci tout spécial Guillaume, mon colocataire bavard de la pièce 32, bon voyage au SGLS. Merci également à mes camarades doctorants de l'INRIA qui m'ont adoptée en cette fin de thèse : Ari, Clément, Romain, j'espère revenir grimper avec vous. Il n'y a pas que des doctorants dans nos labos, et je remercie tout le laboratoire DYN pour ces pauses café et repas où nous avons pu discuter de boulot mais surtout d'autres choses. Un grand merci en particulier à la team MANTA, qui m'a aidé à développer le code nécessaire aux résultats de cette thèse. Merci également aux membres de l'équipe projet SERENA de l'INRIA, je suis souvent passée en coup de vent, mais je me suis toujours sentie chez moi.

Enfin, la réussite de cette thèse tient aussi à l'équilibre entre le travail et le reste. Merci à mes amis qui m'ont soutenue pendant ces trois ans : Alexandre pour m'avoir montré la voie en soutenant quelque mois avant, Albane, Emma, Maëlle, Marine, depuis l'ENSTA, merci d'avoir toujours été là pour tester mes tentatives culinaires sans gluten, partager une raclette ou une partie de Seven Wonders. Merci à Bastien, Théo et Maxime d'être souvent venus en lointaine banlieue grimper à Massy. Un immense merci à ma famille pour avoir su me rappeler pourquoi je faisais une thèse dans les moments de doute, avoir écouté tant de tentatives de vulgarisation de mon sujet, et d'avoir toujours été là quand j'en avais besoin. Enfin, un grand merci à l'homme qui partage ma vie, Mathieu, pour m'avoir accompagnée chaque jour dans ce projet et dans tous les autres.



# Contents

<b>1</b>	<b>Introduction</b>	<b>7</b>
1.1	Contexte industriel . . . . .	7
1.2	Enjeux numériques . . . . .	9
1.3	La méthode HHO . . . . .	14
1.4	Objectifs de la thèse et plan du manuscrit . . . . .	18
<b>2</b>	<b>The HHO method in a nutshell</b>	<b>22</b>
2.1	Some background on the HHO method . . . . .	22
2.2	Exemple: Discretizing the Poisson problem . . . . .	23
2.3	Error analysis . . . . .	29
<b>3</b>	<b>Semi-implicit HHO method for the acoustic wave equation</b>	<b>33</b>
3.1	Optimal convergence of the HHO space semi-discretization . . . . .	34
3.2	Full discretization with leapfrog scheme . . . . .	39
3.3	Numerical experiments . . . . .	52
3.4	Other possible schemes . . . . .	54
3.5	Conclusion . . . . .	55
<b>4</b>	<b>Explicit HHO method for the acoustic wave equation</b>	<b>57</b>
4.1	Study of the operator splitting . . . . .	57
4.2	Parametric study of the splitting procedure . . . . .	64
4.3	Comparison to finite elements . . . . .	74
4.4	Conclusion on the splitting procedure . . . . .	80
<b>5</b>	<b>Explicit HHO method for the nonlinear acoustic wave equation</b>	<b>82</b>
5.1	Splitting procedure for the nonlinear wave equation . . . . .	82
5.2	Nonlinear wave equation with a p-structure potential . . . . .	86
5.3	Vibrating membrane problem . . . . .	91
5.4	Conclusions on the efficiency of the splitting procedure . . . . .	95
<b>6</b>	<b>Explicit HHO method for structural dynamics</b>	<b>96</b>
6.1	The HHO method for the elastic wave equation . . . . .	96
6.2	Splitting for the elastic wave equation . . . . .	105
6.3	Numerical experiments on the elastic wave equation . . . . .	111



6.4	Full discretization of structural dynamics with finite-strain plasticity . . . . .	120
6.5	Numerical experiments with finite-strain plasticity . . . . .	126
<b>7</b>	<b>Conclusion and perspectives</b>	<b>136</b>

# Chapter 1

## Introduction

Dans ce chapitre, nous présentons tout d'abord le contexte industriel et informatique de la thèse. Nous détaillons ensuite les principaux enjeux en simulation numérique pour la dynamique des structures : verrouillage volumique, intégration temporelle, raffinement de maillage adaptatif. Nous présentons également les principales méthodes de discrétisation spatiale utilisées et leur capacité à résoudre ces enjeux. Puis, nous introduisons la méthode hybride d'ordre élevé (*Hybrid High Order, HHO*) avec une revue bibliographique et une présentation des principaux ingrédients de cette méthode. Enfin, nous présentons le plan du manuscrit.

### 1.1 Contexte industriel

Dans de nombreux secteurs industriels, des études de scénarios accidentels sont menées afin de contrôler l'intégrité des structures mécaniques suites à un phénomène très violent (choc, explosion...). Dans le cadre de l'industrie nucléaire, des exemples de scénarios accidentels sont la chute d'avion de ligne sur le bâtiment réacteur d'une centrale nucléaire, le risque de chute de colis transportant du combustible usagé lors de son transport et de sa manutention, les accidents de type BORAX (explosion, voire fusion rapide d'une partie du cœur dans certains types de réacteurs), le risque H2 (risque d'explosion comme à Fukushima), l'Accident de Perte de Réfrigérant Primaire (APRP) ou encore l'Accident de Dimensionnement du Confinement dans le cas du prototype de réacteur à neutrons rapides comme le prototype Astrid. Lors de ces scénarios accidentels, il est impératif de s'assurer qu'aucune substance radioactive n'est émise. Ces accidents donnent lieu à d'importantes sollicitations des structures qu'il faut déterminer avec précision afin de garantir leur intégrité (enceinte de confinement, tuyauterie par exemple).

Un des objectifs du CEA est de fournir un soutien à l'industrie nucléaire, en participant aux recherches de ses partenaires (parmi lesquels on peut citer EDF, Framatome et l'IRSN) ou en lançant ses propres projets (réacteurs expérimentaux dans le domaine civil comme le réacteur à neutrons rapides Astrid). Parmi les accidents graves étudiés au CEA, l'Accident de Dimensionnement du Confinement (ADC) pour le prototype Astrid est un exemple caractéristique des enjeux de sûreté nucléaire. Il s'agit d'une explosion au sein même du réacteur qui crée une *bulle* de sodium chargée énergétiquement au centre du réacteur. Cette bulle s'étend, générant ainsi des ondes de pression. Ces ondes de pression se propagent dans le cœur et rencontrent la structure d'enceinte. Sous cette pression, les structures se déforment plastiquement, en particulier le sommet de la cuve du réacteur. Si la cuve est trop déformée, une partie du sodium peut fuir dans l'enceinte de confinement, prenant alors feu. Dans un tel scénario, la conception d'une voûte et d'une enceinte de confinement résistantes à différentes sollicitations mécaniques et thermiques est indispensable. Pour plus de détails sur l'ADC et les conséquences sur la déformation des structures, voir [58].

Les expériences réelles pouvant être difficiles, voire impossibles à réaliser dans le cas de l'industrie

nucléaire, la simulation numérique en mécanique se révèle être un outil essentiel pour le CEA dans l'étude de tels phénomènes. La simulation numérique est également employée pour la conception de prototypes dans le cadre de projets. En effet, le CEA, qui exploite un certain nombre d'installations nucléaires de recherche, doit en garantir la sûreté.

Le CEA est doté d'un certain nombre de logiciels de simulations, développés en son sein, afin de répondre au besoin croissant de simulation numérique. Dans le domaine de la dynamique des structures, le logiciel actuellement utilisé au CEA est le logiciel Europlexus, avec lequel des simulations d'ADC sont réalisées, voir par exemple [102]. Pour le CEA, les enjeux de la simulation numérique sont de fournir des simulations les plus fiables et précises possible. Les logiciels actuels permettent de faire des simulations crédibles, mais sont limités en performance par certaines contraintes technologiques. En effet, Europlexus est un logiciel dont le développement a commencé il y a une quarantaine d'années. Les architectures de calcul moderne ont depuis largement évolué et il est maintenant difficile d'adapter le logiciel à ces nouvelles architectures. Le passage au calcul massivement parallèle et la versatilité des logiciels sont par exemple des enjeux majeurs pour le CEA. Ils permettent en effet de produire des simulations de plus grande précision en des temps plus courts, facilitant ainsi les calculs multiphysiques. Ce dernier aspect est essentiel pour les simulations accidentelles. En effet, dans le cas de l'ADC, il faut être capable de coupler la résolution du problème de neutronique, de thermique, de mécanique des structures et de mécanique des fluides pour simuler complètement l'accident.

Afin de moderniser ses logiciels de simulation existants, le Service d'Études de Mécanique et de Thermique (SEMT) du CEA développe le code de simulation nouvelle génération MANTA (*Mechanical Analysis Numerical Toolbox for advanced Applications*). Il s'agit d'un code C++ *open-source*, tourné vers le calcul haute performance, qui doit permettre d'augmenter les capacités de simulation numérique des logiciels existants (en particulier, Europlexus pour la dynamique rapide et Cast3m pour la mécanique quasi-statique). MANTA est un code nativement parallèle, ce qui permet d'utiliser les architectures des calculateurs modernes (aujourd'hui, le calcul CPU massivement parallèle et, à terme, le calcul GPU également). Un des objectifs de MANTA est de permettre une utilisation facile de méthodes de simulation numérique modernes et performantes. Un autre enjeu est de pouvoir faire des simulations mettant en jeu de nombreuses physiques différentes : dynamique des structures, dynamique des fluides, thermodynamique par exemple, et leur interaction. Le logiciel MANTA a une visée industrielle, puisqu'il remplacera les logiciels actuellement utilisés dans les collaborations du CEA à l'horizon 2030.

Dans le cadre de la simulation en dynamique des structures, la méthode la plus couramment utilisée est celle des éléments finis. Cette méthode permet de résoudre la plupart des équations de la physique et est aisée à mettre en œuvre. Cependant, lors de fortes déformations des structures, les simulations éléments finis peuvent produire dans certains cas une solution non-physique qui sous-estime considérablement les déformations réellement observées. Ce phénomène est appelé verrouillage volumique et apparaît plus précisément lors de fortes déformations plastiques ou de déformations élastiques incompressibles. La solution obtenue n'est alors pas utilisable dans le cadre d'études de sûreté.

Un autre désavantage de la méthode des éléments finis se manifeste lors d'un raffinement de maillage adaptatif (AMR pour *Adaptive Mesh Refinement*). Il s'agit d'un processus qui permet d'augmenter la précision de la simulation uniquement en certaines zones d'intérêt de la structure en raffinant individuellement et récursivement les cellules d'un maillage initial. La méthode des éléments finis n'est pas optimale pour l'AMR, car elle nécessite de propager le raffinement au-delà des zones où il est strictement nécessaire (voir section 1.2.2 pour plus de détails). Il existe des solutions pour utiliser du raffinement de maillage adaptatif avec les éléments finis, mais celles-ci ont tendance à avoir un impact négatif sur la performance ou la précision des calculs (même si des résultats théoriques permettent de tempérer cet impact dans le régime asymptotique [145]).

L'objectif du CEA étant de produire les simulations les plus fiables et performantes possible (en termes de temps de calcul), la recherche d'une méthode de simulation alternative à la méthode des éléments finis est ainsi justifiée. Or, il se trouve que la méthode HHO permet une gestion naturelle de l'AMR et n'est pas sensible au problème de verrouillage volumique. L'enjeu de cette thèse est d'évaluer l'utilisation de la méthode HHO pour la simulation numérique en dynamique des structures et notamment sa compétitivité

face à la méthode des éléments finis en termes de précision et de temps de calcul.

## 1.2 Enjeux numériques

La performance et la fiabilité des logiciels de simulation numérique est essentielle. Nous détaillons ici trois enjeux numériques impactant les performances et la fiabilité des simulations en dynamique des structures.

### 1.2.1 Intégration explicite en temps

Un premier aspect crucial des simulations en dynamique des structures est le choix de l'intégration temporelle, puisque celle-ci impacte aussi bien la performance des simulations que leur fiabilité.

On peut classer les méthodes d'intégration en temps en deux catégories : les méthodes implicites et les méthodes explicites. Ce que désignent les termes explicites et implicites peut varier selon le contexte : analyse mathématique, vision physique ou encore aspects informatiques. Nous choisissons ici la distinction suivante : toutes les méthodes s'écrivent sous la forme  $g(U^n) = f((U^k)_{1 \leq k \leq n-1})$ , où  $U$  est l'inconnue,  $n$  le pas de temps discret et  $f, g$  des fonctions arbitraires. Une méthode est considérée explicite si le coût pour inverser la fonction  $g$  est du même ordre de grandeur que le coût pour calculer  $f(U^k)$ . Ainsi, si  $g$  s'écrit comme une matrice diagonale ou diagonale par bloc, la méthode est considérée explicite. Au contraire, si  $g$  est une fonction linéaire avec une matrice non diagonale par bloc, ou une fonction non linéaire, le problème est considéré comme implicite, car cela revient à résoudre un système linéaire, et pour les équations non-linéaires, à trouver un zéro d'une fonction. Le surcoût pour un pas de temps implicite est souvent compensé par la possibilité de pouvoir faire des pas de temps aussi grands que souhaités car les schémas implicites sont en général inconditionnellement stables. Ainsi, le pas de temps n'est contraint que par la précision numérique recherchée. En revanche, les méthodes explicites sont souvent contraintes par une relation de stabilité CFL de la forme  $\Delta t \leq c_{\text{CFL}} h^\alpha$  où  $h$  représente une longueur caractéristique des cellules du maillage spatial, typiquement leur diamètre et  $\alpha = 1$  pour les problèmes hyperboliques et  $\alpha = 2$  pour les problèmes paraboliques. Cette condition impose donc une réduction du pas de temps lorsque le maillage est raffiné. Ceci peut donc conduire à effectuer de nombreuses itérations en temps lors de simulations avec un grand nombre de degrés de liberté. Lors de simulations éléments finis avec un schéma d'intégration aux différences finies centrées explicites pour la dynamique des structures, on obtient une équation algébrique de la forme  $\mathcal{M}U^n = f((U^k)_{1 \leq k \leq n-1})$ . Résoudre ce problème demanderait donc d'inverser la matrice  $\mathcal{M}$ . Pour obtenir un schéma complètement explicite, la matrice  $\mathcal{M}$  est diagonalisée, i.e., remplacée par une matrice  $\mathcal{M}'$  diagonale. Son inversion est alors triviale et le calcul des inconnues  $U^n$  ne demande que l'évaluation de la fonction  $f$ . Ces méthodes de diagonalisation, dites de *mass lumping*, sont classiquement mises en œuvre à l'ordre 1 [62, 115, 104, 103], et des extensions aux ordres plus élevés existent, voir par exemple [74, 73] pour l'application à l'équation des ondes.

Le choix de la méthode d'intégration temporelle des équations de la dynamique dépend du type de problème et de réponse étudiés et de la régularité des données. Il est possible de classer les types de problèmes de dynamique en deux catégories : les problèmes de propagation d'ondes et les problèmes inertiels [24]. Dans ce deuxième cas, les ondes présentes dans la solution sont de basse fréquence (par rapport au temps d'observation) et facilement représentées par une discrétisation grossière (grands pas de temps), si bien que les méthodes d'intégration implicite sont préférées. En effet, elles permettent de ne faire que le nombre de pas de temps nécessaires afin de représenter les fréquences temporelles de la solution sans avoir à en faire davantage pour respecter une condition de stabilité CFL. Au contraire, les cas de propagation d'ondes désignent les problèmes pour lesquels la fréquence des ondes présentes dans la solution est très élevée (par rapport au temps d'observation). Pour représenter ces ondes, le pas de temps requis est très petit, si bien que les méthodes d'intégration explicite sont souvent préférées. En effet, la condition de stabilité CFL n'est pas le seul facteur incitant à prendre un pas de temps faible. Comme la résolution d'un pas de temps explicite est moins coûteuse que celle d'un pas de temps implicite, à pas de temps fixé, les méthodes explicites sont plus efficaces.

Les méthodes explicites sont également favorisées lorsque les données (chargement, conditions aux limites, propriétés matériaux) sont très peu régulières, par exemple lors de déformations plastiques fortement non-linéaires, de très grandes déformations ou de problèmes de contact. Les méthodes implicites manquent alors de robustesse et de fiabilité, car la convergence du solveur linéaire peut devenir très délicate. Voir par exemple [146] pour une comparaison des méthodes d'intégration implicite et explicite en dynamique des structures.

Dans le cadre de simulations de scénarios accidentels, les temps caractéristiques des déformations sont très rapides et les non-linéarités, notamment dues aux grandes déformations plastiques, sont très fortes. Il est également fréquent que la résolution temporelle souhaitée soit très précise, nécessitant un grand nombre de pas de temps. Ainsi, Europlexus et MANTA utilisent des méthodes d'intégration explicite pour les simulations en dynamique rapide.

### 1.2.2 Raffinement de maillage adaptatif

Le raffinement de maillage adaptatif est une méthode permettant d'optimiser les ressources de calcul. Il répond donc à un objectif de performance. Les calculs en simulation numérique requièrent une discrétisation spatiale de la structure étudiée par un maillage. Par défaut, ce maillage peut être raffiné de la même façon partout, c'est-à-dire que les cellules ont à peu près la même taille dans toute la structure. La taille des cellules peut également être contrainte par la géométrie afin de représenter la structure correctement. La taille des cellules impacte également la qualité de la solution. Plus le maillage est fin (plus les cellules sont petites), plus la solution obtenue est précise, mais plus le calcul est coûteux. Il est très fréquent que la zone d'intérêt de la simulation ne soit composée que d'une sous-partie de la structure étudiée. Il est alors naturel de souhaiter une plus grande précision dans les zones d'intérêt. Raffiner tout le maillage uniformément représenterait un coût important, pour un bénéfice faible dans certaines zones. Il est alors plus économe de raffiner le maillage seulement dans les zones d'intérêt.

Les méthodes de raffinement de maillage adaptatif permettent d'estimer au cours de la simulation quelles sont les zones nécessitant une plus grande précision et de raffiner le maillage localement en fonction de cette estimation. Nous ne rentrons pas dans le détail de l'estimation dans cette thèse, car le raffinement de maillage adaptatif n'y est pas exploité. Mais il s'agit d'un des outils majeurs pour la performance des logiciels de simulation numérique modernes. Les méthodes développées doivent donc être compatibles avec l'AMR.

Il est également important de classer les types de raffinements de maillages locaux en deux catégories : conforme et non-conforme. Le raffinement conforme évite l'apparition de *hanging nodes* (nœuds orphelins), c'est-à-dire de points du maillage au milieu d'une arête ou d'une face en 3D. Une illustration des deux types de raffinement de maillages est présentée dans la figure 1.1. On souhaite raffiner localement le maillage de la figure 1.1a autour du coin inférieur gauche. Un premier exemple de raffinement conforme, illustré sur la figure 1.1b, consiste à remailler complètement le domaine. Un deuxième exemple, illustré à la figure 1.1c, consiste à gérer localement la présence de nœuds orphelins en raffinant les mailles voisines. Le raffinement non-conforme, illustré sur la figure 1.1d, raffine seulement une partie des mailles, mais entraîne l'apparition des nœuds orphelins matérialisés par des points rouges sur la figure. L'avantage du raffinement non-conforme est qu'il conserve une hiérarchie entre le maillage initial et le maillage raffiné. Ceci permet donc une projection relativement simple des champs définis sur le maillage initial sur le maillage raffiné.

Les éléments finis sont compatibles avec le raffinement local conforme. En revanche, lors de l'utilisation de raffinement non-conforme, les nœuds orphelins demandent un traitement particulier. Ce n'est pas le cas pour la méthode HHO, dont le caractère non-conforme la rend compatible avec tout type de maillages. Ceci permet donc une plus grande souplesse dans le choix de la méthode d'AMR, ce qui fournit un argument supplémentaire pour l'étude de la méthode HHO dans cette thèse.

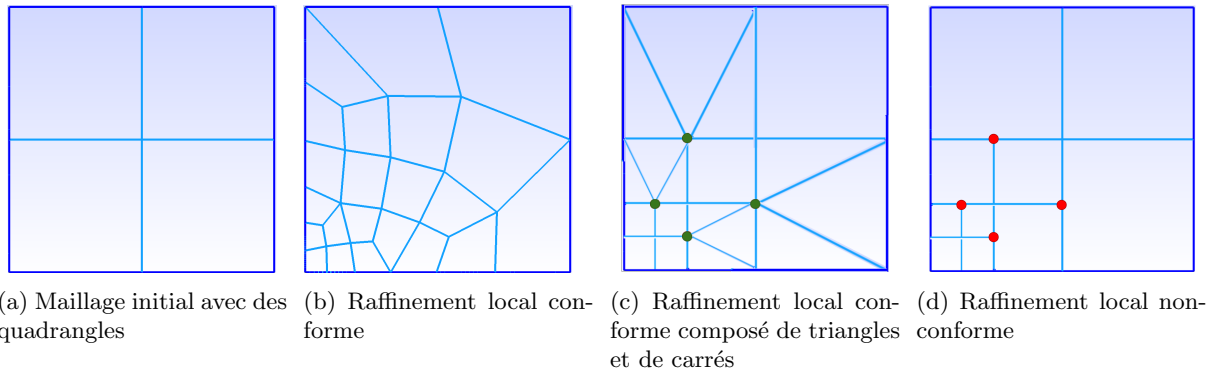


Figure 1.1: Raffinement de maillage adaptatif - Illustration des types de raffinement locaux (conforme et non-conforme) sur un maillage de quadrangles.

### 1.2.3 Le verrouillage volumique : un enjeu de fiabilité des simulations

Nous avons évoqué précédemment le phénomène du verrouillage volumique, qui nuit à la fiabilité des simulations. Nous décrivons ici plus précisément ce problème, avant de détailler les méthodes existantes permettant de le résoudre.

#### 1.2.3.1 Description du phénomène

Le verrouillage volumique est un phénomène pouvant apparaître lors de simulations numériques en mécanique avec des éléments  $H^1$ -conformes en présence de déformations élastiques ou plastiques incompressibles. Il s'agit d'un phénomène générant une oscillation du tenseur des contraintes aux points d'intégration, et potentiellement un déplacement ne convergeant pas vers la bonne solution lorsque le maillage est raffiné. Les conséquences de ce phénomène sont multiples. Tout d'abord, il existe des cas où la solution ne converge pas vers la solution théorique ou expérimentale. La simulation est donc grossièrement fautive, quel que soit le maillage utilisé. Dans d'autres cas, la solution en déplacement n'est pas impactée, mais le tenseur des contraintes aux points d'intégration oscille fortement. Des algorithmes basés sur la valeur de ce tenseur comme les algorithmes d'endommagement sont alors affectés.

Une des raisons expliquant l'apparition de verrouillage volumique est que les éléments finis  $H^1$ -conformes d'ordre bas ne sont pas assez riches pour approcher correctement les champs de vecteurs à divergence nulle (il n'y a pas assez de degrés de liberté dans l'espace d'approximation pour de tels champs) [137]. Le lien entre l'approximation d'un espace à divergence nulle avec le verrouillage volumique est détaillé dans [148]. Une étude mathématique du phénomène de verrouillage volumique avec de nombreux exemples est présentée dans [7]. Cet article relie la dépendance de la constante de stabilité inf-sup au coefficient de Poisson, la dégradation de cette constante impliquant le manque de robustesse à l'incompressibilité. Plus précisément, cette constante peut être bornée inférieurement indépendamment de la taille des mailles et du coefficient de Poisson uniquement pour les éléments finis d'ordre élevé (ordre polynomial  $p \geq 4$ ). Pour les ordres plus bas, cette constante se dégrade lorsque le coefficient de Poisson est proche de 0.5. L'étude de la robustesse de l'ordre élevé des éléments finis est conduite dans [16, 17], où sont également présentés des exemples numériques de verrouillage volumique.

Une illustration minimale de ce phénomène peut être faite en considérant une structure carrée pleine, maillée par deux triangles rectangles, comme les maillages illustrés sur les figures 1.2a et 1.2b. Le carré est fixé au bâti sur les segments  $[AB]$  et  $[AC]$ . On considère des déformations incompressibles, c'est-à-dire que l'aire de chaque triangle doit être conservée. Dans la configuration de la figure 1.2a, si on applique une force sur le point  $D$ , l'élément 1 ne peut pas se déformer, mais l'élément 2 peut se déformer, tant que la distance entre le point  $D$  et le segment  $[BC]$  est conservée. Dans la configuration de la figure 1.2b, il faut que les aires des deux éléments soient conservées, c'est-à-dire que la distance des points  $B$  et  $C$  au

segment  $[AD]$  soit préservée. On voit ainsi que le point  $D$  ne peut être déplacé, quelle que soit la force appliquée. Il s'agit donc bien d'un cas de verrouillage, car sous une force non nulle, la structure doit se déformer. Ce verrouillage perdure lorsque le maillage est raffiné comme sur la figure 1.2c. En effet, en utilisant le raisonnement précédent, les points  $E$  et  $G$  sont immobiles. Le raisonnement s'étend donc aux quadrilatères inférieur gauche. De la même façon, les points  $E$ ,  $H$  et  $B$  étant fixés, le même raisonnement empêche le déplacement du point  $F$ . Finalement, le point  $D$  est également bloqué par les points  $E$ ,  $F$  et  $G$ . Ainsi, quel que soit le raffinement du maillage de ce carré, aucun déplacement n'est possible.

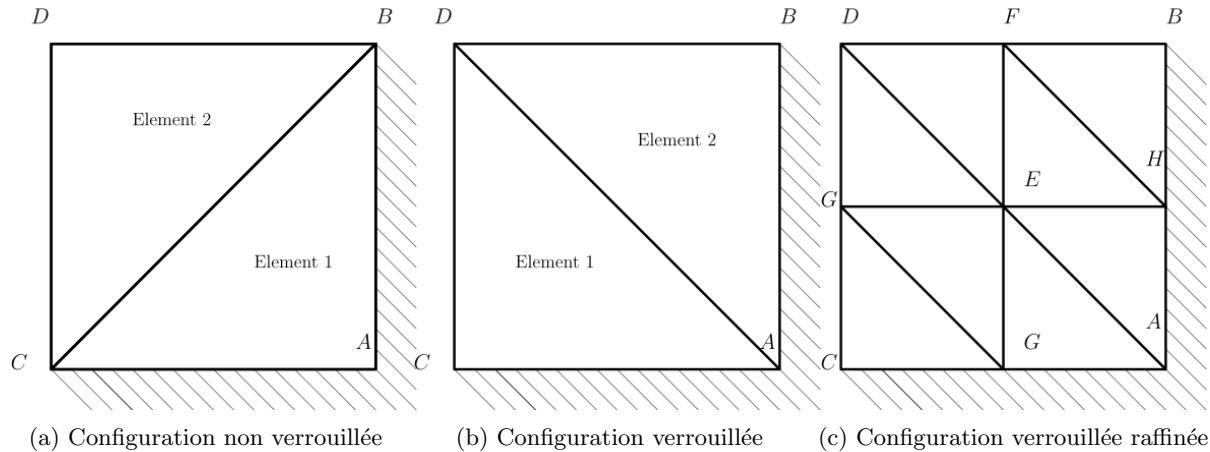


Figure 1.2: Configurations possibles pour mailler une structure carrée avec des triangles rectangles : démonstration du phénomène de verrouillage volumique.

Le phénomène de verrouillage volumique se retrouve lors de simulations éléments finis plus réalistes, voir par exemple [24, § 8.4.3] pour une série d'exemples. Un cas classique de verrouillage volumique en élasticité linéaire est celui de la membrane de Cook élastique quasi-incompressible. On considère la structure de la figure 1.3a en déformations planes, dont la partie gauche est fixe. On applique une traction verticale  $F_y$  sur la partie droite du panneau, on considère des déformations élastiques, quasi-incompressibles (coefficient de Poisson  $\nu = 0.49$ ) et quasi-statiques. La figure 1.3b montre la trace du tenseur des déformations de Cauchy sur cette membrane, grandeur appelée pression hydrostatique. Cette pression présente de fortes oscillations entre deux points d'intégration d'une même cellule. Il s'agit de la matérialisation du phénomène de verrouillage volumique en éléments finis. Lorsque le maillage est raffiné, le déplacement du point  $A$  converge vers la valeur de référence, mais les oscillations de pression en forme de damier persistent.

### 1.2.3.2 État de l'art des méthodes de résolution du verrouillage volumique

Comme mentionné précédemment, la méthode des éléments finis est la méthode classiquement utilisée en simulation numérique pour la mécanique des structures. Parmi les améliorations de cette méthode conçues pour résoudre le verrouillage volumique, nous évoquons les suivantes :

- **Méthodes mixtes** : ces méthodes ont été introduites en élasto-plasticité dans [142] et développées par la suite notamment pour les petites déformations plastiques dans [47, 48, 49, 59, 60] et pour les grandes déformations plastiques dans [123, 4, 9]. Elles consistent à introduire des variables supplémentaires pour imposer la condition d'incompressibilité de manière faible, relâchant ainsi la contrainte sur le déplacement. On introduit couramment la variable supplémentaire  $p$  décrivant la contrainte moyenne ou pression hydrostatique. Dans les cas de lois de comportement complexes couplant les composantes déviatoriques et moyennes des contraintes, on ajoute également une variable décrivant le changement de volume. Les inconnues peuvent être définies aux points d'intégration

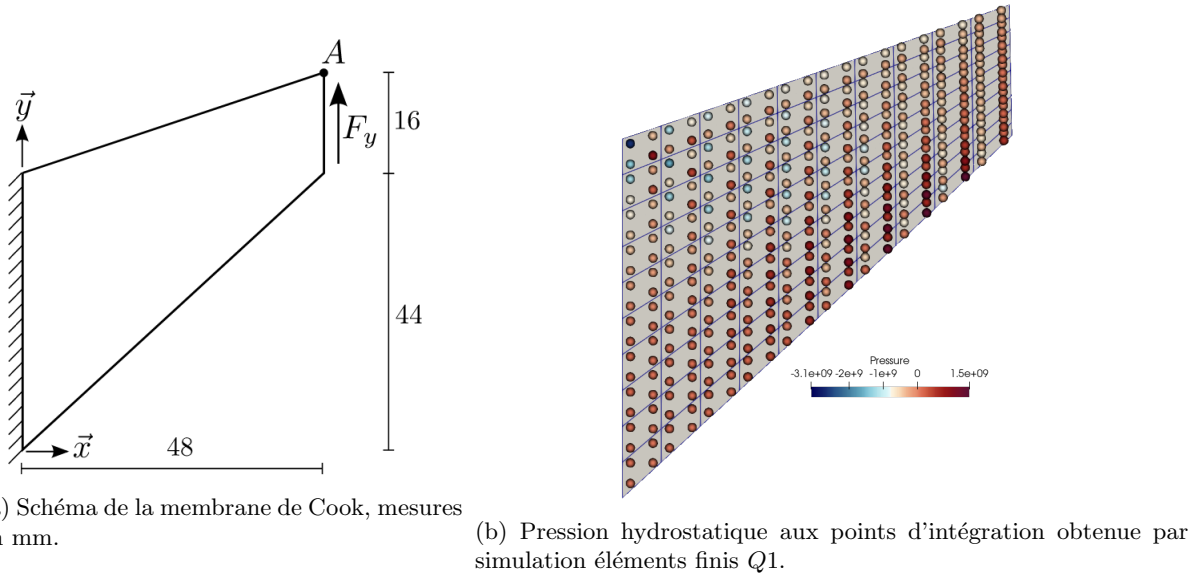


Figure 1.3: Membrane de Cook élastique quasi-incompressible : illustration du phénomène de verrouillage volumique.

ou au niveau des cellules du maillage. Bien que l'on ajoute des inconnues supplémentaires, il peut être possible de résoudre une partie des équations sur chacune des cellules, limitant alors le surcoût de l'approche mixte. Un des désavantages est toutefois que la formulation est dépendante de la loi de comportement choisie. Ceci demande donc une implémentation différenciée selon le comportement des matériaux. De plus, les méthodes mixtes sont difficilement extensibles aux domaines polyédriques, et demandent une stabilisation pour les discrétisations à l'ordre le plus bas [147].

- **Intégration réduite ou sélective** : Les méthodes d'intégration réduite [152] consistent à utiliser une quadrature d'ordre plus faible que l'ordre nécessaire pour intégrer de façon exacte les fonctions de forme. Les méthodes d'intégration sélective [91] utilisent le même procédé, mais seulement sur une partie du tenseur des contraintes. Celui-ci est décomposé en deux, la composante déviatorique est intégrée complètement, et la composante hydrostatique est sous-intégrée. Ceci induit donc une approximation dans le calcul des intégrales sur les cellules du maillage, ce qui a pour effet de relâcher la contrainte d'incompressibilité locale. Ces deux méthodes ont été rapprochées des méthodes mixtes dans [129], mais les méthodes d'intégration réduite ou sélective conservent une formulation uniquement sur le déplacement. Ces méthodes permettent de résoudre le problème de verrouillage volumique, mais elles peuvent nécessiter une pénalisation de certains modes d'énergie nulle (dits *hour-glass*), voir par exemple [135] pour une application à la plasticité en grandes déformations. De plus, leur extension à des matériaux anisotropes peut se révéler difficile [114].
- **Méthodes Enhanced Assumed Strain (EAS)** : Ces méthodes, présentées dans [141], consistent à enrichir le tenseur des déformations. Par exemple, dans le cas de l'élasticité non-linéaire, on définit  $\epsilon := \nabla^s \mathbf{u} + \tilde{\epsilon}$ , où le gradient symétrisé  $\nabla^s \mathbf{u}$  est la mesure de déformation classiquement utilisée, appelée déformation admissible, et  $\tilde{\epsilon}$  est l'augmentation, dont la discrétisation peut être discontinue entre les cellules du maillage. Ceci permet notamment d'éliminer le verrouillage volumique pour les éléments finis d'ordre bas, en concentrant la contrainte d'incompressibilité sur la nouvelle variable  $\tilde{\epsilon}$ . On observe cependant l'apparition de modes *hour-glass* pour les grandes déformations. Ces modes peuvent être éliminés par stabilisation, par modification de la méthode d'enrichissement [105] ou par association de la méthode d'enrichissement aux méthodes mixtes présentées précédemment [118, 119]. L'application de la méthode EAS aux grandes déformations élastiques avec une stabilisation



optimisée est présentée dans [121, 122] et son application aux grandes déformations plastiques est présentée dans [140, 139].

Il existe d'autres méthodes de discrétisations spatiales que la méthode des éléments finis. Certaines de ces méthodes résolvent le problème du verrouillage volumique. Il s'agit souvent d'une conséquence des propriétés de ces méthodes comme un ordre élevé ou la non-conformité. Parmi les méthodes de discrétisation alternatives aux éléments finis, on peut citer les suivantes :

- **Méthodes isogéométriques (IGA)** : Ces méthodes consistent à utiliser des fonctions NURBS (Non-Uniform Rational B-Splines) d'ordre élevé. La discrétisation spatiale se fonde sur des fonctions plus régulières que les polynômes par morceaux utilisés en éléments finis. Ceci améliore la précision de la méthode, mais étend son *stencil*. En effet, le support des fonctions de base s'étend à des *patches*, c'est-à-dire à des sous-domaines composés de plusieurs cellules du maillage. L'analyse mathématique des méthodes IGA peut être trouvée dans [20], notamment les propriétés de stabilité et les ordres de convergence. Ces méthodes ont été introduites pour la plasticité en petites déformations dans [95] et étendues aux grandes déformations dans [96]. Leur implémentation dans des logiciels de calcul conçus pour les éléments finis peut s'avérer complexe du fait du couplage entre la géométrie et la discrétisation du domaine.
- **Méthodes de type éléments virtuels (Virtual Element Method, VEM)** : Les méthodes VEM [23] sont une extension des éléments finis aux maillages généraux (polygonaux et polyédriques) qui repose sur l'introduction de fonctions de formes non-polynomiales, mais qui ne sont jamais explicitement calculées (d'où le qualificatif virtuel). Ces fonctions de forme sont typiquement des solutions d'EDP locales dans chaque cellule du maillage avec des données et des conditions aux limites polynomiales. Ces méthodes VEM reposent sur une formulation primale et sont robustes au verrouillage volumique, voir [21] pour l'application à des cas d'élasticité linéaire quasi-incompressible. L'application aux petites déformations plastiques est présentée dans [22] et celle aux grandes déformations dans [151, 113].
- **Méthodes de Galerkin discontinu (dG)** [14, 67] : Contrairement aux méthodes éléments finis classiques ( $H^1$ -conformes) et aux autres méthodes évoquées ci-dessus, les méthodes dG reposent sur une approximation polynomiale par morceaux. La solution discrète est donc en général discontinue aux interfaces entre les cellules du maillage. Les méthodes dG conservent la formulation primale des équations. Elles sont adaptées à l'utilisation de maillages polygonaux et polyédriques. Quelques exemples récents d'application à des problèmes de propagation d'ondes sont [12, 10, 11, 29]. Dans le cadre de la dynamique des structures, les méthodes DG génèrent une matrice de masse bloc-diagonale, ce qui rend les schémas explicites en temps faciles à mettre en place. Cependant, dans le cas de la plasticité, la loi de comportement doit être évaluée sur les faces du maillage en plus des cellules, et la matrice de rigidité n'est plus symétrique. L'application aux déformations élastiques quasi-incompressibles et incompressibles est présentée dans [109]. Pour les grandes déformations plastiques, on peut se référer à [127, 128]. Enfin, une variante prometteuse de la méthode DG est la méthode Trefftz-DG espace-temps, où les fonctions de base sont des solutions locales de l'équation des ondes. On pourra consulter [19] pour un exemple.

Une autre classe de méthodes permet de résoudre le problème de verrouillage volumique. Il s'agit de la famille des méthodes hybrides à laquelle appartient la méthode HHO, que nous allons maintenant présenter plus en détail.

### 1.3 La méthode HHO

La méthode HHO est une méthode qui permet l'utilisation d'AMR non-conforme, et qui est robuste au verrouillage volumique, ce qui la rend particulièrement attractive pour les simulations en dynamique des

structures. La méthode HHO offre en outre de bonnes performances computationnelles et des propriétés de conservation locale (sur chaque cellule du maillage). Nous présentons ici la méthode HHO, ses liens avec les autres méthodes de la littérature, ainsi que les ingrédients principaux de sa conception.

### 1.3.1 Historique et lien avec les autres méthodes

La méthode HHO a été introduite pour la diffusion linéaire dans [84] et pour l'élasticité linéaire dans [82]. Elle a depuis été étendue à de nombreux autres champs d'application. Dans le cadre des problèmes linéaires, des travaux ont depuis eu lieu sur la diffusion linéaire à coefficients variables [6, 83], le problème de Stokes [86], l'advection-diffusion-réaction [79], le problème de Cahn–Hilliard [57], et les écoulements en milieux poreux fracturés [55, 56]. En ce qui concerne les problèmes non-linéaires, des travaux ont été menés sur le problème de Leray–Lions [76, 77] et de minimisation convexe [43, 42], sur les équations de Navier–Stokes stationnaires [30, 87, 46], sur les problèmes spectraux [38, 41], et sur les fluides de Bingham [44]. Plus particulièrement en mécanique, la méthode HHO a été appliquée au problème de Biot [26], à l'élasticité non-linéaire [31], à l'hyperélasticité [1], à l'élastoplasticité [2, 3], aux équations de plaques de Kirchhoff–Love [27], aux problèmes de contact [61] avec frottement de Tresca et à l'opérateur biharmonique [92, 93]. En outre, la méthode HHO a également été utilisée avec des méthodes de domaines fictifs [37], des méthodes multi-échelles [65] et pour des problèmes d'interface [45]. L'implémentation de la méthode HHO est décrite dans [63]. Pour finir, des revues de la méthode HHO sont disponibles dans [85, 88] et deux livres ont été récemment consacrés au sujet [66, 78].

Tout comme les méthodes de dG, la méthode HHO est une méthode non-conforme. Cependant, contrairement aux méthodes dG, la méthode HHO n'est pas formulée en termes de flux entre les cellules, mais en termes d'inconnues attachées aux faces du maillage, d'où le terme *hybride*. La méthode HHO utilise également, comme les méthodes DG, des degrés de liberté attachés aux cellules du maillage. Par la suite, nous noterons  $l$  le degré polynomial des inconnues de cellule et  $k$  celui des inconnues de face. La méthode HHO est dite de même ordre (equal-order,  $l = k$ ) si les deux inconnues ont le même degré polynomial, et d'ordre mixte (mixed-order,  $l = k + 1$ ) si le degré des inconnues des cellules est supérieur d'un ordre à celui des inconnues des faces. Un point pratique important est que les degrés de liberté de cellule peuvent être éliminés localement par condensation statique (au niveau algébrique, cela revient à invoquer un complément de Schur). Le terme *ordre élevé* vient de la facilité avec laquelle la méthode HHO peut être utilisée avec des ordres polynomiaux arbitraires. La formulation de la méthode HHO repose sur deux ingrédients essentiels : (i) un opérateur de reconstruction locale du gradient à partir des inconnues locales sur la cellule et ses faces ; (ii) un opérateur de stabilisation qui pénalise au sens des moindres carrés la différence entre la trace des inconnues de cellule et les inconnues de face. Plus de détails sont donnés à la section 1.3.2 ainsi que dans les chapitres suivants.

La méthode HHO appartient à la classe des méthodes de Galerkin discontinues hybridisables (*Hybridizable Discontinuous Galerkin*, HDG) [70], comme le montre [71]. La méthode HDG approche un triplet composé de la paire d'inconnues utilisées dans la méthode HHO et de polynômes par morceaux à valeurs vectorielles pour la variable duale (typiquement, le gradient). Dans la méthode HDG, la variable duale discrète peut être exprimée localement en termes des deux autres variables, et cette formule correspond à la reconstruction du gradient HHO. En outre, la trace numérique du flux, qui est l'un des ingrédients principaux de la méthode HDG, peut être explicitement liée à la composante normale du gradient reconstruit par HHO et à la stabilisation de la méthode HHO. Ainsi, la méthode HHO présente les mêmes avantages que les méthodes HDG comme l'utilisation naturelle de maillages polyédriques, tout en proposant une formulation primale. Par ailleurs, la méthode HHO partage les mêmes principes de conception que les méthodes de Galerkin faibles (*Weak Galerkin*, WG) [149]. Le gradient reconstruit dans la méthode HHO est appelé gradient faible dans la méthode WG. La stabilisation HHO de même ordre n'a pas encore été considérée dans la méthode WG, où la stabilisation est soit une simple pénalisation aux moindres carrés. Dans le cas d'ordre mixte, la stabilisation HHO utilise la projection  $L^2$  des inconnues de cellule sur des polynômes d'ordre  $k$  sur les faces, comme dans la stabilisation HDG de Lehrenfeld–Schöberl [124, 125]. Notons que cette dernière conduit à des estimations d'erreur optimales,

alors qu'un ordre de convergence est en général perdu avec la simple stabilisation par moindres carrés sans projection. Une comparaison détaillée des méthodes HHO et WG pour le problème biharmonique peut être trouvée dans [93]. Comme indiqué dans [71, 80, 66], la méthode HHO avec le choix  $l = k - 1$ ,  $k \geq 1$ , est étroitement liée à la méthode des éléments virtuels non-conformes (ncVEM) ; voir aussi [54] dans le contexte des problèmes multi-échelles. Enfin, en ce qui concerne l'application à la mécanique, la méthode HHO présente un avantage supplémentaire qui est la robustesse face au problème de verrouillage volumique tout en conservant une formulation primale, voir [82, Rem. 9] pour l'estimation d'erreur sans verrouillage pour l'élasticité linéaire. Dans le cas de la mécanique statique, cet avantage a été exploité par exemple dans [31, 1, 2, 3].

Dans ce manuscrit, nous allons étendre la méthode HHO à des équations dynamiques traitées par intégration temporelle explicite : l'équation des ondes acoustiques et l'équation de la dynamique des structures. L'intégration temporelle sera faite par le schéma des différences finies centrées, qui est un schéma explicite du second ordre, conditionnellement stable, et très populaire pour discrétiser en temps l'équation des ondes combinée avec une discrétisation par éléments finis ou la méthode dG en espace. La discrétisation spatiale de l'équation des ondes acoustiques linéaires par la méthode HHO a été conçue dans [34] et une version de cette méthode sur maillages immergés est introduite dans [35]. L'analyse de convergence spatiale dans le cas continu en temps est effectuée dans [72] pour la méthode HDG et dans le chapitre 3 de cette thèse pour la méthode HHO (voir également [36]). Les schémas entièrement discrets avec des méthodes hybrides en espace n'ont été considérés jusqu'à présent que dans le contexte des schémas implicites en temps, comme le schéma de Störmer–Numerov du quatrième ordre combiné à la méthode HDG dans [69] et la formule de différenciation implicite du second ordre avec le schéma WG dans [112]. Dans cette thèse, nous développons un schéma de marche en temps explicite combiné à la méthode HHO, en effectuons l'analyse de convergence, et le testons sur de nombreux cas. En raison des liens évoqués ci-dessus, les résultats de cette thèse s'étendent aux méthodes HDG et WG (ainsi que ncVEM).

**Remarque 1.3.1** (Formulation à l'ordre un). Il est également possible de considérer la formulation à l'ordre un en temps de l'équation des ondes. Dans ce cas, on peut considérer soit la formulation de type hamiltonien dans laquelle les inconnues sont la variable primale et sa vitesse, soit la formulation mixte (ou de type Friedrichs) dans laquelle les inconnues sont la vitesse et la variable duale (typiquement, le gradient de la variable primale). Ces deux formulations conduisent, après une semi-discrétisation en espace, à un système d'EDO couplées du premier ordre qui peuvent être discrétisées à l'aide, par exemple, de méthodes de Runge–Kutta (RK) sous différentes formes. La formulation de type hamiltonien est examinée dans le contexte des méthodes dG dans [13], des méthodes HDG dans [136], et des méthodes WG (avec le schéma de Crank–Nicolson) dans [116]. Pour la formulation mixte, nous mentionnons les méthodes HDG avec des schémas RK implicites [133, 132] ou explicites [143] et la méthode HHO avec des schémas RK implicites ou explicites [34]. Dans cette thèse, nous nous focaliserons sur la formulation d'ordre deux en temps de l'équation des ondes.

### 1.3.2 Principaux ingrédients de la méthode HHO

La définition mathématique des inconnues et des opérateurs spécifiques à la méthode HHO sera détaillée dans le chapitre 2. Nous donnons ici une illustration du caractère hybride de la méthode par une visualisation des degrés de liberté. Nous esquissons brièvement les principaux ingrédients de la méthode, et indiquons les ordres de convergences spatiaux attendus. Enfin, nous présentons l'allure des solutions HHO pour un problème de Poisson afin d'illustrer le caractère non-conforme de la méthode (la solution pouvant être discontinue aux interfaces entre les cellules du maillage).

**Inconnues HHO.** Une des caractéristiques de la méthode HHO est la définition des inconnues. Ce sont des polynômes attachés aux cellules et aux faces du maillage. La figure 1.4 illustre les degrés de liberté HHO d'ordre  $(l, k) = (0,0)$ ,  $(1,1)$ ,  $(1,0)$  et  $(2,1)$  sur des quadrangles. Le sigle DdL signifie "Degrés de Liberté", les  $DdL$  de cellules sont indiqués par des étoiles et les  $DdL$  de face par des cercles. Ces

symboles ne matérialisent pas l'emplacement des  $DdL$  comme dans les représentations pour les éléments finis de Lagrange, mais permettent uniquement de les comptabiliser et de les associer aux différentes parties du maillage (cellule ou face).

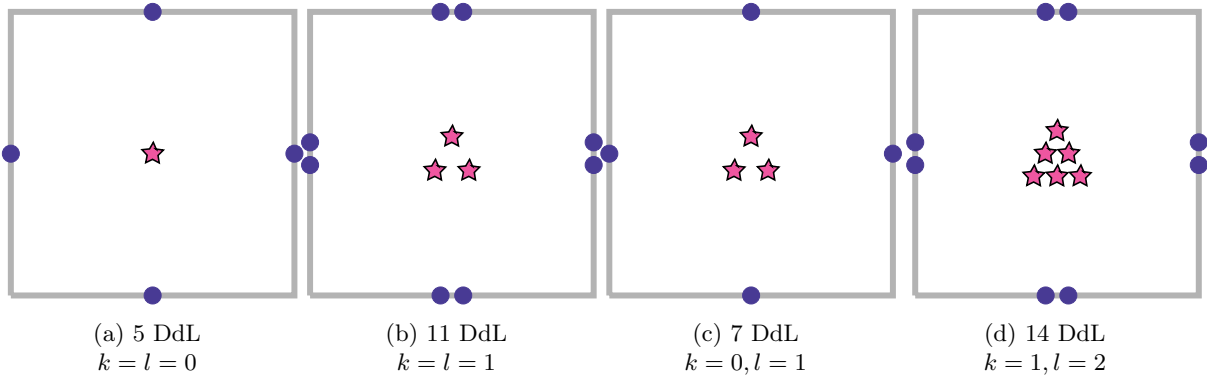
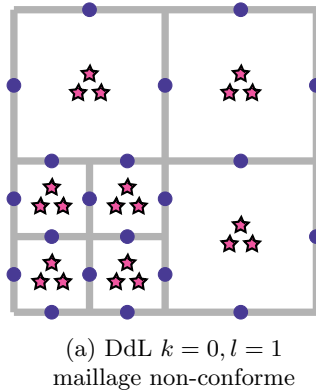


Figure 1.4: DdL dans une maille carrée,  $k \in 0, 1$ ,  $l \in \{k, k + 1\}$

Une fois assemblés, les degrés de liberté de face sont partagés par les deux cellules adjacentes, comme l'illustre la figure 1.5a. Sur ce maillage non-conforme, la cellule en bas à droite comporte 5 faces, il s'agit donc d'un pentagone, alors que la cellule en haut à droite est un bien carré comportant 4 faces. Ceci ne pose aucun problème pour la définition des DdL. En effet, la cellule pentagonale possède 5 DdL liés à ses faces. Le fait qu'il y ait un nœud orphelin au milieu d'une face n'a aucune influence sur la définition des DdL, ni sur les définitions présentées dans les chapitres suivants. C'est ce qui explique que la méthode HHO est naturellement adaptée au raffinement de maillage local non-conforme.



(a) DdL  $k = 0, l = 1$   
maillage non-conforme

Figure 1.5: DdL HHO(1,0) pour un maillage non-conforme cartésien raffiné localement.

**Opérateurs discrets.** L'idée principale de la méthode HHO est de mettre à profit les inconnues de cellule et de face afin de reconstruire localement dans chaque cellule un gradient en prenant en compte les DdL de cellule et de face. Comme nous le verrons au chapitre 2, cette reconstruction conduit à une représentation plus riche du gradient local, i.e. ayant de meilleures propriétés d'approximation, par rapport au simple gradient des DdL de cellule.

Par exemple, si  $l = 1$  et  $k = 1$ , c'est-à-dire si les inconnues de cellule et de face sont des fonctions affines par morceaux sur chaque cellule et face, respectivement, le gradient des inconnues cellule est constant par morceaux sur chaque cellule. Si l'on se réfère à la figure 1.4b, ceci signifie que les 11 degrés de liberté hybrides d'une cellule carrée ne génèrent qu'un gradient ayant 2 DdL (une constante pour

chaque composante). L'idée de la méthode HHO est de reconstruire un gradient affine, c'est-à-dire ayant 6 DdL (3 pour chaque composante). Ce gradient est déterminé en résolvant un problème local au sein de chaque cellule. Ce problème local ne dépend pas de l'équation étudiée, il est linéaire et il est calculé indépendamment pour chaque cellule. Ceci permet donc de calculer le gradient reconstruit sur chaque cellule de manière efficace. Comme ce gradient est plus riche, on peut s'attendre à ce que l'approximation soit d'un ordre plus élevé. L'idée présentée ici s'étend à tous les ordres polynomiaux.

Afin que le problème discret formulé en utilisant ce gradient reconstruit soit bien posé, il convient de stabiliser la méthode. En effet, même si le gradient reconstruit est précis, il n'est en général pas stable au sens où sa nullité n'implique pas que les DdL de cellule et de face ayant servi à sa construction soient tous égaux à une même constante. La stabilisation a pour objectif, sur chaque face du maillage, de contrôler la différence entre la trace des inconnues des cellules adjacentes et celles de la face. La stabilisation pénalise cette différence de manière linéaire (au sens des moindres carrés). Les méthodes HDG et DG sont également des méthodes comportant une stabilisation par pénalisation au sens des moindres carrés.

**Ordres de convergences attendus.** La méthode HHO permet d'utiliser des inconnues polynomiales d'ordre élevé, ce qui assure une vitesse de convergence rapide. On définit l'erreur  $L^2$  commise par la méthode HHO comme la différence entre la solution HHO et la projection de la solution exacte sur l'espace polynomial brisé des inconnues de cellule. Pour l'erreur  $H^1$ , on compare le gradient reconstruit avec le gradient de la solution exacte. On montre (voir le chapitre 2 pour plus de détails) que pour une solution exacte suffisamment régulière, l'erreur de discrétisation converge comme  $h^{k+1}$  en norme d'énergie, où  $h$  est le pas de discrétisation spatiale. En norme  $L^2$ , pour un problème elliptique offrant un shift de régularité maximal (solution dans  $H^2$  pour une donnée dans  $L^2$ ), l'erreur de discrétisation converge comme  $h^{k+2}$ .

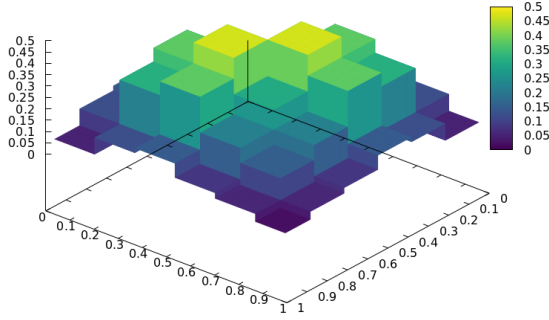
**Allure des solutions HHO.** La méthode HHO est une méthode non-conforme, c'est-à-dire que la solution peut être discontinue entre les cellules, ainsi qu'entre les cellules et les faces. Ce phénomène est illustré par les solutions présentées à la figure 1.6. Il s'agit de la résolution d'un problème de diffusion stationnaire avec un coefficient de diffusion variable, en forme de damier, variant entre 1 et 5. Le domaine considéré est  $\Omega := (0, 1)^2$  et on impose une source sinusoïdale. Il n'y a pas de solution analytique facilement calculable dans ce cas.

Les figures 1.6a et 1.6b montrent les degrés de liberté cellule par une représentation 3D pour  $k = 0$  ou  $k = 1$ , dans le cas de même ordre (c.-à-d.  $l = k$ ). Dans le cas  $k = 0$ , la solution est discontinue, car constante par morceaux. Dans le cas  $k = 1$ , la solution est affine par morceaux, et la solution n'est pas continue non plus. Les figures 1.6c et 1.6d montrent la solution HHO en coupe, le long de la droite  $y = 0.25$  et  $y = 0.5$ , respectivement. Les degrés de liberté cellule sont représentés par des lignes et des degrés de liberté face par les points. Sur les deux coupes, la solution d'ordre 0 est discontinue, et les degrés de liberté de face sont situés entre les valeurs des deux cellules adjacentes. La solution d'ordre 1 est affine par morceaux et bien plus proche de la référence en pointillés. Les degrés de liberté cellule semblent former une ligne brisée continue, bien qu'il y ait un saut modéré entre les cellules. Ce saut est visible au niveau du point  $x = 0.5$  sur la figure 1.6d.

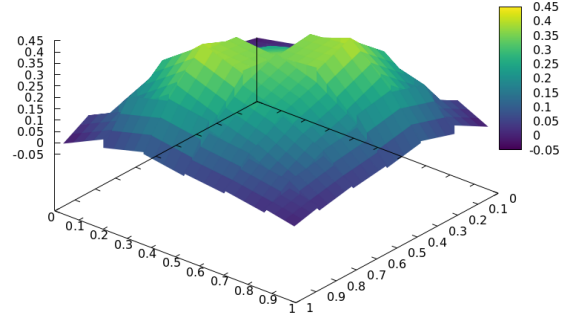
En conclusion, ces figures illustrent le caractère non-conforme de la solution HHO, en particulier lorsque les coefficients physiques varient fortement. Bien sûr, les sauts illustrés ci-dessus tendent vers zéro lorsqu'on raffine le maillage.

## 1.4 Objectifs de la thèse et plan du manuscrit

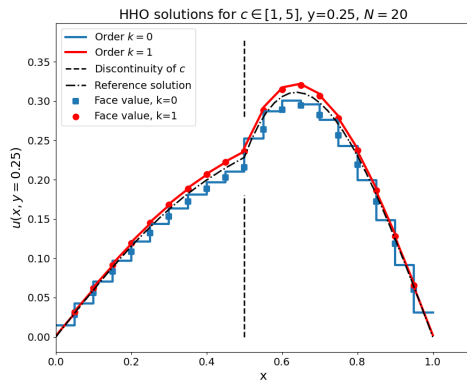
La méthode HHO a déjà été utilisée pour des simulations en mécanique statique ou quasi-statique avec des déformations élastiques ou plastiques [82, 1, 3, 2]. Il reste à étudier la possibilité de l'utiliser dans le cadre de la dynamique rapide. La similarité entre l'équation de la dynamique des structures et celle des ondes permet d'exploiter le travail déjà existant sur l'équation des ondes linéaires [34]. Ce travail a



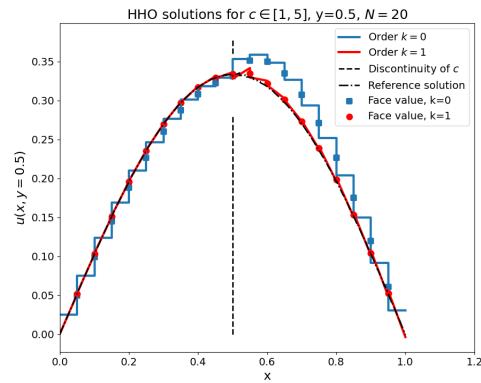
(a) Solution cellule  $(k, l) = (0, 0)$ , maillage  $h = 0.125$



(b) Solution cellule  $(k, l) = (1, 1)$ , maillage  $h = 0.125$



(c) Solution cellule et face, tranche  $y = 0.25$ , maillage  $h = 0.05$



(d) Solution cellule et face, tranche  $y = 0.5$ , maillage  $h = 0.05$ .

Figure 1.6: Diffusion stationnaire avec coefficient de diffusion variable, allure de la solution HHO pour  $k = l \in \{0, 1\}$ . Ligne du haut : solution cellule sur un maillage comprenant  $8 \times 8$  mailles, ligne du bas : tranche de solution cellule et face, maillage comprenant  $20 \times 20$  mailles.

mis en avant la difficulté de mettre en place des schémas d'intégration explicites en temps. L'objectif de cette thèse est donc de développer un schéma d'intégration explicite en temps pour les équations de propagation d'ondes acoustiques et de la dynamique des structures, toutes deux discrétisées en espace avec la méthode HHO. Ce schéma doit être performant en temps de calcul par rapport à la méthode élément finis, tout en conservant les avantages de la méthode HHO concernant les maillages polyédriques et le verrouillage volumique.

Le reste du manuscrit se compose de six chapitres. Le chapitre 2 rappelle les bases mathématiques de la méthode HHO sur l'exemple du problème de Poisson. Ce chapitre commence par une mise en perspective de la méthode HHO dans le contexte de la simulation numérique en mécanique, et en particulier pour l'équation des ondes. Les principaux concepts de la méthode HHO sont ensuite introduits : inconnues discrètes, opérateur de reconstruction de gradient et opérateur de stabilisation. Enfin, les éléments clés pour l'analyse d'erreur HHO sont présentés. Le chapitre 2 contient essentiellement des rappels de la littérature.

Le chapitre 3 est consacré à l'étude de la méthode HHO pour l'équation des ondes linéaires. Dans un premier temps, l'équation des ondes sous sa formulation du second ordre en temps est semi-discrétisée

en espace avec la méthode HHO, et une étude de convergence est menée. Les techniques de preuve sont différentes de celles de [72] pour les méthodes HDG. Elles exploitent plutôt le point de vue primal au cœur des méthodes HHO et s'appuient pour l'erreur d'énergie sur l'analyse d'erreur de [94, 18] pour les éléments finis continus et de [107] pour les méthodes dG. Pour l'erreur en norme  $L^2$ , nous utilisons un argument de dualité dans le même esprit que [150], ce qui nous amène à introduire le concept d'opérateur solution HHO. Les ordres de convergence optimaux rencontrés dans le cas statique sont préservés, c.-à-d.  $h^{k+1}$  pour la norme d'énergie et  $h^{k+2}$  pour la norme  $L^2$ . Ces résultats sont présentés dans un article publié en 2021 dans *Journal of Scientific Computing* et intitulé *Convergence Analysis of Hybrid High-Order Methods for the Wave Equation* [36]. Puis, nous présentons la discrétisation complète de l'équation des ondes avec le schéma des différences finies centrées en temps. Des ordres de convergence optimaux en espace et en temps sont prouvés. Des expériences numériques illustrant ces ordres de convergence sont également menées. La preuve repose sur des idées relativement classiques pour l'analyse du schéma des différences finies centrées lorsque celui-ci est combiné avec la méthode des éléments finis ou la méthode de Galerkin discontinu : dérivation de l'équation de l'erreur faisant apparaître les erreurs de consistance en temps et en espace comme des termes source, obtention d'une identité d'énergie sur les erreurs et, enfin, calcul d'une borne sur les termes de consistance. Il y a cependant deux différences importantes dans le cas des méthodes HHO. La première vient du fait que l'expression naïve de l'erreur d'énergie ne fournit pas une fonctionnelle convexe des dérivées en temps et en espace de l'erreur, du fait de l'erreur de consistance non nulle sur les faces du maillage. Cette difficulté est surmontée en introduisant une énergie modifiée, conduisant à un terme source supplémentaire dans le bilan d'énergie modifiée. La seconde différence vient de la nécessité de borner ce terme source additionnel, ce qui est possible car l'énergie modifiée contrôle sur la dérivée en temps de l'erreur sur les faces. Ces travaux sont présentés dans un article soumis pour publication [101]. Enfin, d'autres discrétisations en temps sont brièvement introduites à des fins de comparaison, à la fois pour la formulation d'ordre deux en temps de l'équation des ondes et celle d'ordre un. Une conclusion importante du chapitre 3, déjà mise en avant dans [34], est que la méthode HHO combinée au schéma des différences finies centrées en temps ne conduit pas à un schéma complètement explicite en temps, du fait de l'existence d'un couplage statique entre les inconnues de cellule et de face.

Le chapitre 4 propose une méthode de splitting pour résoudre ce problème de couplage statique. Cette méthode de splitting repose sur la linéarité de la stabilisation, afin de remplacer une inversion matricielle par un algorithme de point fixe. L'analyse de stabilité de cet algorithme itératif est menée, et assure une convergence du splitting à condition que la stabilisation soit multipliée par un coefficient  $\gamma$  suffisamment grand. Une borne inférieure sur  $\gamma$  est donnée, et une étude de l'impact de la valeur de  $\gamma$  sur la qualité de la solution et la condition de stabilité CFL est menée sur un cas test analytique. Enfin, la méthode HHO avec splitting est comparée, dans un cas de propagation d'ondes dans un milieu hétérogène, à la méthode des éléments finis P1 en termes d'erreur et de temps de calcul. La plupart des résultats du chapitre 4 sont présentés dans un article paru en 2023 dans *Mathematical Modelling and Numerical Analysis* et intitulé *Time-explicit Hybrid High-Order Method for the Nonlinear Acoustic Wave Equation* [144].

Le chapitre 5 étend la méthode de splitting à l'équation des ondes acoustiques non-linéaires. L'opérateur de stabilisation restant linéaire, cela permet de réutiliser la procédure de splitting de la même manière que pour l'équation des ondes linéaires. Dans le cas non-linéaire, le splitting remplace un algorithme de résolution de système non-linéaire comme un algorithme de Newton, qui devrait être utilisé à chaque pas de temps. L'efficacité computationnelle du splitting comparée à celle d'un algorithme de Newton, ainsi qu'aux méthodes éléments finis P1 et P2, est vérifiée sur deux cas tests, l'un avec un potentiel dit de "p-structure", l'autre avec une membrane vibrante non-linéaire inspirée du modèle présenté dans [53]. L'utilisation de l'accélération de point fixe d'Aitken est également étudiée. L'ensemble des résultats numériques du chapitre 5 est également présenté dans [144].

Le chapitre 6 détaille l'adaptation de la méthode HHO avec splitting à l'équation de la dynamique des structures. Nous présentons dans un premier temps la discrétisation HHO pour l'équation des ondes élastiques, ainsi que l'application de la méthode de splitting à cette équation. Une étude de la valeur des paramètres de splitting (coefficient  $\gamma$  pour la stabilisation et pas de temps critique) est ensuite menée, en prenant notamment en compte l'impact des paramètres matériaux. Une étude du paramètre de splitting

optimal est également présentée. Une expérience numérique de propagation d'ondes élastiques dans un milieu hétérogène permet ensuite de tester la méthode proposée. Enfin, une extension à l'équation de la dynamique des structures avec plasticité en grandes transformations est réalisée. Une vérification de l'efficacité de la méthode proposée en comparaison avec une résolution par un schéma semi-implicite avec algorithme de Newton est faite *via* plusieurs expériences numériques, qui illustrent notamment la robustesse de la méthode HHO au verrouillage volumique. Ces travaux font l'objet d'un article en cours de rédaction.

Enfin, le chapitre 7 synthétise les principales conclusions de ces travaux de thèse et donne des perspectives mathématiques et applicatives pour de futurs développements.



## Chapter 2

# The HHO method in a nutshell

This chapter contains a brief introduction to the HHO method. The first goal is to put the HHO method in perspective with other methods in the context of computational mechanics and, in particular, the wave equation (Section 2.1). The second goal is to introduce the key concepts underlying the HHO method using as an example of application the Poisson model problem: discrete hybrid unknowns, gradient reconstruction and stabilization (Section 2.2). Finally, the third goal is to present the main ideas underlying the error analysis of the HHO discretization (Section 2.3).

### 2.1 Some background on the HHO method

The HHO method was introduced for linear diffusion in [84] and for linear elasticity in [82]. Since then, the method has been extended to many other fields of applications involving both linear and nonlinear PDEs.

In the context of linear problems, work has since been carried out on linear diffusion with variable coefficients [6, 83], the Stokes problem [5, 86], advection-diffusion-reaction [79], the Cahn-Hilliard problem [57], and flows in fractured porous media [55, 56, 100]. In the context of nonlinear problems, work has been carried out on the Leray-Lions problem [76, 77] and more general convex minimization problems [43, 42], the stationary Navier-Stokes equations [30, 87, 46], spectral problems [38, 41], and Bingham fluids [44]. More particularly in solid mechanics, the HHO method has been applied to the Biot problem [26], nonlinear elasticity [31], hyperelasticity [1], elastoplasticity [2, 3], Kirchhoff-Love plate equations [27], biharmonic problems [93], and contact problems [45, 61, 64]. Moreover, further developments of the HHO method have been made with unfitted meshes [37, 33, 32], for multiscale methods [65] and on an a posteriori error control [25]. Details on the implementation of the HHO method are given in [63]. Finally, two books have recently been devoted to the HHO method [66, 78].

The HHO method belongs to the class of hybridizable discontinuous Galerkin methods (HDG) [70], as shown in [71]. There are essentially two differences between HHO and HDG: the HHO reconstruction operator replaces the HDG equation for the dual variable, and the stabilization of the two methods is different, at least for the equal-order setting for the cell and face degrees of freedom. Indeed, in this case, the HHO stabilization uses a nonlocal (linear) operator on the faces of the mesh attached to a given mesh cell, and this turns out to be a critical device to obtain higher-order convergence. In addition, as also shown in [71], see also [80, 66], the HHO method is closely related to the nonconforming virtual element method (ncVEM) introduced in [15] and, as also shown in [68, 93], the HHO method shares the same design principles as the weak Galerkin method (WG) introduced in [149]. Moreover, the multiscale HHO method has been bridged in [54] to the hybrid mixed multiscale method introduced in [110]. Like discontinuous Galerkin (dG) methods, the HHO method is a nonconforming method. However, unlike dG methods, the HHO method is not formulated in terms of cell unknowns and fluxes between adjacent

mesh cells, but in terms of unknowns attached to the cells and the faces of the mesh. The term *hybrid* comes from the combined use of these degrees of freedom. The cell degrees of freedom can be eliminated locally by static condensation for static equations such as stationary diffusion. The term *high-order* comes from the ease with which the HHO method can be used with arbitrary polynomial orders, without any conceptual modification of the method.

The HHO method offers many advantages: support for polyhedral meshes, optimal convergence rates, local conservation principles and computational efficiency. The HHO unknowns are polynomials attached to the cells and the faces of a mesh, of order  $l$  for the cell unknowns and of order  $k$  for the face unknowns. One must take  $k \geq 1$  for (linear) elasticity and  $k \geq 0$  for scalar-valued diffusion. Moreover, the setting is said to be of equal-order if  $l = k$ , and of mixed-order if  $l = k + 1$ . It is also possible to consider the setting with  $l = k - 1, k \geq 1$ , as in ncVEM.

In this Thesis, we study the HHO method for the space semi-discretization of wave propagation problems. We consider the acoustic wave equation and structural dynamics problems. The time discretization is performed by classical methods, mainly the leapfrog (centered finite difference) scheme.

The spatial semi-discretization of the linear acoustic wave equation by the HHO method has been introduced in [34], considering second- and first-order formulations in time. Space-optimal convergence rates in the semi-discrete case have been established in [36], and an unfitted version of the HHO method has been devised in [35]. As shown in [34], the combination of the leapfrog scheme in time with the HHO method for the second-order time formulation of the wave equation leads to a semi-implicit scheme. Indeed, at each time step, the equation for the cell unknowns is explicit, but a static coupling between the face unknowns and the cell unknowns persists. Alternatively, for the first-order formulation in time, a fully explicit scheme can be designed by combining an explicit Runge–Kutta scheme in time with the HHO method formulated in the mixed-order setting. The same difficulty is encountered when using HDG and WG schemes for the space semi-discretization of the wave equation. For example, for the first-order formulation, implicit HDG schemes [133, 132] and explicit HDG schemes [143] have been designed for the wave equation, with better computational efficiency for the explicit version reported in [120]. A symplectic HDG discretization is given in [136]. The second-order formulation in time of the wave equation with HDG space semi-discretization is carried out in [69] using a Störmer–Numerov scheme in time, leading to a fully implicit scheme (both on the cell and the face unknowns). Similarly, in the context of WG, a fully implicit scheme is developed in [112, 116].

For mechanical applications, the HHO method has the additional advantage of being robust with respect to volume locking arising in the quasi-incompressible limit and in cases of strong plasticity. Instead,  $H^1$ -conforming methods are usually prone to volume locking in the above situations. The HHO method solves this problem while still hinging on a primal formulation. In the case of static mechanics, this advantage has been exploited, for example, in [1, 2, 3, 31]. Studies on time-dependent mechanical problems are still to be carried out and the present Thesis fills (at least partly) this gap.

## 2.2 Exemple: Discretizing the Poisson problem

We consider the Poisson problem posed on a domain  $\Omega \subset \mathbb{R}^d$ ,  $d \geq 1$ , with homogeneous Dirichlet conditions for simplicity. The problem consists in finding  $u : \Omega \rightarrow \mathbb{R}$  such that

$$B(u) := -\nabla \cdot (\mu^2 \nabla u) = f, \quad \text{in } \Omega, \tag{2.1a}$$

$$u = 0, \quad \text{on } \partial\Omega. \tag{2.1b}$$

The diffusion coefficient is written  $\mu^2$  instead of  $\mu$  for consistency with the notation adopted for the wave equation.

Standard notation is used for Lebesgue and Sobolev spaces. Let  $(\cdot, \cdot)_\Omega$  denote the  $L^2$ -inner product on  $\Omega$  and  $\|\cdot\|_\Omega$  the associated norm. Boldface notation is used for  $\mathbb{R}^d$ -valued vectors and vector-valued fields, as well as for  $\mathbb{R}^{d \times d}$ -valued matrices and matrix-valued fields. Assuming  $f \in L^2(\Omega)$ , we seek the

solution of (2.1) in  $H_0^1(\Omega)$  such that

$$(\mu^2 \nabla u, \nabla w)_\Omega = (f, w)_\Omega, \quad \forall w \in H_0^1(\Omega). \quad (2.2)$$

For dimensional consistency, we consider a length scale  $\ell_\Omega$  representative of  $\Omega$ , e.g. its diameter.

### 2.2.1 Discrete setting

Let  $(\mathcal{T}_h)_{h>0}$  be a sequence of polyhedral meshes of  $\Omega$ , such that each mesh  $\mathcal{T}_h$  covers exactly  $\Omega$ . For all  $h > 0$ , let  $T$  denote a generic mesh cell in  $\mathcal{T}_h$ ,  $h_T$  its diameter and  $\mathbf{n}_T$  its unit outward normal. We set  $h := \max_{T \in \mathcal{T}_h} h_T$  for the mesh size. We say that the  $(d-1)$ -dimensional set  $F$  is a mesh face if there is a hyperplane  $H_F$  such that either  $F = H_F \cap \partial T_- \cap \partial T_+$  for two distinct mesh cells  $T_-$  and  $T_+$  (and  $F$  is called mesh interface) or  $F = H_F \cap \partial T_- \cap \partial \Omega$  (and  $F$  is called mesh boundary face). The collection of all the mesh faces is denoted  $\mathcal{F}_h$ . For all  $T \in \mathcal{T}_h$ , we denote by  $\mathcal{F}_T$  the collection of the mesh faces composing the boundary  $\partial T$ , and for all  $F \in \mathcal{F}_T$ , we set  $\mathbf{n}_{TF} := \mathbf{n}_T|_F$ .

The mesh sequence is assumed to be shape-regular. Here, we only consider meshes composed of cells with simple shape (triangles, tetrahedra, quadrangles and hexahedra) so that the classical notion of regularity by Ciarlet is sufficient for our purposes; for mesh regularity with more general shapes, we refer the reader, e.g., to [81, 39, 78, 66]. Mesh regularity implies, in particular, that for all  $h > 0$ , all  $T \in \mathcal{T}_h$  and all  $F \in \mathcal{F}_T$ , the diameter  $h_F$  of  $F$  is uniformly comparable to  $h_T$ .

The coefficient  $\mu$  is assumed to be piecewise constant on a polyhedral partition of  $\Omega$ , and all the meshes of  $\mathcal{T}_h$  are assumed to be compatible with this partition. As a consequence,  $\mu$  takes a constant value denoted by  $\mu_T$  in each mesh cell  $T \in \mathcal{T}_h$ . We assume that there are  $0 < \mu_b \leq \mu_\sharp < \infty$  such that  $\mu_b \leq \mu_T \leq \mu_\sharp$  for all  $T \in \mathcal{T}_h$ , and we assume that  $\frac{\mu_b}{\mu_\sharp}$  is not too large, so that it can be hidden in the generic constants used in the error analysis.

Let the integer  $k \geq 0$  be the polynomial order of the face unknowns and let  $l \in \{k, k+1\}$  be the order of the cell unknowns. Recall that the setting is said to be of equal-order if  $l = k$  and of mixed-order if  $l = k+1$ . Let  $\mathbb{P}_d^l(T)$  (resp.  $\mathbb{P}_{d-1}^k(F)$ ) denote the set of  $d$ -variate (resp.  $(d-1)$ -variate) polynomials of degree at most  $l$  (resp.  $k$ ) restricted to the cell  $T \in \mathcal{T}_h$  (resp. to the face  $F \in \mathcal{F}_h$ ). The linear space composed of all the cell degrees of freedom is denoted  $\mathcal{U}_\mathcal{T}^l$ , and the linear space composed of all the face degrees of freedom is denoted  $\mathcal{U}_\mathcal{F}^k$ . These spaces are defined as Cartesian products in the form

$$\mathcal{U}_\mathcal{T}^l := \bigtimes_{T \in \mathcal{T}_h} \mathbb{P}_d^l(T), \quad \mathcal{U}_\mathcal{F}^k := \bigtimes_{F \in \mathcal{F}_h} \mathbb{P}_{d-1}^k(F), \quad (2.3)$$

and we slightly abuse the notation by viewing an element  $w_\mathcal{T} = (w_T)_{T \in \mathcal{T}_h} \in \mathcal{U}_\mathcal{T}^l$  as a function defined a.e. over  $\Omega$  such that  $w_\mathcal{T}|_T := w_T$  for all  $T \in \mathcal{T}_h$ . The collection of all the cell and face degrees of freedom is the hybrid space

$$\widehat{\mathcal{U}}_h^{l,k} := \mathcal{U}_\mathcal{T}^l \times \mathcal{U}_\mathcal{F}^k. \quad (2.4)$$

A generic element of  $\widehat{\mathcal{U}}_h^{l,k}$  is denoted  $\widehat{w}_h := (w_\mathcal{T}, w_\mathcal{F}) \in \mathcal{U}_\mathcal{T}^l \times \mathcal{U}_\mathcal{F}^k$  and, in what follows, variables with hats refer to hybrid variables. For a given cell  $T \in \mathcal{T}_h$ , we also define a local hybrid space of degrees of freedom

$$\widehat{\mathcal{U}}_T^{l,k} := \mathbb{P}_d^l(T) \times \mathcal{U}_{\partial T}^k, \quad \mathcal{U}_{\partial T}^k := \bigtimes_{F \in \mathcal{F}_T} \mathbb{P}_{d-1}^k(F). \quad (2.5)$$

Then  $\widehat{w}_T := (w_T, w_{\partial T} = (w_F)_{F \in \mathcal{F}_T}) \in \widehat{\mathcal{U}}_T^{l,k}$  denotes a generic local hybrid unknown in  $T$ , composed of one cell unknown and the collection of the face unknowns for all the faces in  $\mathcal{F}_T$ . As above, we slightly abuse the notation by viewing an element  $w_{\partial T} = (w_F)_{F \in \mathcal{F}_T} \in \mathcal{U}_{\partial T}^k$  as a function defined a.e. over  $\partial T$  such that  $w_{\partial T}|_F := w_F$  for all  $F \in \mathcal{F}_T$ . Let  $\mathcal{U}_{\mathcal{F},0}^k := \{v_\mathcal{F} \in \mathcal{U}_\mathcal{F}^k, \text{ s.t. } v_F = 0, \forall F \subset \partial \Omega\}$  be the subspace of face unknowns respecting the homogeneous Dirichlet conditions. The subspace of hybrid unknowns respecting the homogeneous Dirichlet conditions is denoted

$$\widehat{\mathcal{U}}_{h,0}^{l,k} := \mathcal{U}_\mathcal{T}^l \times \mathcal{U}_{\mathcal{F},0}^k. \quad (2.6)$$

$L^2$ -orthogonal projections onto polynomial spaces are denoted with the symbol  $\Pi$ . For instance, for all  $T \in \mathcal{T}_h$ ,  $\Pi_T^l$  is the projection onto  $\mathbb{P}_d^l(T)$  and for all  $F \in \mathcal{F}_h$ ,  $\Pi_F^k$  the projection onto  $\mathbb{P}_{d-1}^k(F)$ . The  $L^2$ -orthogonal projection onto the broken polynomial spaces  $\mathcal{U}_T^l$  and  $\mathcal{U}_F^k$  is denoted by  $\Pi_T^l$  and  $\Pi_F^k$  respectively. Let  $(\cdot, \cdot)_T$  and  $(\cdot, \cdot)_F$  respectively denote the  $L^2$ -inner product in the cell  $T \in \mathcal{T}_h$  and the face  $F \in \mathcal{F}_h$ . For all  $v_{\partial T}, w_{\partial T} \in \mathcal{U}_{\partial T}^k$ , we also define  $(v_{\partial T}, w_{\partial T})_{\partial T} := \sum_{F \in \mathcal{F}_T} (v_F, w_F)_F$ . Let  $\|\cdot\|_T, \|\cdot\|_F$  and  $\|\cdot\|_{\partial T}$  denote the norms associated respectively with the  $L^2$ -inner products  $(\cdot, \cdot)_T, (\cdot, \cdot)_F$  and  $(\cdot, \cdot)_{\partial T}$ .

For the comparison of the discrete HHO solution with the continuous global solution, a projection operator is needed. One first defines the local projector  $\hat{I}_T^{k,l} : H^1(T) \rightarrow \hat{\mathcal{U}}_T^{l,k}$  for all  $T \in \mathcal{T}_h$  such that, for all  $w \in H^1(T)$ ,

$$\hat{I}_T^{k,l}(w) := (\Pi_T^l(w), (\Pi_F^k(w))_{F \in \partial T}). \quad (2.7)$$

Then the global projection operator can be defined as  $\hat{I}_h^{k,l} : H^1(\Omega) \rightarrow \hat{\mathcal{U}}_h^{l,k}$  such that, for all  $w \in H^1(\Omega)$ ,

$$\hat{I}_h^{k,l}(w) := ((\Pi_T^l(w))_{T \in \mathcal{T}_h}, (\Pi_F^k(w))_{F \in \mathcal{F}_h}). \quad (2.8)$$

Notice that this global projection operator maps functions from  $H_0^1(\Omega)$  onto  $\hat{\mathcal{U}}_{h,0}^{l,k}$ . The definition of  $\hat{I}_h^{k,l}$  is meaningful since a function  $v \in H^1(\Omega)$  does not jump across the mesh interfaces. Let us record here for later use the following discrete inverse inequality (see, e.g. [97, Chap. 12]).

**Lemma 2.2.1** (Discrete inverse inequalities). *Let the polynomial degree  $l$  be fixed. There is  $C_{\text{disc}} > 0$  such that, for all  $h > 0$ , all  $T \in \mathcal{T}_h$  and all  $w \in \mathbb{P}_d^l(T)$ ,*

$$\|\nabla w\|_T \leq C_{\text{disc}} h_T^{-1} \|w\|_T, \quad (2.9a)$$

$$\|w\|_{\partial T} \leq C_{\text{disc}} h_T^{-\frac{1}{2}} \|w\|_T. \quad (2.9b)$$

## 2.2.2 HHO operators and basic properties

The HHO discretization relies on two key operators: a gradient reconstruction operator and a stabilization operator. Both operators are local, i.e. they are defined independently in every mesh cell  $T \in \mathcal{T}_h$ .

**Gradient reconstruction.** The local gradient reconstruction operator builds a gradient in the cell  $T \in \mathcal{T}_h$  from the local cell and face unknowns in  $\hat{\mathcal{U}}_T^{l,k}$ . This operator  $\mathbf{G}_T^k : \hat{\mathcal{U}}_T^{l,k} \rightarrow \mathbb{P}_d^k(T; \mathbb{R}^d)$  is evaluated by solving the following problem: For all  $\hat{v}_T \in \hat{\mathcal{U}}_T^{l,k}$ ,

$$(\mathbf{G}_T^k(\hat{v}_T), \mathbf{q})_T = (\nabla v_T, \mathbf{q})_T + (v_{\partial T} - v_T, \mathbf{q} \cdot \mathbf{n}_T)_{\partial T}, \quad \forall \mathbf{q} \in \mathbb{P}_d^k(T; \mathbb{R}^d), \quad (2.10)$$

where  $\mathbb{P}_d^k(T; \mathbb{R}^d)$  denotes the space of  $\mathbb{R}^d$ -valued  $d$ -variate polynomials of degree  $k$  in the cell  $T$ . In practice, each component of the reconstructed gradient is found independently by inverting the mass matrix associated with a chosen scalar-valued basis of  $\mathbb{P}_d^k(T)$ . Using the gradient reconstruction operator in every mesh cell  $T \in \mathcal{T}_h$ , we define the local stiffness bilinear form  $b_T$  such that, for all  $\hat{v}_T, \hat{w}_T \in \hat{\mathcal{U}}_T^{l,k}$ ,

$$b_T(\hat{v}_T, \hat{w}_T) := \mu_T^2 (\mathbf{G}_T^k(\hat{v}_T), \mathbf{G}_T^k(\hat{w}_T))_T. \quad (2.11)$$

**Potential reconstruction.** One can also build a potential reconstruction operator  $R_T^{k+1} : \hat{\mathcal{U}}_T^{l,k} \rightarrow \mathbb{P}_d^{k+1}(T)$  by solving, for all  $\hat{v}_T \in \hat{\mathcal{U}}_T^{l,k}$ , the following Neumann problem:

$$(\nabla R_T^{k+1}(\hat{v}_T), \nabla q)_T = (\nabla v_T, \nabla q)_T + (v_{\partial T} - v_T, \nabla q \cdot \mathbf{n}_T)_{\partial T}, \quad \forall q \in \mathbb{P}_d^{k+1}(T), \quad (2.12)$$

and the mean-value condition

$$(R_T^{k+1}(\hat{v}_T), 1)_T = (v_T, 1)_T. \quad (2.13)$$

The computation of  $R_T^{k+1}(\hat{v}_T)$  requires inverting the stiffness matrix for a chosen basis of  $\mathbb{P}_d^{k+1,0}(T) := \{q \in \mathbb{P}_d^{k+1}(T), \text{ s.t. } (q, 1)_T = 0\}$ . In practice, the Neumann problem (2.12) is solved first by looking for

$R_{T,0}^{k+1}(\hat{v}_T) \in \mathbb{P}_d^{k+1,0}(T)$  having zero mean-value on  $T$ . Then the reconstructed potential can be found by using the average condition (2.13) and setting

$$R_T^{k+1}(\hat{v}_T) = R_{T,0}^{k+1}(\hat{v}_T) + \frac{\int_T v_T dT}{\int_T 1 dT}. \quad (2.14)$$

Notice that we have

$$\nabla R_T^{k+1}(\hat{v}_T) = \Pi_{\nabla \mathbb{P}_d^{k+1}(T)} \mathbf{G}_T^k(\hat{v}_T) \quad (2.15)$$

for all  $\hat{v}_T \in \widehat{\mathcal{U}}_T^{l,k}$ .

**Stabilization.** The role of the stabilization is to weakly enforce the matching between the cell and the face unknowns at each mesh face. Let  $T \in \mathcal{T}_h$ . For all  $\hat{w}_T \in \widehat{\mathcal{U}}_T^{l,k}$ , set  $\delta_T(\hat{w}_T) := w_{\partial T} - w_T|_{\partial T}$  on  $\partial T$  and  $\delta_{TF}(\hat{w}_T) := \delta_T(\hat{w}_T)|_F$  for all  $F \in \mathcal{F}_T$ . In the mixed-order setting, the local stabilization operator  $S_{TF}$  is defined as

$$S_{TF}(\hat{w}_T) := \Pi_F^k(\delta_{TF}(\hat{w}_T)), \quad \forall \hat{w}_T \in \widehat{\mathcal{U}}_T^{l,k}, \quad (2.16)$$

and leads to the so-called Lehrenfeld-Schöberl stabilization in the HDG context (see, e.g. [124, 125]). In the equal-order setting, the definition of  $S_{TF}$  requires the computation of  $R_T^{k+1}$ , which is for instance motivated in [66, Lem. 2.7], and writes

$$S_{TF}(\hat{w}_T) := \Pi_F^k(\delta_{TF}(\hat{w}_T) + (I - \Pi_T^k)R_T^{k+1}(0, \delta_T(\hat{w}_T))|_F), \quad \forall \hat{w}_T \in \widehat{\mathcal{U}}_T^{l,k} =: \widehat{\mathcal{U}}_T^k. \quad (2.17)$$

In both equal- and mixed-order settings, we define, for all  $\hat{v}_T, \hat{w}_T \in \widehat{\mathcal{U}}_T^{l,k}$ , the local stabilization form as

$$s_T(\hat{v}_T, \hat{w}_T) = \gamma \mu_T^2 \sum_{F \in \mathcal{F}_T} \eta_{TF}^{-1} (S_{TF}(\hat{v}_T), S_{TF}(\hat{w}_T))_F, \quad (2.18)$$

with the scaling factor  $\eta_{TF}$  equal to  $h_F$  or  $h_T$ , and  $\gamma > 0$  a scaling parameter. Choosing  $\eta_{TF} := h_T$  is often more relevant, since it is independent of the size of the faces. In what follows, we set  $S_{\partial T}(\hat{v}_T)|_F := S_{TF}(\hat{v}_T)$  for all  $F \in \mathcal{F}_T$ .

**Remark 2.2.2** (Choice of  $\gamma$ ). Any positive value can be chosen for the scaling parameter  $\gamma$ . For the plain HHO method, the choice  $\gamma = 1$  is made in [82, 84]. In the present setting, this parameter is introduced to tune the scale of the stabilization compared to the stiffness bilinear form  $b_T$ . If  $\gamma \ll 1$ , the stabilization is much smaller than the stiffness, and the problem is close to being singular. The case  $\gamma \gg 1$  does not entail singularity issues but may lead to larger errors than when taking  $\gamma \simeq 1$ . Taking larger values of  $\gamma$  also lowers the CFL stability restriction on the time step, as we shall see below in the context of explicit time-marching schemes for wave propagation problems. Thus, very large values of  $\gamma$  are not computationally efficient due to the large number of time steps. We refer the reader to Section 4.2.2 for further insight.

**Stability and boundedness.** A direct verification shows that the map  $\|\cdot\|_{\text{HHO}} : \widehat{\mathcal{U}}_h^{l,k} \rightarrow \mathbb{R}_+$  such that

$$\|\hat{v}_h\|_{\text{HHO}}^2 := \sum_{T \in \mathcal{T}_h} \mu_T^2 \{ \|\nabla v_T\|_T^2 + h_T^{-1} \|v_{\partial T} - v_T\|_{\partial T}^2 \}, \quad \forall \hat{v}_h \in \widehat{\mathcal{U}}_h^{l,k}, \quad (2.19)$$

defines a norm on  $\widehat{\mathcal{U}}_{h,0}^{l,k}$  (and a seminorm on  $\widehat{\mathcal{U}}_h^{l,k}$ ), and we have the following important result [84, Lem. 4], [82, Lem. 4].

**Lemma 2.2.3** (Stability and boundedness). *There are  $0 < \alpha \leq \varpi < \infty$  such that, for all  $\hat{v}_h \in \widehat{\mathcal{U}}_{h,0}^{l,k}$  and all  $h > 0$ ,*

$$\alpha \|\hat{v}_h\|_{\text{HHO}}^2 \leq \|\mu \mathbf{G}_{\mathcal{T}}^k(\hat{v}_h)\|_{\Omega}^2 + |\hat{v}_h|_S^2 \leq \varpi \|\hat{v}_h\|_{\text{HHO}}^2, \quad (2.20)$$

with the seminorm  $|\hat{v}_h|_S^2 := s_h(\hat{v}_h, \hat{v}_h)$  and the global gradient reconstruction such that  $\mathbf{G}_{\mathcal{T}}^k(\hat{v}_h)|_T := \mathbf{G}_T^k(\hat{v}_T)$  for all  $T \in \mathcal{T}_h$ .

We also record here for later use the following discrete Poincaré inequality.

**Lemma 2.2.4** (Discrete Poincaré inequality). *There is  $C_{\text{dP}}$  such that for all  $h > 0$  and all  $\hat{w}_h \in \hat{\mathcal{U}}_{h,0}^{l,k}$ ,*

$$\|w_{\mathcal{T}}\|_{\Omega} \leq C_{\text{dP}} \mu_{\sharp}^{-1} \ell_{\Omega} \|\hat{w}_h\|_{\text{HHO}}. \quad (2.21)$$

*Proof.* Let  $\hat{w}_h \in \hat{\mathcal{U}}_{h,0}^{l,k}$ . There is  $\mathbf{v} \in \mathbf{H}^1(\Omega)$  such that  $\nabla \cdot \mathbf{v} = w_{\mathcal{T}}$  and  $\|\mathbf{v}\|_{\mathbf{H}^1(\Omega)} \leq C \ell_{\Omega} \|w_{\mathcal{T}}\|_{\Omega}$ , with  $\|v\|_{\mathbf{H}^1(\Omega)} := \{\|v\|_{\Omega}^2 + \ell_{\Omega}^2 \|\nabla v\|_{\Omega}^2\}^{\frac{1}{2}}$ , and where  $C$  only depends on  $\Omega$ .

We integrate by parts and use the fact that  $\mathbf{v}$  is single-valued at all the mesh interfaces and  $w_{\partial T}$  vanishes at all the mesh boundary faces to obtain

$$\|w_{\mathcal{T}}\|_{\Omega}^2 = (w_{\mathcal{T}}, \nabla \cdot \mathbf{v})_{\Omega} = \sum_{T \in \mathcal{T}_h} (\nabla w_T, \mathbf{v})_T + (w_T - w_{\partial T}, \mathbf{v} \cdot \mathbf{n}_{TF})_{\partial T}.$$

Using the definition of the  $\|\cdot\|_{\text{HHO}}$  norm, the Cauchy-Schwarz inequality and a multiplicative trace inequality in  $\mathbf{H}^1(T)$  giving  $h_T^{\frac{1}{2}} \|\mathbf{v}\|_{\partial T} \leq C(\|\mathbf{v}\|_T + h_T \|\nabla \mathbf{v}\|_T)$  (see [97, Lem. 12.15]), we infer that

$$\|w_{\mathcal{T}}\|_{\Omega}^2 \leq C \mu_{\sharp}^{-1} \|\hat{w}_h\|_{\text{HHO}} \left\{ \sum_{T \in \mathcal{T}_h} \|\mathbf{v}\|_T^2 + h_T \|\mathbf{v}\|_{\partial T}^2 \right\}^{\frac{1}{2}} \leq C' \mu_{\sharp}^{-1} \|\hat{w}_h\|_{\text{HHO}} \|\mathbf{v}\|_{\mathbf{H}^1(\Omega)}.$$

Using the inequality  $\|\mathbf{v}\|_{\mathbf{H}^1(\Omega)} \leq C \ell_{\Omega} \|w_{\mathcal{T}}\|_{\Omega}$  yields the claim.  $\square$

**Approximation properties.** Let  $\mathcal{E}_{\mathcal{T}}^{k+1} : H^1(\Omega) \rightarrow \mathcal{U}_{\mathcal{T}}^{k+1}$  denote the broken elliptic projection onto  $\mathcal{U}_{\mathcal{T}}^{k+1}$ , so that for all  $v \in H^1(\Omega)$  and all  $T \in \mathcal{T}_h$ ,  $\mathcal{E}_{\mathcal{T}}^{k+1}(v)|_T$  is uniquely defined by the relations

$$(\nabla(\mathcal{E}_T^{k+1}(v) - v), \nabla q)_T = 0, \quad \forall q \in \mathbb{P}_d^{k+1,0}(T), \quad \text{and} \quad (\mathcal{E}_T^{k+1}(v) - v, 1)_T = 0.$$

The following result is of great importance in view of error analysis and is established in [84, Lem. 3], [66, Lem. 3.1].

**Lemma 2.2.5** (Approximation properties for  $R_T^{k+1} \circ \hat{I}_T^{k,l}$  and  $\mathbf{G}_T^k \circ \hat{I}_T^{k,l}$ ). *The following holds true for all  $h > 0$  and all  $T \in \mathcal{T}_h$ :*

$$\mathbf{G}_T^k(\hat{I}_T^{k,l}(v)) = \Pi_T^k(\nabla v), \quad \nabla R_T^{k+1}(\hat{I}_T^{k,l}(v)) = \nabla \mathcal{E}_T^{k+1}(v), \quad \forall v \in H^1(T). \quad (2.22)$$

Moreover, there exists a real number  $C > 0$  such that, for all  $h > 0$ , all  $T \in \mathcal{T}_h$  and all  $v \in H^{k+2}(T)$ , we have

$$\|v - R_T^{k+1}(\hat{I}_T^{k,l}(v))\|_T + h_T^{\frac{1}{2}} \|v - R_T^{k+1}(\hat{I}_T^{k,l}(v))\|_{\partial T} \leq C h_T^{k+2} |v|_{H^{k+2}(T)}, \quad (2.23)$$

and

$$\|\nabla(v - R_T^{k+1}(\hat{I}_T^{k,l}(v)))\|_T + h_T^{\frac{1}{2}} \|\nabla(v - R_T^{k+1}(\hat{I}_T^{k,l}(v)))\|_{\partial T} \leq C h_T^{k+1} |v|_{H^{k+2}(T)}, \quad (2.24a)$$

$$\|\nabla v - \mathbf{G}_T^k(\hat{I}_T^{k,l}(v))\|_T + h_T^{\frac{1}{2}} \|\nabla v - \mathbf{G}_T^k(\hat{I}_T^{k,l}(v))\|_{\partial T} \leq C h_T^{k+1} |v|_{H^{k+2}(T)}. \quad (2.24b)$$

The following result is established in [84, Eq. (45)], [66, Lem. 2.7].

**Lemma 2.2.6** (Approximation property for  $S_{\partial T} \circ \hat{I}_T^{k,l}$ ). *There is a constant  $C > 0$  such that for all  $h > 0$ , all  $T \in \mathcal{T}_h$ , and all  $v \in H^1(T)$ , we have*

$$\|S_{\partial T}(\hat{I}_T^{k,l}(v))\|_{\partial T} \leq C h_T^{\frac{1}{2}} \|\nabla(v - P_T^{k+1}(v))\|_T, \quad (2.25)$$

where  $P_T^{k+1} := \mathcal{E}_T^{k+1}$  if  $l = k$  and  $P_T^{k+1} := \Pi_T^{k+1}$  if  $l = k + 1$ .

**Variants.** The space where the gradient is reconstructed can be chosen differently. The choice made initially in [82, 84] is to set

$$\mathbf{G}_T^{k,\nabla}(\hat{v}_T) := \nabla R_T^{k+1}(\hat{v}_T). \quad (2.26)$$

The choice made in (6.21a) is more relevant in nonlinear problems [76, 31, 1], but the advantage of (2.26) is that  $\mathbf{G}_T^{k,\nabla}(\hat{v}_T)$  is curl-free, which gives the direct relation (2.26) and not a projection like in (2.15). One can also reconstruct the gradient in a larger space than  $\mathbb{P}_d^k(T; \mathbb{R}^d)$  by simply posing the problem (6.21a) in a larger space. However, taking a too large space will eventually impact the optimal convergence of the method. This happens when the normal components of the reconstructed gradient on each face  $F \in \mathcal{F}_T$  sit in a larger space than  $\mathbb{P}_d^k(F)$ , since the first identity in (2.22) then fails. Thus, the largest space ensuring optimal convergence is the Raviart–Thomas–Nédelec space of order  $k$ ,  $\mathbf{RTN}^k(T) := \mathbb{P}_d^k(T; \mathbb{R}^d) \oplus \left( \mathbf{X} \mathbb{P}_{d,H}^k(T) \right)$ , where  $\mathbf{X} = (x, y)$  in 2D and  $\mathbf{X} = (x, y, z)$  in 3D and  $\mathbb{P}_{d,H}^k(T)$  is composed of homogenous  $d$ -variate polynomials of order  $k$ . This reconstruction (and any reconstruction in a larger space, e.g.  $\mathbb{P}_d^{k+1}(T; \mathbb{R}^d)$ ) leads to a stable gradient (i.e. having a one-dimensional kernel in  $\widehat{\mathcal{U}}_T^{l,k}$ ), thus making the stabilization dispensable for (2.20) to hold. This leads to an *unstabilized* HHO setting. It must, however, be noticed that when reconstructing the gradient in  $\mathbb{P}_d^{k+1}(T; \mathbb{R}^d)$ , the approximation property is only of order  $h_T^k$  in Lemma 2.2.5. The above discussion is summarized in Table 2.1.

Reconstruction space	$\nabla \mathbb{P}_d^{k+1}(T)$	$\mathbb{P}_d^k(T; \mathbb{R}^d)$	$\mathbf{RTN}^k(T)$	$\mathbb{P}_d^{k+1}(T; \mathbb{R}^d)$
Stabilization needed	yes	yes	no	no
Upper bound (2.24)	$k + 1$	$k + 1$	$k + 1$	$k$
Curl-free $\mathbf{G}_T^k$	yes	no	no	no

Table 2.1: Comparison of reconstruction spaces for the HHO gradient.

### 2.2.3 Global discretization

**Functional formulation.** The global bilinear forms  $b_h$  and  $s_h$  are defined, for all  $\hat{v}_h, \hat{w}_h \in \widehat{\mathcal{U}}_h^{l,k}$ , as

$$b_h(\hat{v}_h, \hat{w}_h) := \sum_{T \in \mathcal{T}_h} b_T(\hat{v}_T, \hat{w}_T), \quad s_h(\hat{v}_h, \hat{w}_h) := \sum_{T \in \mathcal{T}_h} s_T(\hat{v}_T, \hat{w}_T), \quad (2.27)$$

and the global stiffness bilinear form  $a_h$  is defined as

$$a_h(\hat{v}_h, \hat{w}_h) := b_h(\hat{v}_h, \hat{w}_h) + s_h(\hat{v}_h, \hat{w}_h). \quad (2.28)$$

The discrete scheme for the Poisson equation (2.1) consists of finding  $\hat{u}_h := (u_{\mathcal{T}}, u_{\mathcal{F}}) \in \widehat{\mathcal{U}}_{h,0}^{l,k}$  such that,

$$a_h(\hat{u}_h, \hat{w}_h) = (f, w_{\mathcal{T}})_{\Omega}, \quad \forall \hat{w}_h := (w_{\mathcal{T}}, w_{\mathcal{F}}) \in \widehat{\mathcal{U}}_{h,0}^{l,k}. \quad (2.29)$$

Notice that the homogeneous Dirichlet boundary condition is enforced by the condition  $\hat{u}_h \in \widehat{\mathcal{U}}_{h,0}^{l,k}$ .

**Algebraic formulation.** Let  $N_{\mathcal{T}} := \dim(\mathcal{U}_{\mathcal{T}}^l)$  and  $N_{\mathcal{F}} := \dim(\mathcal{U}_{\mathcal{F},0}^k)$  and let  $\{\phi_i\}_{1 \leq i \leq N_{\mathcal{T}}}, \{\psi_i\}_{1 \leq i \leq N_{\mathcal{F}}}$  be bases of  $\mathcal{U}_{\mathcal{T}}^l$  and  $\mathcal{U}_{\mathcal{F},0}^k$ , respectively. Let also  $(\mathbf{U}_{\mathcal{T}}, \mathbf{U}_{\mathcal{F}}) \in \mathbb{R}^{N_{\mathcal{T}}} \times \mathbb{R}^{N_{\mathcal{F}}}$  be the vector of degrees of freedom of the solution  $\hat{u}_h$  on these bases, and  $\mathbf{F}_{\mathcal{T}} \in \mathbb{R}^{N_{\mathcal{T}}}$  the vector having components  $((f, \phi_i)_{\Omega})_{1 \leq i \leq N_{\mathcal{T}}}$ . The algebraic form of (2.29) is

$$\begin{bmatrix} \mathcal{A}_{\mathcal{T}\mathcal{T}} & \mathcal{A}_{\mathcal{T}\mathcal{F}} \\ \mathcal{A}_{\mathcal{F}\mathcal{T}} & \mathcal{A}_{\mathcal{F}\mathcal{F}} \end{bmatrix} \begin{pmatrix} \mathbf{U}_{\mathcal{T}} \\ \mathbf{U}_{\mathcal{F}} \end{pmatrix} = \begin{bmatrix} \mathbf{F}_{\mathcal{T}} \\ 0 \end{bmatrix}, \quad (2.30)$$

with  $\mathcal{A}$  the symmetric definite positive stiffness matrix associated with the bilinear form  $a_h$ . The stiffness matrix can be written as the sum of two matrices

$$\mathcal{A} = \mathcal{B} + \mathcal{S}, \quad (2.31)$$

where  $\mathcal{B}$  is associated with the stiffness bilinear form  $b_h$  and  $\mathcal{S}$  with the stabilization bilinear form  $s_h$ . The sub-matrix  $\mathcal{A}_{\mathcal{T}\mathcal{T}}$  is block-diagonal since both  $b_h$  and  $s_h$  do not couple cell degrees of freedom from different cells. Instead,  $\mathcal{A}_{\mathcal{F}\mathcal{F}}$  does not have a block-diagonal structure since the gradient reconstruction operator couples the face degrees of freedom of all the faces of a given mesh cell. It is useful to notice that in the mixed-order setting,  $\mathcal{S}_{\mathcal{F}\mathcal{F}}$  is block-diagonal, since it reduces on each mesh face to the mass matrix associated with the degrees of freedom on that face.

**Static condensation.** Considering the structure of  $\mathcal{A}_{\mathcal{T}\mathcal{T}}$ , a static condensation can be performed. This consists in computing the Schur complement of the matrix  $\mathcal{A}_{\mathcal{T}\mathcal{T}}$  and reducing the hybrid linear system (2.30) to a linear system on the face unknowns only:

$$(\mathcal{A}_{\mathcal{F}\mathcal{F}} - \mathcal{A}_{\mathcal{F}\mathcal{T}}\mathcal{A}_{\mathcal{T}\mathcal{T}}^{-1}\mathcal{A}_{\mathcal{T}\mathcal{F}})\mathbf{U}_{\mathcal{F}} = -\mathcal{A}_{\mathcal{F}\mathcal{T}}\mathcal{A}_{\mathcal{T}\mathcal{T}}^{-1}\mathbf{F}_{\mathcal{T}}. \quad (2.32)$$

The cell unknowns can then be obtained by a post-processing:

$$\mathbf{U}_{\mathcal{T}} = \mathcal{A}_{\mathcal{T}\mathcal{T}}^{-1}(\mathbf{F}_{\mathcal{T}} - \mathcal{A}_{\mathcal{T}\mathcal{F}}\mathbf{U}_{\mathcal{F}}). \quad (2.33)$$

This operation can be performed locally in each mesh cell owing to the block-diagonal structure of  $\mathcal{A}_{\mathcal{T}\mathcal{T}}$ . The static condensation procedure reduces substantially the cost of solving the linear system (2.30). Indeed, since the face unknowns are  $(d-1)$ -variate polynomials, the total number of face unknowns,  $\mathcal{O}(k^{d-1})$ , grows slower than the number of cell unknowns,  $\mathcal{O}(k^d)$ , when increasing the polynomial order. Notice, however, that the Schur complement matrix  $\mathcal{A}_{\mathcal{F}\mathcal{F}} - \mathcal{A}_{\mathcal{F}\mathcal{T}}\mathcal{A}_{\mathcal{T}\mathcal{T}}^{-1}\mathcal{A}_{\mathcal{T}\mathcal{F}}$  is not as sparse as  $\mathcal{A}$ .

## 2.3 Error analysis

In what follows, the symbol  $C$  denotes a generic positive constant whose value can change at each occurrence as long as it is independent of the mesh size.

### 2.3.1 Consistency error

In view of the error analysis for the time-dependent wave equation, a useful notion is the HHO solution map  $\hat{J}_h = (J_{\mathcal{T}}, J_{\mathcal{F}}) : Y \rightarrow \hat{\mathcal{U}}_{h,0}^{l,k}$ , with  $Y := \{v \in H_0^1(\Omega), B(v) = -\nabla \cdot (\mu^2 \nabla v) \in L^2(\Omega)\}$ , such that for all  $v \in Y$ ,  $\hat{J}_h(v) \in \hat{\mathcal{U}}_{h,0}^{l,k}$  is uniquely defined by the following equations:

$$a_h(\hat{J}_h(v), \hat{q}_h) = (B(v), q_{\mathcal{T}})_{\Omega}, \quad \forall \hat{q}_h \in \hat{\mathcal{U}}_{h,0}^{l,k}, \quad (2.34)$$

The coercivity of  $a_h$  on  $\hat{\mathcal{U}}_{h,0}^{l,k}$  (see Lemma 2.2.3) ensures that  $\hat{J}_h(v)$  is well-defined by means of (2.34). Notice that (2.34) allows us to rewrite the discrete HHO solution of (2.29) as

$$\hat{u}_h = \hat{J}_h(u).$$

For all  $v \in H^{1+\nu}(\Omega)$ ,  $\nu > \frac{1}{2}$ , we consider the seminorm

$$\begin{cases} |v|_{*,h}^2 := \sum_{T \in \mathcal{T}_h} \mu_T^2 \{ \|\gamma(v)\|_T^2 + h_T \|\gamma(v) \cdot \mathbf{n}_T\|_{\partial T}^2 + \|\nabla \eta(v)\|_T^2 \}, \\ \text{with } \gamma(v) := \nabla v - \mathbf{G}_{\mathcal{T}}^k(\hat{I}_h^{k,l}(v)), \quad \eta(v) := v - P_{\mathcal{T}}^{k+1}(v). \end{cases} \quad (2.35)$$



Recall that  $P_T^{k+1} := \mathcal{E}_T^{k+1}$  if  $l = k$  and  $P_T^{k+1} := \Pi_T^{k+1}$  if  $l = k+1$ , see Lemma 2.2.6, and  $P_T^{k+1}$  denotes the broken version of  $P_T^{k+1}$  defined elementwise. Notice also that  $\|\gamma(v)\|_T \leq \|\nabla\eta(v)\|_T$  owing to Lemma 2.2.5, but the two terms are kept to better track some terms in the proofs. For a linear form  $\phi \in (\widehat{\mathcal{U}}_{h,0}^{l,k})'$ , we set

$$\|\phi\|_{(\text{HHO})'} := \sup_{\hat{q}_h \in \widehat{\mathcal{U}}_{h,0}^{l,k}} \frac{|\phi(\hat{q}_h)|}{\|\hat{q}_h\|_{\text{HHO}}},$$

with the norm  $\|\cdot\|_{\text{HHO}}$  defined in (2.19).

**Lemma 2.3.1** (Consistency error). *Let  $u \in Y$  be the unique solution to (2.2) and let  $\hat{u}_h \in \widehat{\mathcal{U}}_{h,0}^{l,k}$  be the unique solution to (2.29). Assume  $u \in H^{1+\nu}(\Omega)$ ,  $\nu > \frac{1}{2}$ . We define the HHO error as*

$$\hat{e}_h = (e_{\mathcal{T}}, e_{\mathcal{F}}) := \hat{J}_h(u) - \hat{I}_h^{k,l}(u) \in \widehat{\mathcal{U}}_{h,0}^{l,k}. \quad (2.36)$$

Let us define the consistency error as the linear form  $\psi_h := a_h(\hat{e}_h, \cdot) \in (\widehat{\mathcal{U}}_{h,0}^{l,k})'$ . Then, there is a constant  $c_*$  such that, for all  $h > 0$ ,

$$\|\psi_h\|_{(\text{HHO})'} \leq c_* |u|_{*,h}. \quad (2.37)$$

*Proof.* Let us sketch the proof from [82, 84]. We observe that for all  $\hat{q}_h \in \widehat{\mathcal{U}}_{h,0}^{l,k}$ ,

$$\psi_h(\hat{q}_h) = a_h(\hat{e}_h, \hat{q}_h) = (B(u), q_{\mathcal{T}})_{\Omega} - a_h(\hat{I}_h^{k,l}(u), \hat{q}_h). \quad (2.38)$$

A direct calculation using the definition of the gradient reconstruction operator on the test function shows that

$$\psi_h(\hat{q}_h) = \sum_{T \in \mathcal{T}_h} \mu_T^2 (\gamma(u) \cdot \mathbf{n}_T, q_{\partial T} - q_T)_{\partial T} - s_h(\hat{I}_h^{k,l}(u), \hat{q}_h),$$

with  $\gamma(u)$  defined in (2.35) (notice that we used  $(\gamma(u), \nabla q_T)_T = 0$  since  $\nabla q_T \in \nabla \mathbb{P}_d^l(T) \subset \mathbb{P}_d^k(T; \mathbb{R}^d)$ ). The first term on the right-hand side is bounded by the Cauchy–Schwarz inequality. For the second term, invoking the Cauchy–Schwarz inequality and the bound (2.25) on the stabilization gives

$$\begin{aligned} |s_h(\hat{I}_h^{k,l}(u), \hat{q}_h)| &\leq \sum_{T \in \mathcal{T}_h} s_T(\hat{I}_T^{k,l}(u), \hat{I}_T^{k,l}(u))^{\frac{1}{2}} s_T(\hat{q}_T, \hat{q}_T)^{\frac{1}{2}} \\ &\leq C \sum_{T \in \mathcal{T}_h} \|\nabla\eta(u)\|_T s_T(\hat{q}_T, \hat{q}_T)^{\frac{1}{2}} \\ &\leq C' \|\nabla\eta(u)\|_{\Omega} \|\hat{q}_h\|_{\text{HHO}}. \end{aligned}$$

This proves the claim.  $\square$

## 2.3.2 Approximation error

A direct consequence of Lemma 2.3.1 is a bound on the discrete energy error.

**Corollary 2.3.2** (Discrete energy error). *Assuming  $u \in H^{1+\nu}(\Omega) \cap Y$ , there is  $C$  such that, for all  $h > 0$ ,*

$$\|\hat{J}_h(u) - \hat{I}_h^{k,l}(u)\|_{\text{HHO}} \leq C |u|_{*,h}. \quad (2.39)$$

*Proof.* Invoking the coercivity of  $a_h$  on  $\widehat{\mathcal{U}}_{h,0}^{l,k}$  and Lemma 2.3.1 yields

$$\alpha \|\hat{e}_h\|_{\text{HHO}}^2 \leq a_h(\hat{e}_h, \hat{e}_h) = \psi_h(\hat{e}_h) \leq c_* |u|_{*,h} \|\hat{e}_h\|_{\text{HHO}},$$

which proves the bound (2.39).  $\square$

To estimate the  $L^2$ -norm of the error, we assume that an elliptic regularity pickup is available, i.e. there is  $s \in (\frac{1}{2}, 1]$  and a (non-dimensional) constant  $c_{\text{ell}}$  such that for all  $g \in L^2(\Omega)$ , the unique function  $\zeta_g \in H_0^1(\Omega)$  such that  $B(\zeta_g) = g$  in  $\Omega$  satisfies

$$\|\zeta_g\|_{H^{1+s}(\Omega)} \leq c_{\text{ell}} \mu_b^{-2} \ell_\Omega^2 \|g\|_\Omega. \quad (2.40)$$

(To be dimensionally consistent we set  $\|\zeta_g\|_{H^{1+s}(\Omega)}^2 := \|\zeta_g\|_\Omega^2 + \ell_\Omega^2 \|\nabla \zeta_g\|_\Omega^2 + \ell_\Omega^{2(1+s)} |\zeta_g|_{H^{1+s}(\Omega)}^2$ .) For all  $u \in H^{1+\nu}(\Omega) \cap Y$ , we consider the additional seminorm

$$|u|_{**,h} := |u|_{*,h} + \mu_b^{-2} \ell_\Omega^\delta h^{1-\delta} \|B(u) - \Pi_{\mathcal{T}}^l(B(u))\|_\Omega, \quad (2.41)$$

with  $\delta := 0$  if  $l \geq 1$  and  $\delta := s$  if  $l = 0$ .

**Lemma 2.3.3** (Discrete  $L^2$ -error). *Assuming  $u \in H^{1+\nu}(\Omega) \cap Y$ , there is  $C$  such that for all  $h > 0$ ,*

$$\|J_{\mathcal{T}}(u) - \Pi_{\mathcal{T}}^l(u)\|_\Omega \leq C \ell_\Omega^{1-s} h^s |u|_{**,h}. \quad (2.42)$$

*Proof.* Let  $\zeta \in H_0^1(\Omega)$  be such that  $B(\zeta) = e_{\mathcal{T}}$  in  $\Omega$ . Proceeding as in the proof of [84, Thm. 10] or [82, Thm. 11] yields

$$\|e_{\mathcal{T}}\|_\Omega^2 = (e_{\mathcal{T}}, B(\zeta))_\Omega = \sum_{T \in \mathcal{T}_h} \mu_T^2 \left\{ (\nabla e_T, \nabla \zeta)_T + (e_{\partial T} - e_T, \nabla \zeta \cdot \mathbf{n}_T)_{\partial T} \right\},$$

where we used that  $\nabla \zeta \cdot \mathbf{n}_T$  can be localized to the mesh interfaces and is continuous across them, owing to the elliptic regularity pickup. Let  $\boldsymbol{\xi}(\zeta) := \nabla \zeta - \mathbf{G}_{\mathcal{T}}^k(\hat{I}_h^{k,l}(\zeta))$  and observe that  $(\nabla e_T, \boldsymbol{\xi}(\zeta))_T = 0$  for all  $T \in \mathcal{T}_h$ . Plugging this and the definition (2.27) of  $b_h$  in the previous computation, we obtain

$$\begin{aligned} \|e_{\mathcal{T}}\|_\Omega^2 &= b_h(\hat{e}_h, \hat{I}_h^{k,l}(\zeta)) + \sum_{T \in \mathcal{T}_h} \mu_T^2 (e_{\partial T} - e_T, \boldsymbol{\xi}(\zeta) \cdot \mathbf{n}_T)_{\partial T} \\ &= a_h(\hat{e}_h, \hat{I}_h^{k,l}(\zeta)) - s_h(\hat{e}_h, \hat{I}_h^{k,l}(\zeta)) + \sum_{T \in \mathcal{T}_h} \mu_T^2 (e_{\partial T} - e_T, \boldsymbol{\xi}(\zeta) \cdot \mathbf{n}_T)_{\partial T}. \end{aligned}$$

Since  $a_h(\hat{e}_h, \hat{I}_h^{k,l}(\zeta)) = (B(u), \Pi_{\mathcal{T}}^l(\zeta))_\Omega - a_h(\hat{I}_h^{k,l}(u), \hat{I}_h^{k,l}(\zeta))$  by definition of the discrete solution and since  $(\mu^2 \nabla u, \nabla \zeta)_\Omega = (B(u), \zeta)_\Omega$  by definition of the exact solution, we add  $(\mu^2 \nabla u, \nabla \zeta)_\Omega - (B(u), \zeta)_\Omega = 0$  to the previous result and infer that  $\|e_{\mathcal{T}}\|_\Omega^2 = \epsilon_1 + \epsilon_2 + \epsilon_3$  with

$$\begin{aligned} \epsilon_1 &:= \sum_{T \in \mathcal{T}_h} \mu_T^2 (e_{\partial T} - e_T, \boldsymbol{\xi}(\zeta) \cdot \mathbf{n}_T)_{\partial T} - s_h(\hat{e}_h, \hat{I}_h^{k,l}(\zeta)), \\ \epsilon_2 &:= (\mu^2 \nabla u, \nabla \zeta)_\Omega - a_h(\hat{I}_h^{k,l}(u), \hat{I}_h^{k,l}(\zeta)), \\ \epsilon_3 &:= (B(u), \zeta - \Pi_{\mathcal{T}}^l(\zeta))_\Omega. \end{aligned}$$

Using (6.17), (2.25) and the inequality (2.40) from the elliptic regularity pickup, we infer that

$$\begin{aligned} |\epsilon_1| &\leq C \|\hat{e}_h\|_{\text{HHO}} \mu_\#^2 h^s |\zeta|_{H^{1+s}(\Omega)} \\ &\leq C' \|\hat{e}_h\|_{\text{HHO}} \mu_\#^2 \ell_\Omega^{-1-s} h^s \|\zeta\|_{H^{1+s}(\Omega)} \\ &\leq C'' \|\hat{e}_h\|_{\text{HHO}} c_{\text{ell}} \ell_\Omega^{1-s} h^s \|e_{\mathcal{T}}\|_\Omega \\ &\leq C''' |u|_{*,h} c_{\text{ell}} \ell_\Omega^{1-s} h^s \|e_{\mathcal{T}}\|_\Omega. \end{aligned}$$

where the last bound follows from (2.39). Since  $\epsilon_2 = (\mu^2 (\nabla u - \mathbf{G}_{\mathcal{T}}^k(\hat{I}_h^{k,l}(u))), \boldsymbol{\xi}(\zeta))_\Omega - s_h(\hat{I}_h^{k,l}(u), \hat{I}_h^{k,l}(\zeta))$ , we obtain by invoking similar arguments

$$|\epsilon_2| \leq C |u|_{*,h} c_{\text{ell}} \ell_\Omega^{1-s} h^s \|e_{\mathcal{T}}\|_\Omega.$$

Furthermore, we have  $\epsilon_3 = (B(u) - \Pi_{\mathcal{T}}^l(B(u)), \zeta - \Pi_{\mathcal{T}}^l(\zeta))_{\Omega}$ , and since  $1 + s - \delta = 1 + s$  if  $l \geq 1$  and  $1 + s - \delta = 1$  if  $l = 0$ , we obtain

$$\begin{aligned} \|\zeta - \Pi_{\mathcal{T}}^l(\zeta)\|_{\Omega} &\leq C \ell_{\Omega}^{-1+\delta-s} h^{1+s-\delta} \|\zeta\|_{H^{1+s}(\Omega)} \\ &\leq C' \mu_b^{-2} \ell_{\Omega}^{1+\delta-s} h^{1+s-\delta} \|e_{\mathcal{T}}\|_{\Omega}. \end{aligned}$$

We infer that

$$|\epsilon_3| \leq C \mu_b^{-2} \ell_{\Omega}^{\delta} h^{1-\delta} \|B(u) - \Pi_{\mathcal{T}}^l(B(u))\|_{\Omega} \ell_{\Omega}^{1-s} h^s \|e_{\mathcal{T}}\|_{\Omega}.$$

Finally, the bound (2.42) follows by putting together the above three estimates.  $\square$

**Lemma 2.3.4** (Approximation). *Assume that  $u \in H^{r+1}(\Omega)$  with  $r \in [\nu, k+1]$ ,  $\nu \geq s > \frac{1}{2}$ , (this additional regularity assumption can be localized to the mesh cells). There is  $C$  such that, for all  $h > 0$ ,*

$$|u|_{*,h} \leq C h^r |u|_{H^{r+1}(\Omega)}. \quad (2.43)$$

*Assuming additionally that  $f = B(u) \in H^{r'}(\Omega)$  with  $r' := \max(r-1+\delta, 0)$  (so that  $r' = \max(r-1, 0)$  if  $l \geq 1$  and  $r' = r-1+s$  if  $l = 0$  since  $r-1+s \geq \nu-1+s > 0$ ), we have*

$$|u|_{**,h} \leq C (h^r |u|_{H^{r+1}(\Omega)} + \mu_b^{-2} \ell_{\Omega}^{\delta} h^{1-\delta+r'} |B(u)|_{H^{r'}(\Omega)}). \quad (2.44)$$

*Proof.* The estimate (2.43) results from Lemma 2.3.2 combined with the approximation properties of the  $\mathbf{L}^2$ -orthogonal and elliptic projections. To prove (2.44), we only need to bound the additional term  $\mu_b^{-2} \ell_{\Omega}^{\delta} h^{1-\delta} \|B(u) - \Pi_{\mathcal{T}}^l(B(u))\|_{\Omega}$ . If  $l = 0$ , then  $\delta = s, k = 0$  and  $r' \in (0, 1]$  so that

$$h^{1-\delta} \|B(u) - \Pi_{\mathcal{T}}^0(B(u))\|_{\Omega} \leq C h^{1-\delta+r'} |B(u)|_{H^{r'}(\Omega)},$$

and we obtain the expected estimate. If  $l \geq 1$ , then  $\delta = 0$  and  $l+1 \geq k+1 \geq r'$ , so that

$$h^{1-\delta} \|B(u) - \Pi_{\mathcal{T}}^l(B(u))\|_{\Omega} = h \|B(u) - \Pi_{\mathcal{T}}^l(B(u))\|_{\Omega} \leq C h h^{r'} |B(u)|_{H^{r'}(\Omega)},$$

yielding again the expected estimate.  $\square$

**Remark 2.3.5** (Regularity assumption). If  $\mu$  is smooth (e.g. constant) and  $r \geq 1$ , then  $B(u) \in H^{r-1}(\Omega)$  if  $u \in H^{r+1}(\Omega)$  so that the regularity assumption on  $B(u)$  follows from that on  $u$  whenever  $l \geq 1$ . For  $l = 0$  and full elliptic regularity pickup, the additional regularity assumption is  $f = B(u) \in H^1(\Omega)$ .

## Chapter 3

# Semi-implicit HHO method for the acoustic wave equation

The goal of this chapter is to present the space semi-discretization with the HHO method of the linear acoustic wave equation in its second-order form. Error estimates are derived using energy and duality arguments for the continuous-in-time setting and using energy arguments for the fully discrete setting using the leapfrog (centered finite difference) scheme in time. This chapter also illustrates the difficulty to design an explicit scheme in time with the HHO method. Solutions to overcome this issue are presented in the next chapter.

The linear acoustic wave equation is posed on the space domain  $\Omega \subset \mathbb{R}^d$  and the time interval  $J := [0, \mathfrak{T}]$ , with  $\mathfrak{T} > 0$ . The problem consists in finding  $u : \Omega \times J \rightarrow \mathbb{R}$  such that

$$\partial_t^2 u - \nabla \cdot (\mu^2 \nabla u) = f, \quad \text{in } \Omega \times J, \quad (3.1a)$$

$$u|_{t=0} = u_0, \quad \text{in } \Omega, \quad (3.1b)$$

$$\partial_t u|_{t=0} = v_0, \quad \text{in } \Omega, \quad (3.1c)$$

$$u = 0, \quad \text{in } \partial\Omega \times J, \quad (3.1d)$$

with  $f : \Omega \times J \rightarrow \mathbb{R}$  the source term,  $\mu : \Omega \rightarrow \mathbb{R}_+$  a function representing the speed of sound, which, in the linear case considered in this chapter, is assumed to depend only on the position and for simplicity, we assume as in Chapter 2 that  $\mu$  is piecewise constant on a polyhedral partition of  $\Omega$ . Homogeneous Dirichlet boundary conditions are considered for simplicity, and we assume that the initial data  $u_0, v_0 \in H_0^1(\Omega)$  satisfy these conditions.

Standard notation is used for Bochner spaces. Assuming  $f \in L^2(J; L^2(\Omega))$ , we seek the solution of (3.1) in  $U := H^2(J; L^2(\Omega)) \cap L^2(J; H_0^1(\Omega))$  such that

$$(\partial_t^2 u, w)_\Omega + (\mu^2 \nabla u, \nabla w)_\Omega = (f, w)_\Omega, \quad \text{a.e. } t \in J, \quad \forall w \in H_0^1(\Omega), \quad (3.2)$$

with the initial conditions

$$u(0) = u_0, \quad \partial_t u(0) = v_0. \quad (3.3)$$

The convergence analysis is first done in the space semi-discrete setting, yielding the same optimal convergence rates as in the static case, both for the energy norm and for the  $L^2$ -norm at (say) the final time (Section 3.1). These results are published in [36], where an analysis of the discretization of the first-order formulation in time of the wave equation is also presented. The leapfrog scheme is then used to discretize in time the space semi-discrete equation. The convergence analysis is made in the energy norm at (say) the final time, showing that optimal convergence in space and in time is obtained (Section 3.2). These optimal convergence rates in space and in time are then verified via numerical experiments (Section 3.3). Finally, time discretizations with alternative time schemes are briefly discussed (Section

3.4), followed by a conclusion on the space semi-discretization of the acoustic wave equation with the HHO method (Section 3.5).

We refer the reader to Section 2.2.1 for the local discrete setting, and to Sections 2.2.2 and 2.2.3 for the definition of the HHO operators (gradient reconstruction and stabilization) and the global formulation in space. In this chapter,  $C$  denotes a generic positive constant whose value can change at each occurrence as long as it is independent of the mesh size and the time step. To avoid technicalities with the time stepping and the error analysis, we assume in the entire chapter that  $f \in C^0(\bar{J}; L^2(\Omega))$  with  $\bar{J} = [0, \mathfrak{T}]$ , and that  $u \in C^1(\bar{J}; H^{1+\nu}(\Omega)) \cap C^2(\bar{J}; L^2(\Omega))$  with  $\nu > \frac{1}{2}$ .

### 3.1 Optimal convergence of the HHO space semi-discretization

The results of this section are published in [36]. The main results (Theorems 3.1.1 and 3.1.2) concern the second-order formulation in time of the wave equation for which  $H^1$ - and  $L^2$ -error estimates are established at (say) the final time. The techniques of proof are different from those of [72] for HDG methods. Indeed, the present analysis exploits the primal viewpoint at the core of the HHO method and draws on the error analysis from [94, 18] for continuous finite elements and [107] for dG methods. There are, however, substantial differences with respect to [94, 18, 107]. The key issue is that the HHO method relies on a pair of discrete unknowns, so that it is not possible to proceed simply by extending the discrete bilinear form to an infinite-dimensional functional space that contains the exact solution. This leads us to use the notion of HHO solution map for the steady differential operator in space (see Section 2.3 and, in particular, (2.34) for its definition). We then use this map to perform a suitable error decomposition in the context of the wave equation, in the spirit of the seminal work [150], where an elliptic projector is introduced to derive optimal  $L^2$ -error estimates for the heat equation approximated by continuous finite elements (the same idea is re-used in [18] for the wave equation).

#### 3.1.1 HHO space semi-discretization of the wave equation

We use the HHO method presented in Chapter 2 for the space semi-discretization. Recall that  $\widehat{\mathcal{U}}_{h,0}^{l,k} = \mathcal{U}_{\mathcal{T}}^l \times \mathcal{U}_{\mathcal{F},0}^k$  and that the discrete bilinear form  $a_h$  is defined in (2.28). The space semi-discrete scheme for the linear wave equation (3.1) consists of finding  $\hat{u}_h := (u_{\mathcal{T}}, u_{\mathcal{F}}) \in C^2(\bar{J}; \widehat{\mathcal{U}}_{h,0}^{l,k})$  such that, for all  $t \in \bar{J}$  and all  $\hat{w}_h := (w_{\mathcal{T}}, w_{\mathcal{F}}) \in \widehat{\mathcal{U}}_{h,0}^{l,k}$ ,

$$(\partial_t^2 u_{\mathcal{T}}(t), w_{\mathcal{T}})_{\Omega} + a_h(\hat{u}_h(t), \hat{w}_h) = (f(t), w_{\mathcal{T}})_{\Omega}. \quad (3.4)$$

Notice that the homogeneous Dirichlet boundary condition is enforced by the condition  $\hat{u}_h(t) \in \widehat{\mathcal{U}}_{h,0}^{l,k}$  at all times  $t \in \bar{J}$ . The initial conditions are enforced on the cell degrees of freedom as follows:

$$u_{\mathcal{T}}(0) := \Pi_{\mathcal{T}}^l(u_0), \quad \partial_t u_{\mathcal{T}}(0) := \Pi_{\mathcal{T}}^l(v_0). \quad (3.5)$$

Notice that initial conditions on the face degrees of freedom are not needed.

For all  $t \in \bar{J}$ , let  $(\mathbf{U}_{\mathcal{T}}(t), \mathbf{U}_{\mathcal{F}}(t)) \in \mathbb{R}^{N_{\mathcal{T}}} \times \mathbb{R}^{N_{\mathcal{F}}}$  be the vector of time-dependent degrees of freedom of the solution  $\hat{u}_h(t)$  on bases of  $\mathcal{U}_{\mathcal{T}}^l$  and  $\mathcal{U}_{\mathcal{F},0}^k$ , and  $\mathbf{F}_{\mathcal{T}} \in \mathbb{R}^{N_{\mathcal{T}}}$  the vector having components  $((f(t), \phi_i)_{\Omega})_{1 \leq i \leq N_{\mathcal{T}}}$ . The algebraic form of (3.4) is

$$\begin{bmatrix} \mathcal{M} & 0 \\ 0 & 0 \end{bmatrix} \begin{pmatrix} \partial_{tt} \mathbf{U}_{\mathcal{T}}(t) \\ \cdot \end{pmatrix} + \begin{bmatrix} \mathcal{A}_{\mathcal{T}\mathcal{T}} & \mathcal{A}_{\mathcal{T}\mathcal{F}} \\ \mathcal{A}_{\mathcal{F}\mathcal{T}} & \mathcal{A}_{\mathcal{F}\mathcal{F}} \end{bmatrix} \begin{pmatrix} \mathbf{U}_{\mathcal{T}}(t) \\ \mathbf{U}_{\mathcal{F}}(t) \end{pmatrix} = \begin{bmatrix} \mathbf{F}_{\mathcal{T}}(t) \\ 0 \end{bmatrix}, \quad \forall t \in \bar{J}, \quad (3.6)$$

with  $\mathcal{A}$  the symmetric definite positive stiffness matrix associated with the bilinear form  $a_h$  and  $\mathcal{M}$  the cell mass matrix. As the structure of the global mass matrix makes its definition irrelevant,  $\mathbf{U}_{\mathcal{F}}(t)$  is replaced

by a “.” in the acceleration term in (3.6). Let us note that the cell mass matrix  $\mathcal{M}$  is block-diagonal. Moreover, the second line in (3.6) writes

$$\mathcal{A}_{\mathcal{F}\mathcal{T}}\mathbf{U}_{\mathcal{T}}(t) + \mathcal{A}_{\mathcal{F}\mathcal{F}}\mathbf{U}_{\mathcal{F}}(t) = 0, \quad (3.7)$$

and requires solving a linear system of size  $N_{\mathcal{F}} \times N_{\mathcal{F}}$  independently of the choice of the time integration scheme. This means that even explicit time integration schemes lead to semi-implicit algorithms, with an explicit problem in the cell unknowns and an implicit problem in the face unknowns. One of the main contributions of this Thesis is to address this issue (see Chapter 4).

### 3.1.2 Energy error estimate

Let us start with the energy error estimate. For all  $t \in \bar{J}$ , we consider the seminorm  $|u(t)|_{*,h}$  defined in (2.35) by

$$|u(t)|_{*,h}^2 := \sum_{T \in \mathcal{T}_h} \mu_T^2 \{ \|\gamma(u(t))\|_T^2 + h_T \|\gamma(u(t)) \cdot \mathbf{n}_T\|_{\partial T}^2 + \|\nabla \eta(u(t))\|_T^2 \}, \quad (3.8)$$

with

$$\gamma(u(t)) := \nabla u(t) - \mathbf{G}_{\mathcal{T}}^k(\hat{I}_h^{k,l}(u(t))), \quad \eta(u(t)) := u(t) - P_{\mathcal{T}}^{k+1}(u(t)). \quad (3.9)$$

Recall that  $P_T^{k+1} := \mathcal{E}_T^{k+1}$  if  $l = k$  and  $P_T^{k+1} := \Pi_T^{k+1}$  if  $l = k+1$ , see Lemma 2.2.6, and  $P_{\mathcal{T}}^{k+1}$  denotes the broken version of  $P_T^{k+1}$  defined elementwise. Using  $\gamma(\partial_t u)$  and  $\eta(\partial_t u)$ , we can define  $|\partial_t u(t)|_{*,h}$  similarly, giving for all  $t \in \bar{J}$ ,

$$|\partial_t u(t)|_{*,h}^2 := \sum_{T \in \mathcal{T}_h} \mu_T^2 \{ \|\gamma(\partial_t u(t))\|_T^2 + h_T \|\gamma(\partial_t u(t)) \cdot \mathbf{n}_T\|_{\partial T}^2 + \|\nabla \eta(\partial_t u(t))\|_T^2 \}. \quad (3.10)$$

We also set

$$|u|_{C^0(0,t;*,h)} := \sup_{s \in [0,t]} |u(s)|_{*,h}, \quad |\partial_t u|_{L^1(0,t;*,h)} := \int_0^t |\partial_t u(s)|_{*,h} \, ds.$$

For a function  $\hat{v}_h \in C^0(\bar{J}; \hat{\mathcal{U}}_{h,0}^{l,k})$ , we set, for all  $t \in \bar{J}$ ,

$$\|\hat{v}_h\|_{C^0(0,t;\text{HHO})} := \sup_{s \in [0,t]} \|\hat{v}_h(s)\|_{\text{HHO}}.$$

The following result shows that the energy error converges as  $\mathcal{O}(h^{k+1})$  for smooth solutions.

**Theorem 3.1.1** (Energy-error estimate). *Let  $u$  solve (3.2) with the initial conditions (3.3), and let  $\hat{u}_h$  solve (3.4) with the initial conditions (3.5). Assume that  $u \in C^1(\bar{J}; H^{1+\nu}(\Omega)) \cap C^2(\bar{J}; L^2(\Omega))$  with  $\nu > \frac{1}{2}$ . There is  $C$  such that for all  $h > 0$  and all  $t \in \bar{J}$ ,*

$$\|\partial_t u_{\mathcal{T}} - \Pi_{\mathcal{T}}^l(\partial_t u)\|_{C^0(0,t;L^2(\Omega))} + \|\hat{u}_h - \hat{I}_h^{k,l}(u)\|_{C^0(0,t;\text{HHO})} \leq C \left( |u|_{C^0(0,t;*,h)} + |\partial_t u|_{L^1(0,t;*,h)} \right). \quad (3.11)$$

Moreover, we have

$$\begin{aligned} & \|\partial_t u_{\mathcal{T}} - \partial_t u\|_{C^0(0,t;L^2(\Omega))} + \|\mu(\mathbf{G}_{\mathcal{T}}^k(\hat{u}_h) - \nabla u)\|_{C^0(0,t;L^2(\Omega))} \\ & \leq C \left( |u|_{C^0(0,t;*,h)} + |\partial_t u|_{L^1(0,t;*,h)} \right) + \|\partial_t u - \Pi_{\mathcal{T}}^l(\partial_t u)\|_{C^0(0,t;L^2(\Omega))}. \end{aligned} \quad (3.12)$$

Finally, if there is  $r \in (\frac{1}{2}, k+1]$  so that  $u \in C^1(\bar{J}; H^{r+1}(\Omega))$ , we have

$$\|\partial_t u_{\mathcal{T}} - \partial_t u\|_{C^0(0,t;L^2(\Omega))} + \|\mu(\mathbf{G}_{\mathcal{T}}^k(\hat{u}_h) - \nabla u)\|_{C^0(0,t;L^2(\Omega))} \leq C \mu_{\#} h^r |u|_{C^1(0,t;H^{r+1}(\Omega))}, \quad (3.13)$$

where  $|u|_{C^1(0,t;H^{r+1}(\Omega))} := |u|_{C^0(0,t;H^{r+1}(\Omega))} + t |\partial_t u|_{C^0(0,t;H^{r+1}(\Omega))}$ .

*Proof. Step 1: Error equation.* Let us set, for all  $t \in \bar{J}$ ,

$$\hat{e}_h(t) := \hat{u}_h(t) - \hat{I}_h^{k,l}(u(t)) \in \widehat{\mathcal{U}}_{h,0}^{l,k}.$$

Then (3.11) can be rewritten as

$$\|\partial_t e_{\mathcal{T}}\|_{C^0(0,t;L^2(\Omega))} + \|\hat{e}_h\|_{C^0(0,t;\text{HHO})} \leq C \left( |u|_{C^0(0,t;*,h)} + |\partial_t u|_{L^1(0,t;*,h)} \right). \quad (3.14)$$

Recall the notation  $B(v) = -\nabla \cdot (\mu^2 \nabla v)$  and that the exact solution verifies  $\partial_t^2 u(t) + B(u(t)) = f(t)$  in  $L^2(\Omega)$  for all  $t \in \bar{J}$ , owing to our regularity assumption on  $f$  and  $\partial_t^2 u$ . We observe that for all  $\hat{q}_h \in \widehat{\mathcal{U}}_{h,0}^{l,k}$  and all  $t \in \bar{J}$ ,

$$\begin{aligned} (\partial_t^2 e_{\mathcal{T}}(t), q_{\mathcal{T}})_{\Omega} + a_h(\hat{e}_h(t), \hat{q}_h) &= (f(t), q_{\mathcal{T}})_{\Omega} - (\partial_t^2 \Pi_{\mathcal{T}}^l(u(t)), q_{\mathcal{T}})_{\Omega} - a_h(\hat{I}_h^{k,l}(u(t)), \hat{q}_h) \\ &= (\partial_t^2 u(t) - \partial_t^2 \Pi_{\mathcal{T}}^l(u(t)), q_{\mathcal{T}})_{\Omega} + (B(u(t)), q_{\mathcal{T}})_{\Omega} - a_h(\hat{I}_h^{k,l}(u(t)), \hat{q}_h) \\ &= (B(u(t)), q_{\mathcal{T}})_{\Omega} - a_h(\hat{I}_h^{k,l}(u(t)), \hat{q}_h) =: \psi_h(t; \hat{q}_h), \end{aligned} \quad (3.15)$$

where we used that  $\partial_t^2 \Pi_{\mathcal{T}}^l(u(t)) = \Pi_{\mathcal{T}}^l(\partial_t^2 u(t))$  (so that  $(\partial_t^2 u(t) - \partial_t^2 \Pi_{\mathcal{T}}^l(u(t)), q_{\mathcal{T}})_{\Omega} = 0$ ). The linear form  $\psi_h(t; \cdot) \in (\widehat{\mathcal{U}}_{h,0}^{l,k})'$  represents the consistency error at time  $t$  associated with the HHO space semi-discretization. Recall from Lemma 2.3.1 that

$$\psi_h(t; \hat{q}_h) = \sum_{T \in \mathcal{T}_h} \mu_T^2 (\gamma(u(t)) \cdot \mathbf{n}_T, q_{\partial T} - q_T)_{\partial T} - s_h(\hat{I}_h^{k,l}(u(t)), \hat{q}_h),$$

with  $\gamma(u(t))$  defined in (3.9) and where we used that  $(\gamma(u(t)), \nabla q_T)_T = 0$  for all  $T \in \mathcal{T}_h$ . We also introduce the linear form  $\dot{\psi}_h(t; \cdot) \in (\widehat{\mathcal{U}}_{h,0}^{l,k})'$  such that

$$\dot{\psi}_h(t; \hat{q}_h) := \sum_{T \in \mathcal{T}_h} \mu_T^2 (\gamma(\partial_t u(t)) \cdot \mathbf{n}_T, q_{\partial T} - q_T)_{\partial T} - s_h(\hat{I}_h^{k,l}(\partial_t u(t)), \hat{q}_h), \quad (3.16)$$

and observe that the product rule for the time derivative implies that for all  $\hat{v}_h \in C^1(\bar{J}; \widehat{\mathcal{U}}_{h,0}^{l,k})$ ,

$$\frac{d}{dt} \psi_h(t; \hat{v}_h(t)) = \psi_h(t; \partial_t \hat{v}_h(t)) + \dot{\psi}_h(t; \hat{v}_h(t)). \quad (3.17)$$

*Step 2: Stability argument.* Let us test the error equation with  $\hat{q}_h := \partial_t \hat{e}_h(t)$  for all  $t \in \bar{J}$ . Since the discrete bilinear form  $a_h$  is symmetric and using (3.17) on the right-hand side leads to

$$\frac{d}{dt} \left\{ \frac{1}{2} \|\partial_t e_{\mathcal{T}}\|_{\Omega}^2 + \frac{1}{2} a_h(\hat{e}_h(t), \hat{e}_h(t)) \right\} = \frac{d}{dt} \psi_h(t; \hat{e}_h(t)) - \dot{\psi}_h(t; \hat{e}_h(t)).$$

Integrating in time from 0 to  $t$ , observing that  $\partial_t e_{\mathcal{T}}(0) = 0$  owing to the initial conditions (we also have  $e_{\mathcal{T}}(0) = 0$  but in general  $e_{\mathcal{F}}(0) \neq 0$ , see below), and using the coercivity and the continuity of the discrete bilinear form  $a_h$  (see Lemma 2.2.3), we infer that

$$\begin{aligned} \frac{1}{2} \|\partial_t e_{\mathcal{T}}(t)\|_{\Omega}^2 + \frac{1}{2} \alpha \|\hat{e}_h(t)\|_{\text{HHO}}^2 &\leq \frac{1}{2} \|\partial_t e_{\mathcal{T}}(t)\|_{\Omega}^2 + \frac{1}{2} a_h(\hat{e}_h(t), \hat{e}_h(t)) \\ &= \frac{1}{2} a_h(\hat{e}_h(0), \hat{e}_h(0)) + \int_0^t \left\{ \frac{d}{ds} \psi_h(s; \hat{e}_h(s)) - \dot{\psi}_h(s; \hat{e}_h(s)) \right\} ds \\ &\leq \frac{1}{2} \varpi \|\hat{e}_h(0)\|_{\text{HHO}}^2 - \psi_h(0; \hat{e}_h(0)) + \psi_h(t; \hat{e}_h(t)) - \int_0^t \dot{\psi}_h(s; \hat{e}_h(s)) ds. \end{aligned}$$

We rewrite this inequality as  $\frac{1}{2} \|\partial_t e_{\mathcal{T}}(t)\|_{\Omega}^2 + \frac{1}{2} \alpha \|\hat{e}_h(t)\|_{\text{HHO}}^2 \leq \frac{1}{2} \varpi \|\hat{e}_h(0)\|_{\text{HHO}}^2 - A_2 + A_3 - A_4$  with

$$\begin{aligned} A_2 &:= \psi_h(0; \hat{e}_h(0)), \\ A_3 &:= \psi_h(t; \hat{e}_h(t)), \\ A_4 &:= \int_0^t \dot{\psi}_h(s; \hat{e}_h(s)) ds. \end{aligned}$$

Using Young's inequality with coefficient  $\frac{\alpha}{2}$  ( $\alpha$  resulting from Lemma 2.2.3), we obtain

$$\begin{aligned} |A_2| &\leq |\psi_h|_{C^0(0,t;(\text{HHO})')} \|\hat{e}_h(0)\|_{\text{HHO}} \leq \frac{1}{\alpha} |\psi_h|_{C^0(0,t;(\text{HHO})')}^2 + \frac{\alpha}{4} \|\hat{e}_h(0)\|_{\text{HHO}}^2, \\ |A_3| &\leq |\psi_h|_{C^0(0,t;(\text{HHO})')} \|\hat{e}_h(t)\|_{\text{HHO}} \leq \frac{1}{\alpha} |\psi_h|_{C^0(0,t;(\text{HHO})')}^2 + \frac{\alpha}{4} \|\hat{e}_h(t)\|_{\text{HHO}}^2. \end{aligned} \quad (3.18)$$

In the same spirit, using Young's inequality with coefficient  $\frac{\alpha}{4}$  as well as Hölder's inequality, we have

$$|A_4| \leq |\dot{\psi}_h|_{L^1(0,t;(\text{HHO})')} \|\hat{e}_h\|_{C^0(0,t;\text{HHO})} \leq \frac{2}{\alpha} |\dot{\psi}_h|_{C^0(0,t;(\text{HHO})')}^2 + \frac{\alpha}{8} \|\hat{e}_h\|_{C^0(0,t;\text{HHO})}^2.$$

Gathering these three estimates, simplifying by  $\frac{\alpha}{4} \|\hat{e}_h(0)\|_{\text{HHO}}^2$  on each side and using the fact that  $\frac{\varpi}{2} + \frac{\alpha}{4} \leq \frac{3}{4} \varpi$  (from the definitions of  $\alpha$  and  $\varpi$  in Lemma 2.2.3) gives

$$\begin{aligned} \frac{1}{2} \|\partial_t e_{\mathcal{T}}(t)\|_{\Omega}^2 + \frac{1}{4} \alpha \|\hat{e}_h(t)\|_{\text{HHO}}^2 &\leq \frac{2}{\alpha} (|\psi_h|_{C^0(0,t;(\text{HHO})')}^2 + |\dot{\psi}_h|_{L^1(0,t;(\text{HHO})')}^2) \\ &\quad + \frac{1}{8} \alpha \|\hat{e}_h\|_{C^0(0,t;\text{HHO})}^2 + \frac{3}{4} \varpi \|\hat{e}_h(0)\|_{\text{HHO}}^2. \end{aligned}$$

Since the left-hand side evaluated at any  $t' \in [0, t]$  is bounded by the right-hand side, in particular at the time when the maximum over  $[0, t]$  of  $\|\hat{e}_h(t')\|_{\text{HHO}}^2$  is reached, we infer that

$$\|\partial_t e_{\mathcal{T}}\|_{C^0(0,t;L^2(\Omega))}^2 + \|\hat{e}_h\|_{C^0(0,t;\text{HHO})}^2 \leq C (|\psi_h|_{C^0(0,t;(\text{HHO})')}^2 + |\dot{\psi}_h|_{L^1(0,t;(\text{HHO})')}^2 + \|\hat{e}_h(0)\|_{\text{HHO}}^2).$$

Taking the square root gives

$$\|\partial_t e_{\mathcal{T}}\|_{C^0(0,t;L^2(\Omega))} + \|\hat{e}_h\|_{C^0(0,t;\text{HHO})} \leq C (|\psi_h|_{C^0(0,t;(\text{HHO})')} + |\dot{\psi}_h|_{L^1(0,t;(\text{HHO})')} + \|\hat{e}_h(0)\|_{\text{HHO}}). \quad (3.19)$$

*Step 3: Bound on consistency error.* Owing to the proof of Lemma 2.3.2, there is  $c_*$  such that for all  $h > 0$  and all  $t \in \bar{J}$ ,

$$\|\psi_h(t; \cdot)\|_{(\text{HHO})'} \leq c_* |u(t)|_{*,h}, \quad \|\dot{\psi}_h(t; \cdot)\|_{(\text{HHO})'} \leq c_* |\partial_t u(t)|_{*,h}. \quad (3.20)$$

Using these bounds in (3.19) gives

$$\|\partial_t e_{\mathcal{T}}\|_{C^0(0,t;L^2(\Omega))} + \|\hat{e}_h\|_{C^0(0,t;\text{HHO})} \leq C (|u|_{C^0(0,t;*,h)} + |\partial_t u|_{L^1(0,t;*,h)} + \|\hat{e}_h(0)\|_{\text{HHO}}).$$

*Step 4: Bound on initial error.* Since  $e_{\mathcal{T}}(0) = 0$ , the coercivity of the discrete bilinear form  $a_h$  implies that

$$\alpha \|\hat{e}_h(0)\|_{\text{HHO}}^2 = \alpha \|(0, e_{\mathcal{F}}(0))\|_{\text{HHO}}^2 \leq a_h((0, e_{\mathcal{F}}(0)), (0, e_{\mathcal{F}}(0))),$$

with  $e_{\mathcal{F}}(0) = u_{\mathcal{F}}(0) - \Pi_{\mathcal{F}}^k(u_0)$ . The space semi-discrete equation (3.4) and the linearity of  $a_h$  with respect to its first argument imply that

$$a_h((0, u_{\mathcal{F}}(0)), (0, e_{\mathcal{F}}(0))) = -a_h((u_{\mathcal{T}}(0), 0), (0, e_{\mathcal{F}}(0))) = -a_h((\Pi_{\mathcal{T}}^l(u_0), 0), (0, e_{\mathcal{F}}(0))).$$

where we used (3.5) in the last equality. Hence, we have

$$\begin{aligned} \alpha \|\hat{e}_h(0)\|_{\text{HHO}}^2 &= \alpha \|(0, e_{\mathcal{F}}(0))\|_{\text{HHO}}^2 \leq a_h((0, u_{\mathcal{F}}(0) - \Pi_{\mathcal{F}}^k(u_0)), (0, e_{\mathcal{F}}(0))) \\ &= -a_h((\Pi_{\mathcal{T}}^l(u_0), 0), (0, e_{\mathcal{F}}(0))) - a_h((0, \Pi_{\mathcal{F}}^k(u_0)), (0, e_{\mathcal{F}}(0))) \\ &= -a_h(\hat{I}_h^{k,l}(u_0), (0, e_{\mathcal{F}}(0))) \\ &= \psi_h(0, (0, e_{\mathcal{F}}(0))) \leq c_* |u_0|_{*,h} \|(0, e_{\mathcal{F}}(0))\|_{\text{HHO}}, \end{aligned}$$

where we used the consistency bound (3.20) (notice that  $u_0 \in Y \cap H^{1+\nu}(\Omega)$  by assumption). This implies that

$$\|(0, e_{\mathcal{F}}^0)\|_{\text{HHO}} \leq C |u_0|_{*,h}. \quad (3.21)$$



Putting the above estimates together proves the error bound (3.14) and thus (3.11).

*Step 5: Estimate (3.12).* The estimate (3.12) follows from (3.11) after invoking the triangle inequality and observing that for all  $t \in \bar{J}$ ,

$$\begin{aligned} \|\mu(\mathbf{G}_{\mathcal{T}}^k(\hat{u}_h(t)) - \nabla u(t))\|_{\Omega} &\leq \|\mu \mathbf{G}_{\mathcal{T}}^k(\hat{e}_h(t))\|_{\Omega} + \|\mu \gamma(u(t))\|_{\Omega} \\ &\leq \varpi^{\frac{1}{2}} \|\hat{e}_h(t)\|_{\text{HHO}} + |u(t)|_{*,h}, \end{aligned}$$

owing to the upper bound from Lemma 2.2.3 applied to  $\hat{e}_h(t)$  for the first term and the definition of the  $|\cdot|_{*,h}$ -seminorm for the second term. This readily gives (3.12).

*Step 6: Convergence rate (3.13).* Finally, the estimate (3.13) follows from (3.12) after invoking the approximation property (2.43) and bounding the  $L^1$ -norm over  $(0, t)$  by  $t$  times the  $C^0$ -norm over  $[0, t]$ .  $\square$

### 3.1.3 $L^2$ -error estimate

We now establish an improved  $L^2$ -error estimate. For all  $t \in \bar{J}$ , we consider the seminorm  $|u(t)|_{**,h}$  defined in (2.41) by

$$|u(t)|_{**,h} := |u(t)|_{*,h} + \mu_b^{-2} \ell_{\Omega}^{\delta} h^{1-\delta} \|B(u(t)) - \Pi_{\mathcal{T}}^l(B(u(t)))\|_{\Omega},$$

as well as

$$\begin{aligned} |\partial_t u(t)|_{**,h} &:= \left\{ \sum_{T \in \mathcal{T}_h} \mu_T^2 \{ \|\gamma(\partial_t u(t))\|_T^2 + h_T \|\gamma(\partial_t u(t)) \cdot \mathbf{n}_T\|_{\partial T}^2 + \|\nabla \partial_t \eta(u(t))\|_T^2 \} \right\}^{\frac{1}{2}} \\ &\quad + \ell_{\Omega}^{\delta} h^{1-\delta} \|B(\partial_t u(t)) - \Pi_{\mathcal{T}}^l(B(\partial_t u(t)))\|_{\Omega}, \end{aligned}$$

where  $\delta = 0$  if  $l \geq 1$ ,  $\delta = s$  if  $l = 0$  and  $s \in (\frac{1}{2}, 1]$  is the index of elliptic regularity pickup (see Section 2.3.2). We also set

$$|u|_{C^0(0,t;**,h)} := \sup_{s \in [0,t]} |u(s)|_{**,h}, \quad |\partial_t u|_{L^1(0,t;**,h)} := \int_0^t |\partial_t u(s)|_{**,h} \, ds \quad \text{for all } t \in \bar{J}.$$

The following result shows that the  $C^0(0, t; L^2(\Omega))$ -error converges as  $\mathcal{O}(h^{k+2})$  for smooth solutions and if full elliptic regularity pickup ( $s = 1$ ) is available.

**Theorem 3.1.2** ( $L^2$ -error estimate). *Let  $u$  solve (3.2) with the initial conditions (3.3), and let  $\hat{u}_h$  solve (3.4) with the initial conditions (3.5). Assume that  $u \in C^1(\bar{J}; H^{1+\nu}(\Omega)) \cap C^2(\bar{J}; L^2(\Omega))$  with  $\nu > \frac{1}{2}$ . Assume that there is an elliptic regularity pickup with index  $s \in (\frac{1}{2}, 1]$ . There is  $C$  such that for all  $h > 0$  and all  $t \in \bar{J}$ ,*

$$\|u_{\mathcal{T}} - \Pi_{\mathcal{T}}^l(u)\|_{C^0(0,t;L^2(\Omega))} \leq C \ell_{\Omega}^{1-s} h^s (|u|_{C^0(0,t;**,h)} + |\partial_t u|_{L^1(0,t;**,h)}). \quad (3.22)$$

Moreover if there is  $r \in (\frac{1}{2}, k+1]$  so that  $u \in C^1(\bar{J}; H^{r+1}(\Omega))$ , and assuming additionally that  $B(u) \in C^1(\bar{J}; H^{r'}(\Omega))$  with  $r' = \max(r-1+\delta, 0)$  and  $\delta := 0$  if  $l \geq 1$  and  $\delta := s$  otherwise, we have

$$\|u_{\mathcal{T}} - \Pi_{\mathcal{T}}^l(u)\|_{C^0(0,t;L^2(\Omega))} \leq C \ell_{\Omega}^{1-s} h^{r+s} (|u|_{C^1(0,t;H^{r+1}(\Omega))} + \ell_{\Omega}^{\delta} |B(u)|_{C^1(0,t;H^{r'}(\Omega))}). \quad (3.23)$$

*Proof. Step 1: Error equation.* We consider a different error decomposition than in Theorem 3.1.1, i.e we now set for all  $t \in \bar{J}$ ,

$$\hat{e}_h(t) := \hat{u}_h(t) - \hat{J}_h(u(t)),$$

where  $\hat{J}_h$  is the HHO solution map defined in (2.34) for the static problem. We infer that

$$\begin{aligned} (\partial_t^2 e_{\mathcal{T}}(t), q_{\mathcal{T}})_{\Omega} + a_h(\hat{e}_h(t), \hat{q}_h) &= (f(t), q_{\mathcal{T}})_{\Omega} - (\partial_t^2 J_{\mathcal{T}}(u(t)), q_{\mathcal{T}})_{\Omega} - a_h(\hat{J}_h(u(t)), \hat{q}_h) \\ &= (\partial_t^2 u(t) - \partial_t^2 J_{\mathcal{T}}(u(t)), q_{\mathcal{T}})_{\Omega} + (B(u(t)), q_{\mathcal{T}})_{\Omega} - a_h(\hat{J}_h(u(t)), \hat{q}_h) \\ &= (\partial_t^2 \Pi_{\mathcal{T}}^l(u(t)) - \partial_t^2 J_{\mathcal{T}}(u(t)), q_{\mathcal{T}})_{\Omega} =: (\partial_t^2 \theta(t), q_{\mathcal{T}})_{\Omega}, \end{aligned}$$

with  $\theta(t) := \Pi_{\mathcal{T}}^l(u(t)) - J_{\mathcal{T}}(u(t))$  for all  $t \in \bar{J}$ .

*Step 2: Stability argument.* We adapt an argument from [72]. Let  $t \in \bar{J}$  and let  $\chi \in [0, t]$ . Let us set  $\hat{z}_h(t) := -\int_{\chi}^t \hat{e}_h(s) ds$  so that  $\partial_t \hat{z}_h(t) = -\hat{e}_h(t)$ . Testing the above error equation with  $\hat{q}_h := \hat{z}_h(t)$  for all  $t \in \bar{J}$ , integrating by parts in time, and using the symmetry of the discrete bilinear form  $a_h$ , we infer that

$$\frac{d}{dt} \left\{ (\partial_t e_{\mathcal{T}}(t), z_{\mathcal{T}}(t))_{\Omega} + \frac{1}{2} \|e_{\mathcal{T}}(t)\|_{\Omega}^2 - \frac{1}{2} a_h(\hat{z}_h(t), \hat{z}_h(t)) \right\} = \frac{d}{dt} (\partial_t \theta(t), z_{\mathcal{T}}(t))_{\Omega} + (\partial_t \theta(t), e_{\mathcal{T}}(t))_{\Omega}.$$

Integrating this identity in time from 0 to  $\chi$  and since  $\hat{z}_h(\chi) = 0$ , we infer that

$$\frac{1}{2} \|e_{\mathcal{T}}(\chi)\|_{\Omega}^2 + \frac{1}{2} a_h(\hat{z}_h(0), \hat{z}_h(0)) = \frac{1}{2} \|e_{\mathcal{T}}(0)\|_{\Omega}^2 - (\partial_t(\theta(0) - e_{\mathcal{T}}(0)), z_{\mathcal{T}}(0))_{\Omega} + \int_0^{\chi} (\partial_t \theta(s), e_{\mathcal{T}}(s))_{\Omega} ds.$$

Since  $\theta(0) - e_{\mathcal{T}}(0) = \Pi_{\mathcal{T}}^l(u_0) - u_{\mathcal{T}}(0) = 0$ ,  $\partial_t(\theta(0) - e_{\mathcal{T}}(0)) = \Pi_{\mathcal{T}}^l(v_0) - \partial_t u_{\mathcal{T}}(0) = 0$ , and  $a_h(\hat{z}_h(0), \hat{z}_h(0)) \geq 0$ , we obtain

$$\frac{1}{2} \|e_{\mathcal{T}}(\chi)\|_{\Omega}^2 \leq \frac{1}{2} \|\theta(0)\|_{\Omega}^2 + \int_0^{\chi} (\partial_t \theta(s), e_{\mathcal{T}}(s))_{\Omega} ds.$$

Invoking Hölder's inequality for the last term on the right-hand side followed by Young's inequality yields

$$\frac{1}{4} \|e_{\mathcal{T}}(\chi)\|_{\Omega}^2 \leq \frac{1}{2} \|\theta(0)\|_{\Omega}^2 + \|\partial_t \theta\|_{L^1(0, \chi; L^2(\Omega))}^2, \quad \forall \chi \in [0, t].$$

Considering the time  $\chi \in [0, t]$  such that  $\|e_{\mathcal{T}}(\chi)\|_{\Omega}$  is maximal gives

$$\frac{1}{4} \|e_{\mathcal{T}}\|_{C^0(0, t; L^2(\Omega))}^2 \leq \frac{1}{2} \|\theta(0)\|_{\Omega}^2 + \|\partial_t \theta\|_{L^1(0, t; L^2(\Omega))}^2.$$

Since  $u_{\mathcal{T}} - \Pi_{\mathcal{T}}^l(u) = e_{\mathcal{T}} - \theta$ , invoking the triangle inequality we conclude that

$$\|u_{\mathcal{T}} - \Pi_{\mathcal{T}}^l(u)\|_{C^0(0, t; L^2(\Omega))} \leq C(\|\theta\|_{C^0(0, t; L^2(\Omega))} + \|\partial_t \theta\|_{L^1(0, t; L^2(\Omega))}).$$

*Step 3: Bound on consistency error.* Since  $\partial_t \theta = \Pi_{\mathcal{T}}^l(\partial_t u) - J_{\mathcal{T}}(\partial_t u)$ , we can invoke (2.42) to infer that (3.22) holds true. Finally, (3.23) follows from (3.22) and (2.44).  $\square$

## 3.2 Full discretization with leapfrog scheme

The full discretization of the acoustic wave equation in its second-order form with the leapfrog scheme in time was introduced and studied numerically in [34], but without providing a convergence analysis. The goal of this study is to fill this gap. It is well known that the time discretization with the leapfrog scheme combined with finite elements in space leads to a second-order conditionally stable scheme which preserves a modified energy. Our first finding is that it is possible to generalize this idea to the HHO space semi-discretization. The present proof of convergence relies on this idea, while adapting the arguments presented in the previous section for the space semi-discrete case: the error equation is obtained by testing the fully discrete problem with a discrete speed error, integration by parts to obtain stability is mimicked by means of a discrete sum reorganization, and the bound on the initial error relies on similar arguments. The fully discrete proof requires, however, a few more technical steps than the one for the space semi-discrete case.

### 3.2.1 Time discretization with leapfrog scheme

Let  $N$  be the number of discrete time intervals such that  $(t^n)_{n \in \{0: N\}}$  are the discrete time nodes with  $t^0 = 0$  and  $t^N := \mathfrak{T}$ . We set  $f^n := f(t^n)$  for all  $n \in \{0: N\}$ . For the sake of simplicity, we consider a fixed time step  $\Delta t := \frac{\mathfrak{T}}{N}$ . The time discrete unknown  $\hat{u}_h^n = (u_{\mathcal{T}}^n, u_{\mathcal{F}}^n) \in \hat{\mathcal{U}}_{h,0}^{l,k}$  is meant to be an approximation of the space semi-discrete HHO solution  $\hat{u}_h(t^n) \in \hat{\mathcal{U}}_{h,0}^{l,k}$ .

The leapfrog scheme consists of solving, for all  $n \in \{1:N-1\}$ ,

$$\frac{1}{\Delta t^2}(u_{\mathcal{T}}^{n+1} - 2u_{\mathcal{T}}^n + u_{\mathcal{T}}^{n-1}, w_{\mathcal{T}})_{\Omega} + a_h(\hat{u}_h^n, \hat{w}_h) = (f^n, w_{\mathcal{T}})_{\Omega}, \quad \forall \hat{w}_h \in \hat{\mathcal{U}}_{h,0}^{l,k}, \quad (3.24)$$

where the unknowns are  $u_{\mathcal{T}}^{n+1}$  and  $u_{\mathcal{F}}^n$ , whereas  $u_{\mathcal{T}}^n$  and  $u_{\mathcal{T}}^{n-1}$  are known from prior time steps or given by the initial conditions as follows:

$$u_{\mathcal{T}}^0 = \Pi_{\mathcal{T}}^l(u_0), \quad (3.25a)$$

$$a_h(\hat{u}_h^0, (0, w_{\mathcal{F}})) = 0, \quad \forall w_{\mathcal{F}} \in \mathcal{U}_{\mathcal{F},0}^k, \quad (3.25b)$$

$$(u_{\mathcal{T}}^1, w_{\mathcal{T}})_{\Omega} = (u_{\mathcal{T}}^0 + \Delta t \Pi_{\mathcal{T}}^l(v_0), w_{\mathcal{T}})_{\Omega} + \frac{\Delta t^2}{2} [(f^0, w_{\mathcal{T}})_{\Omega} - a_h(\hat{u}_h^0, (w_{\mathcal{T}}, 0))], \quad \forall w_{\mathcal{T}} \in \mathcal{U}_{\mathcal{T}}^l. \quad (3.25c)$$

Notice that we used the initial conditions (3.5) in (6.43a) and (6.43c).

At each time step  $n \in \{1:N-1\}$ , the problem (3.24) is solved by first finding the face unknown  $u_{\mathcal{F}}^n \in \mathcal{U}_{\mathcal{F},0}^k$  from the cell unknown  $u_{\mathcal{T}}^n \in \mathcal{U}_{\mathcal{T}}^l$ :

$$a_h((0, u_{\mathcal{F}}^n), (0, w_{\mathcal{F}})) = -a_h((u_{\mathcal{T}}^n, 0), (0, w_{\mathcal{F}})), \quad \forall w_{\mathcal{F}} \in \mathcal{U}_{\mathcal{F},0}^k, \quad (3.26)$$

and then the cell unknown  $u_{\mathcal{T}}^{n+1} \in \mathcal{U}_{\mathcal{T}}^l$  is computed by solving

$$\frac{1}{\Delta t^2}(u_{\mathcal{T}}^{n+1} - 2u_{\mathcal{T}}^n + u_{\mathcal{T}}^{n-1}, w_{\mathcal{T}})_{\Omega} = (f^n, w_{\mathcal{T}})_{\Omega} - a_h(\hat{u}_h^n, (w_{\mathcal{T}}, 0)), \quad \forall w_{\mathcal{T}} \in \mathcal{U}_{\mathcal{T}}^l. \quad (3.27)$$

Owing to the linear solve implied by (3.26), the scheme is semi-implicit.

**Remark 3.2.1** (Final step). At the final step  $n = N-1$ , we compute  $u_{\mathcal{F}}^{N-1}$  from (3.26) and  $u_{\mathcal{T}}^N$  from (3.27). Then  $u_{\mathcal{F}}^N$  can be retrieved by solving (3.26) for  $n = N$ .

In algebraic form, equations (3.26) and (3.27) translate into the hybrid semi-implicit scheme

$$\frac{1}{\Delta t^2} \begin{bmatrix} \mathcal{M} & 0 \\ 0 & 0 \end{bmatrix} \begin{pmatrix} \mathbf{U}_{\mathcal{T}}^{n+1} - 2\mathbf{U}_{\mathcal{T}}^n + \mathbf{U}_{\mathcal{T}}^{n-1} \\ \cdot \end{pmatrix} + \begin{bmatrix} \mathcal{A}_{\mathcal{T}\mathcal{T}} & \mathcal{A}_{\mathcal{T}\mathcal{F}} \\ \mathcal{A}_{\mathcal{F}\mathcal{T}} & \mathcal{A}_{\mathcal{F}\mathcal{F}} \end{bmatrix} \begin{pmatrix} \mathbf{U}_{\mathcal{T}}^n \\ \mathbf{U}_{\mathcal{F}}^n \end{pmatrix} = \begin{bmatrix} \mathbf{F}_{\mathcal{T}}^n \\ 0 \end{bmatrix}, \quad (3.28)$$

using the same notation as in (3.6) and  $\mathbf{F}_{\mathcal{T}}^n := \mathbf{F}_{\mathcal{T}}(t^n)$ . Thus, the resolution of the fully discrete wave equation proceeds in two steps for all  $n \in \{1:N-1\}$ :

1. Compute  $\mathbf{U}_{\mathcal{F}}^n$  by solving the linear system  $\mathcal{A}_{\mathcal{F}\mathcal{F}}\mathbf{U}_{\mathcal{F}}^n = -\mathcal{A}_{\mathcal{F}\mathcal{T}}\mathbf{U}_{\mathcal{T}}^n$ , where, as above,  $\mathbf{U}_{\mathcal{T}}^n$  is the data and  $\mathbf{U}_{\mathcal{F}}^n$  the unknown. This requires the inversion of the sparse matrix  $\mathcal{A}_{\mathcal{F}\mathcal{F}}$ , either by a direct inversion or by an iterative process.
2. Compute  $\mathbf{U}_{\mathcal{T}}^{n+1}$  by solving  $\mathcal{M}\mathbf{U}_{\mathcal{T}}^{n+1} = \mathcal{M}(2\mathbf{U}_{\mathcal{T}}^n - \mathbf{U}_{\mathcal{T}}^{n-1} + \Delta t^2 \mathbf{F}_{\mathcal{T}}^n) - \Delta t^2(\mathcal{A}_{\mathcal{T}\mathcal{T}}\mathbf{U}_{\mathcal{T}}^n + \mathcal{A}_{\mathcal{T}\mathcal{F}}\mathbf{U}_{\mathcal{F}}^n)$ . Since the mass matrix  $\mathcal{M}$  is block-diagonal, this linear system is very easy to solve.

For the first time step  $n = 0$ , the initial cell unknown  $\mathbf{U}_{\mathcal{T}}^0$  is computed via the algebraic counterpart of (6.43a), the initial face unknown  $\mathbf{U}_{\mathcal{F}}^0$  via the algebraic counterpart of (6.43b):

$$\mathbf{U}_{\mathcal{F}}^0 = -\mathcal{A}_{\mathcal{F}\mathcal{F}}^{-1}\mathcal{A}_{\mathcal{F}\mathcal{T}}\mathbf{U}_{\mathcal{T}}^0, \quad (3.29)$$

and the cell unknown at the first time step,  $\mathbf{U}_{\mathcal{T}}^1$ , solves

$$\mathcal{M}\mathbf{U}_{\mathcal{T}}^1 = \mathcal{M}(\mathbf{U}_{\mathcal{T}}^0 + \Delta t \mathbf{V}_{\mathcal{T}}^0 + \frac{\Delta t^2}{2} \mathbf{F}_{\mathcal{T}}^0) - \frac{\Delta t^2}{2}(\mathcal{A}_{\mathcal{T}\mathcal{T}}\mathbf{U}_{\mathcal{T}}^0 + \mathcal{A}_{\mathcal{T}\mathcal{F}}\mathbf{U}_{\mathcal{F}}^0), \quad (3.30)$$

with  $\mathbf{V}_{\mathcal{T}}^0$  is the vector of degrees of freedom of  $v_{\mathcal{T}}^0$  computed via the algebraic counterpart of (3.5).

**Remark 3.2.2** (Discrete speed and acceleration). The scheme can also be expressed in terms of discrete speed and acceleration, both quantities being defined on the cells for all  $n \in \{1:N-1\}$  as follows:

$$v_{\mathcal{T}}^n = \frac{1}{2\Delta t}(u_{\mathcal{T}}^{n+1} - u_{\mathcal{T}}^{n-1}) \in \mathcal{U}_{\mathcal{T}}^l, \quad a_{\mathcal{T}}^n = \frac{1}{\Delta t^2}(u_{\mathcal{T}}^{n+1} - 2u_{\mathcal{T}}^n + u_{\mathcal{T}}^{n-1}) \in \mathcal{U}_{\mathcal{T}}^l, \quad (3.31)$$

which leads to the following identities:

$$v_{\mathcal{T}}^n = v_{\mathcal{T}}^{n-1} + \frac{\Delta t}{2}(a_{\mathcal{T}}^n + a_{\mathcal{T}}^{n-1}), \quad u_{\mathcal{T}}^{n+1} = u_{\mathcal{T}}^n + \Delta t v_{\mathcal{T}}^n + \frac{\Delta t^2}{2}a_{\mathcal{T}}^n. \quad (3.32)$$

Moreover, (3.24) becomes

$$(a_{\mathcal{T}}^n, w_{\mathcal{T}})_{\Omega} + a_h(\hat{u}_h^n, \hat{w}_h) = (f^n, w_{\mathcal{T}})_{\Omega}, \quad \forall \hat{w}_h \in \widehat{\mathcal{U}}_{h,0}^{l,k}. \quad (3.33)$$

The algebraic equation (3.28) rewrites

$$\begin{bmatrix} \mathcal{M} & 0 \\ 0 & 0 \end{bmatrix} \begin{pmatrix} \mathbf{A}_{\mathcal{T}}^n \\ \cdot \end{pmatrix} + \begin{bmatrix} \mathcal{A}_{\mathcal{T}\mathcal{T}} & \mathcal{A}_{\mathcal{T}\mathcal{F}} \\ \mathcal{A}_{\mathcal{F}\mathcal{T}} & \mathcal{A}_{\mathcal{F}\mathcal{F}} \end{bmatrix} \begin{pmatrix} \mathbf{U}_{\mathcal{T}}^n \\ \mathbf{U}_{\mathcal{F}}^n \end{pmatrix} = \begin{bmatrix} \mathbf{F}_{\mathcal{T}}^n \\ 0 \end{bmatrix}. \quad (3.34)$$

where  $\mathbf{A}_{\mathcal{T}}^n$  is the algebraic counterpart of  $a_{\mathcal{T}}^n$ , and the resolution is made in the following two steps:

1. Compute  $\mathbf{U}_{\mathcal{F}}^n$  by solving the linear system  $\mathcal{A}_{\mathcal{F}\mathcal{F}}\mathbf{U}_{\mathcal{F}}^n = -\mathcal{A}_{\mathcal{F}\mathcal{T}}\mathbf{U}_{\mathcal{T}}^n$  as above.
2. Solve  $\mathcal{M}\mathbf{A}_{\mathcal{T}}^n = \mathbf{F}_{\mathcal{T}}^n - (\mathcal{A}_{\mathcal{T}\mathcal{T}}\mathbf{U}_{\mathcal{T}}^n + \mathcal{A}_{\mathcal{T}\mathcal{F}}\mathbf{U}_{\mathcal{F}}^n)$  and set  $\mathbf{U}_{\mathcal{T}}^{n+1} = \mathbf{U}_{\mathcal{T}}^n + \Delta t\mathbf{V}_{\mathcal{T}}^n + \frac{\Delta t^2}{2}\mathbf{A}_{\mathcal{T}}^n$  with  $\mathbf{V}_{\mathcal{T}}^n = \mathbf{V}_{\mathcal{T}}^{n-1} + \frac{\Delta t}{2}(\mathbf{A}_{\mathcal{T}}^n + \mathbf{A}_{\mathcal{T}}^{n-1})$ .

The initial values  $\mathbf{U}_{\mathcal{T}}^0$  and  $\mathbf{V}_{\mathcal{T}}^0$  are computed as above using the algebraic counterpart of (3.5) and  $\mathbf{U}_{\mathcal{F}}^0$  is computed using (3.29).  $\square$

**Energy balance.** The notion of discrete energy is central to the proof of convergence of the fully discrete scheme. Let us first introduce the semi-discrete energy balance as in [34, Lem. 3.1]. Recalling that  $\hat{u}_h \in C^2(\bar{J}; \widehat{\mathcal{U}}_{h,0}^{l,k})$  solves the space semi-discrete problem (3.4), we set, for all  $t \in \bar{J}$ ,

$$E_h(t) := \frac{1}{2}\|\partial_t u_{\mathcal{T}}(t)\|_{\Omega}^2 + \frac{1}{2}\hat{a}_h(\hat{u}_h(t), \hat{u}_h(t)). \quad (3.35)$$

This energy verifies

$$E_h(t) = E_h(0) + \int_0^t (f(s), \partial_t u_{\mathcal{T}}(s))_{\Omega} ds, \quad (3.36)$$

as readily follows by testing (3.4) with  $\hat{w}_h := \partial_t \hat{u}_h(t)$  for all  $t \in \bar{J}$  and integrating in time.

The first important step is to derive the fully discrete counterpart of the semi-discrete energy identity (3.36). It is convenient to set, for all  $n \in \{0:N-1\}$ ,

$$\hat{u}_h^{n+\frac{1}{2}} := \frac{1}{2}(\hat{u}_h^n + \hat{u}_h^{n+1}), \quad \delta \hat{u}_h^{n+\frac{1}{2}} := \frac{1}{\Delta t}(\hat{u}_h^{n+1} - \hat{u}_h^n). \quad (3.37)$$

The fully discrete counterpart of  $E_h(t)$  is the discrete energy  $E_h^{n+\frac{1}{2}}$  defined as

$$E_h^{n+\frac{1}{2}} := \frac{1}{2}\|\delta u_{\mathcal{T}}^{n+\frac{1}{2}}\|_{\Omega}^2 - \frac{\Delta t^2}{8}a_h(\delta \hat{u}_h^{n+\frac{1}{2}}, \delta \hat{u}_h^{n+\frac{1}{2}}) + \frac{1}{2}a_h(\hat{u}_h^{n+\frac{1}{2}}, \hat{u}_h^{n+\frac{1}{2}}). \quad (3.38)$$

**Lemma 3.2.3** (Energy balance). *The fully discrete energy verifies, for all  $n \in \{1:N-1\}$ ,*

$$E_h^{n+\frac{1}{2}} = E_h^{\frac{1}{2}} + \frac{1}{2} \sum_{m=1}^n (f^m, u_{\mathcal{T}}^{m+1} - u_{\mathcal{T}}^{m-1})_{\Omega}. \quad (3.39)$$

*Proof.* Let  $n \in \{1:N-1\}$ . For all  $m \in \{1:n\}$ , we test (3.24) with  $\hat{w}_h := \frac{1}{2}(\hat{u}_h^{m+1} - \hat{u}_h^{m-1})$ . We use the identity

$$\begin{aligned} \frac{1}{2}(u_{\mathcal{T}}^{m+1} - 2u_{\mathcal{T}}^m + u_{\mathcal{T}}^{m-1}, u_{\mathcal{T}}^{m+1} - u_{\mathcal{T}}^{m-1})_{\Omega} &= \frac{1}{2}((u_{\mathcal{T}}^{m+1} - u_{\mathcal{T}}^m) - (u_{\mathcal{T}}^m - u_{\mathcal{T}}^{m-1}), (u_{\mathcal{T}}^{m+1} - u_{\mathcal{T}}^m) + (u_{\mathcal{T}}^m - u_{\mathcal{T}}^{m-1}))_{\Omega} \\ &= \frac{1}{2}\|u_{\mathcal{T}}^{m+1} - u_{\mathcal{T}}^m\|_{\Omega}^2 - \frac{1}{2}\|u_{\mathcal{T}}^m - u_{\mathcal{T}}^{m-1}\|_{\Omega}^2 \\ &= \frac{\Delta t^2}{2}\|\delta u_{\mathcal{T}}^{m+\frac{1}{2}}\|_{\Omega}^2 - \frac{\Delta t^2}{2}\|\delta u_{\mathcal{T}}^{m-\frac{1}{2}}\|_{\Omega}^2. \end{aligned}$$

We also use the following computation which exploits the symmetry of  $a_h$ :

$$\begin{aligned} a_h(\hat{u}_h^m, \hat{u}_h^{m+1} - \hat{u}_h^{m-1}) &= \frac{1}{4}a_h(\hat{u}_h^{m+1} + 2\hat{u}_h^m + \hat{u}_h^{m-1}, (\hat{u}_h^{m+1} + \hat{u}_h^m) - (\hat{u}_h^m + \hat{u}_h^{m-1})) \\ &\quad - \frac{1}{4}a_h(\hat{u}_h^{m+1} - 2\hat{u}_h^m + \hat{u}_h^{m-1}, (\hat{u}_h^{m+1} - \hat{u}_h^m) + (\hat{u}_h^m - \hat{u}_h^{m-1})) \\ &= a_h(\hat{u}_h^{m+\frac{1}{2}} + \hat{u}_h^{m-\frac{1}{2}}, \hat{u}_h^{m+\frac{1}{2}} - \hat{u}_h^{m-\frac{1}{2}}) - \frac{\Delta t^2}{4}a_h(\delta\hat{u}_h^{m+\frac{1}{2}} + \delta\hat{u}_h^{m-\frac{1}{2}}, \delta\hat{u}_h^{m+\frac{1}{2}} - \delta\hat{u}_h^{m-\frac{1}{2}}) \\ &= a_h(\hat{u}_h^{m+\frac{1}{2}}, \hat{u}_h^{m+\frac{1}{2}}) - a_h(\hat{u}_h^{m-\frac{1}{2}}, \hat{u}_h^{m-\frac{1}{2}}) \\ &\quad - \frac{\Delta t^2}{4}(a_h(\delta\hat{u}_h^{m+\frac{1}{2}}, \delta\hat{u}_h^{m+\frac{1}{2}}) - a_h(\delta\hat{u}_h^{m-\frac{1}{2}}, \delta\hat{u}_h^{m-\frac{1}{2}})). \end{aligned}$$

This gives  $E_h^{m+\frac{1}{2}} - E_h^{m-\frac{1}{2}} = \frac{1}{2}(f^m, u_{\mathcal{T}}^{m+1} - u_{\mathcal{T}}^{m-1})_{\Omega}$ . Summing this identity for  $m = 1$  to  $m = n$  yields the claim.  $\square$

The second important question is whether  $E_h^{n+\frac{1}{2}}$  defines a strongly convex functional on  $\delta u_{\mathcal{T}}^{n+\frac{1}{2}}$  and  $\hat{u}_h^{n+\frac{1}{2}}$ . This property can be achieved under a CFL restriction on the time step.

**Lemma 3.2.4** (Strong convexity of  $E_h^{n+\frac{1}{2}}$ ). *Under the following CFL restriction on the time step:*

$$\Delta t \leq \eta \mu_{\sharp}^{-1} h, \quad \eta := C_{\text{div}}^{-1} \varpi^{-\frac{1}{2}}, \quad (3.40)$$

$E_h^{n+\frac{1}{2}}$  defines a strongly convex functional on  $\delta u_{\mathcal{T}}^{n+\frac{1}{2}}$  and  $\hat{u}_h^{n+\frac{1}{2}}$ .

*Proof.* Using the fact that  $\hat{u}_h^n$  satisfies (3.26) and the symmetry of  $a_h$ , we infer that, for all  $n \in \{1:N-1\}$ ,

$$\begin{aligned} \Delta t^2 a_h(\delta\hat{u}_h^{n+\frac{1}{2}}, \delta\hat{u}_h^{n+\frac{1}{2}}) &= a_h(\hat{u}_h^{n+1} - \hat{u}_h^n, \hat{u}_h^{n+1} - \hat{u}_h^n) \\ &= a_h(\hat{u}_h^{n+1} - \hat{u}_h^n, (u_{\mathcal{T}}^{n+1} - u_{\mathcal{T}}^n, 0)) \\ &= a_h((u_{\mathcal{T}}^{n+1} - u_{\mathcal{T}}^n, 0), (u_{\mathcal{T}}^{n+1} - u_{\mathcal{T}}^n, 0)) + a_h((0, u_{\mathcal{F}}^{n+1} - u_{\mathcal{F}}^n), (u_{\mathcal{T}}^{n+1} - u_{\mathcal{T}}^n, 0)) \\ &= a_h((u_{\mathcal{T}}^{n+1} - u_{\mathcal{T}}^n, 0), (u_{\mathcal{T}}^{n+1} - u_{\mathcal{T}}^n, 0)) + a_h((u_{\mathcal{T}}^{n+1} - u_{\mathcal{T}}^n, 0), (0, u_{\mathcal{F}}^{n+1} - u_{\mathcal{F}}^n)) \\ &= a_h((u_{\mathcal{T}}^{n+1} - u_{\mathcal{T}}^n, 0), (u_{\mathcal{T}}^{n+1} - u_{\mathcal{T}}^n, 0)) - a_h((0, u_{\mathcal{F}}^{n+1} - u_{\mathcal{F}}^n), (0, u_{\mathcal{T}}^{n+1} - u_{\mathcal{T}}^n)) \\ &= \Delta t^2 a_h((\delta u_{\mathcal{T}}^{n+\frac{1}{2}}, 0), (\delta u_{\mathcal{T}}^{n+\frac{1}{2}}, 0)) - \Delta t^2 a_h((0, \delta u_{\mathcal{F}}^{n+\frac{1}{2}}), (0, \delta u_{\mathcal{F}}^{n+\frac{1}{2}})) \\ &\leq \Delta t^2 a_h((\delta u_{\mathcal{T}}^{n+\frac{1}{2}}, 0), (\delta u_{\mathcal{T}}^{n+\frac{1}{2}}, 0)). \end{aligned}$$

Combining this bound with the coercivity property from Lemma 2.2.3 gives

$$E_h^{n+\frac{1}{2}} \geq \frac{1}{2}\|\delta u_{\mathcal{T}}^{n+\frac{1}{2}}\|_{\Omega}^2 - \frac{\Delta t^2}{8}a_h((\delta u_{\mathcal{T}}^{n+\frac{1}{2}}, 0), (\delta u_{\mathcal{T}}^{n+\frac{1}{2}}, 0)) + \alpha\|\hat{u}_h^{n+\frac{1}{2}}\|_{\text{HNO}}^2.$$

Recalling that  $\|\cdot\|_{\text{HHO}}$  defines a norm on  $\widehat{\mathcal{U}}_{h,0}^{l,k}$  and  $\hat{u}_h^{n+\frac{1}{2}} \in \widehat{\mathcal{U}}_{h,0}^{l,k}$ , the proof is complete if we show that, under the CFL condition (3.40), we have, for all  $w_{\mathcal{T}} \in \mathcal{U}_{\mathcal{T}}^l$ ,

$$\frac{1}{2}\|w_{\mathcal{T}}\|_{\Omega}^2 - \frac{\Delta t^2}{8}a_h((w_{\mathcal{T}}, 0), (w_{\mathcal{T}}, 0)) \geq \frac{1}{4}\|w_{\mathcal{T}}\|_{\Omega}^2. \quad (3.41)$$

To this purpose, we observe using the inverse inequalities from Lemma 2.2.1 and the boundness property from Lemma 2.2.3 that

$$\begin{aligned} a_h((w_{\mathcal{T}}, 0), (w_{\mathcal{T}}, 0)) &\leq \varpi \|(w_{\mathcal{T}}, 0)\|_{\text{HHO}}^2 \\ &\leq \varpi \mu_{\sharp}^2 \left\{ \sum_{T \in \mathcal{T}_h} \|\nabla w_T\|_T^2 + h_T^{-1} \|w_T\|_{\partial T}^2 \right\} \\ &\leq 2C_{\text{div}}^2 \varpi \mu_{\sharp}^2 h^{-2} \|w_{\mathcal{T}}\|_{\Omega}^2. \end{aligned}$$

Using the CFL condition (3.40), we infer that

$$\frac{\Delta t^2}{8}a_h((w_{\mathcal{T}}, 0), (w_{\mathcal{T}}, 0)) \leq \frac{\Delta t^2}{8}2C_{\text{div}}^2 \varpi \mu_{\sharp}^2 h^{-2} \|w_{\mathcal{T}}\|_{\Omega}^2 \leq \frac{1}{4}\|w_{\mathcal{T}}\|_{\Omega}^2.$$

Re-arranging the terms establishes (3.41) and completes the proof.  $\square$

**Remark 3.2.5** (CFL condition). The CFL condition can be slightly sharpened, i.e. multiplying  $\eta$  by some number larger than 1, using the ideas of [52].

## 3.2.2 Convergence analysis in energy norm

As seen in Section 3.1, the HHO space semi-discretization yields a convergence of order  $h^{k+1}$  in the energy norm, while the leapfrog scheme is of second order in time. The expected error upper bound for the fully discrete scheme is then of the form  $C(h^{k+1} + \Delta t^2)$ . The goal of this section is to prove this result under sufficient smoothness of the exact solution. We start the section by stating our main result, Theorem 3.2.6, and devote the rest of the section to its proof.

### 3.2.2.1 Main result

Recall that  $u$  denotes the solution to the continuous wave equation (3.2) with the initial conditions (3.3) and  $(\hat{u}_h^n)_{n \in \{0:N\}}$  the solution to the fully discrete wave equation (3.24) with the initial conditions (3.25). We define the discrete error

$$\hat{e}_h^n := \hat{u}_h^n - \hat{I}_h^{k,l}(u(t^n)), \quad \forall n \in \{0:N\}, \quad (3.42)$$

which represents the difference between the discrete hybrid solution at time step  $n$  and the projection of the continuous solution at the discrete time  $t^n$  onto the hybrid space. Recalling the seminorm  $|\cdot|_{*,h}$  defined in (3.8), we set, for all  $t \in \bar{J}$ ,

$$|u|_{C^0(0,t;*,h)} := \sup_{s \in [0,t]} |u(s)|_{*,h}, \quad |\partial_t u|_{L^1(0,t;*,h)} := \int_0^t |\partial_t u(s)|_{*,h} ds. \quad (3.43)$$

**Theorem 3.2.6** (Energy error estimate). *Assume the CFL condition (3.40), that  $f \in C^1(\bar{J}; L^2(\Omega))$ , and that  $u \in C^4(\bar{J}; L^2(\Omega)) \cap W^{5,1}(J; L^2(\Omega))$  and  $u \in C^0(\bar{J}; H^{1+\nu}(\Omega)) \cap W^{1,1}(J; H^{1+\nu}(\Omega))$ ,  $\nu \in (\frac{1}{2}, 1]$ . The following holds:*

$$\begin{aligned} \max_{n \in \{0:N-1\}} \|\delta e_{\mathcal{T}}^{n+\frac{1}{2}}\|_{\Omega} + \max_{n \in \{0:N-1\}} \|\hat{e}_h^{n+\frac{1}{2}}\|_{\text{HHO}} &\leq C_1 \left\{ \|u\|_{C^0([0;t^{N-1}];*,h)} + \|\partial_t u\|_{L^1([0;t^{N-1}];*,h)} \right\} \\ &+ C_2 \Delta t^2 \left\{ \|\partial_t^3 u\|_{C^0([0;t^1];L^2(\Omega))} + \Theta \left( \|\partial_t^4 u\|_{C^0([0;t^N];L^2(\Omega))} + \|\partial_t^5 u\|_{L^1([0;t^N];L^2(\Omega))} \right) \right\}, \end{aligned} \quad (3.44)$$

with generic constants  $C_1$  and  $C_2$  and the time scale  $\Theta := \mu_{\sharp}^{-1} \ell_{\Omega}$ .

**Remark 3.2.7** (Regularity assumption). The regularity assumption in space hinging on the shift  $\nu \in (\frac{1}{2}, 1]$  is made to simplify the bound on the consistency error in space. The quasi-minimal setting with  $\nu \in (0, 1]$  can be handled by using the tools introduced in [99], see also [98, Sec. 41.5].

**Remark 3.2.8** (Convergence order). Under the above regularity assumption on  $u$ , we can bound  $|u(t)|_{*,h}$  and  $|\partial_t u(t)|_{*,h}$ , for all  $t \in \bar{J}$ , using the approximation results from Lemma 2.2.5. This gives

$$\begin{aligned} \max_{n \in \{0:N-1\}} \|\delta e_{\mathcal{T}}^{n+\frac{1}{2}}\|_{\Omega} + \|\hat{e}_h^{n+\frac{1}{2}}\|_{\text{HMO}} &\leq C_1 h^{k+1} \left\{ |u|_{C^0([0;t^{N-1}]; H^{k+2}(\Omega))} + |\partial_t u|_{L^1([t^0, t^{N-1}]; H^{k+2}(\Omega))} \right\} \\ &+ C_2 \Delta t^2 \left\{ \|\partial_t^3 u\|_{C^0([0;t^1]; L^2(\Omega))} + \Theta(\|\partial_t^4 u\|_{C^0([0;t^N]; L^2(\Omega))} + \|\partial_t^5 u\|_{L^1([0;t^N]; L^2(\Omega)}) \right\}, \end{aligned} \quad (3.45)$$

which is the optimal convergence order  $\mathcal{O}(h^{k+1} + \Delta t^2)$  in the energy norm.

### 3.2.2.2 Consistency errors and error equation

We define the space consistency error,  $\psi_h^n$ , at the discrete time  $t^n$ , for all  $n \in \{0:N\}$ , as the linear form in  $(\hat{\mathcal{U}}_{h,0}^{l,k})'$  such that, for all  $\hat{w}_h \in \hat{\mathcal{U}}_{h,0}^{l,k}$ ,

$$\psi_h^n(\hat{w}_h) := \psi_h(t^n, \hat{w}_h), \quad (3.46)$$

where  $\psi_h$  is defined in the proof of Theorem 3.1.1. We also define the time consistency error,  $\kappa^n$ , to be the error between the centered difference scheme and the exact second-order derivative in time:

$$\kappa^n := \frac{u(t^{n+1}) - 2u(t^n) + u(t^{n-1}))}{\Delta t^2} - \partial_t^2 u(t^n), \quad \forall n \in \{1:N-1\}, \quad (3.47)$$

and we set  $\kappa^0 := 0$ .

**Lemma 3.2.9** (Discrete error equation). *Let  $(\hat{e}_h^n)_{n \in \{0:N\}}$  be the collection of discrete errors defined in (3.42). The following holds for all  $n \in \{1:N-1\}$ :*

$$\frac{1}{\Delta t^2} (e_{\mathcal{T}}^{n+1} - 2e_{\mathcal{T}}^n + e_{\mathcal{T}}^{n-1}, w_{\mathcal{T}})_{\Omega} + a_h(\hat{e}_h^n, \hat{w}_h) = \psi_h^n(\hat{w}_h) - (\kappa^n, w_{\mathcal{T}})_{\Omega}, \quad \forall \hat{w}_h \in \hat{\mathcal{U}}_{h,0}^{l,k}. \quad (3.48)$$

*Proof.* Let us evaluate the left-hand-side of (3.48) for all  $n \in \{1:N-1\}$ :

$$\begin{aligned} &\frac{1}{\Delta t^2} (e_{\mathcal{T}}^{n+1} - 2e_{\mathcal{T}}^n + e_{\mathcal{T}}^{n-1}, w_{\mathcal{T}})_{\Omega} + a_h(\hat{e}_h^n, \hat{w}_h) \\ &= (f^n, w_{\mathcal{T}})_{\Omega} - \left( \Pi_{\mathcal{T}}^l \left[ \frac{u(t^{n+1}) - 2u(t^n) + u(t^{n-1}))}{\Delta t^2} \right], w_{\mathcal{T}} \right)_{\Omega} - a_h(\hat{I}_h^{k,l}(u(t^n)), \hat{w}_h) \\ &= (f^n, w_{\mathcal{T}})_{\Omega} - (\Pi_{\mathcal{T}}^l(\partial_t^2 u(t^n) + \kappa^n), w_{\mathcal{T}})_{\Omega} - a_h(\hat{I}_h^{k,l}(u(t^n)), \hat{w}_h) \\ &= (\partial_t^2 u(t^n), w_{\mathcal{T}})_{\Omega} + (B(u(t^n)), w_{\mathcal{T}})_{\Omega} - (\partial_t^2 u(t^n) + \kappa^n, w_{\mathcal{T}})_{\Omega} - a_h(\hat{I}_h^{k,l}(u(t^n)), \hat{w}_h) \\ &= \psi_h^n(\hat{w}_h) - (\kappa^n, w_{\mathcal{T}})_{\Omega}, \end{aligned}$$

where we used the discrete scheme (3.24), the fact that the exact solution satisfies  $\partial_t^2 u(t^n) + B(u(t^n)) = f^n$ , the definition of the  $L^2$ -projection onto  $\mathcal{U}_{\mathcal{T}}^l$  and the definition (3.46) of  $\psi_h^n$ .  $\square$

### 3.2.2.3 Energy-error identity

It is convenient to define, for all  $n \in \{0:N-1\}$ , the error and velocity error at the half time-steps as follows:

$$\hat{e}_h^{n+\frac{1}{2}} := \frac{1}{2}(\hat{e}_h^{n+1} + \hat{e}_h^n), \quad \delta \hat{e}_h^{n+\frac{1}{2}} := \frac{1}{\Delta t}(\hat{e}_h^{n+1} - \hat{e}_h^n). \quad (3.49)$$

Concerning the consistency errors, we set, for all  $n \in \{0:N-1\}$ ,

$$\delta\kappa^{n+\frac{1}{2}} := \frac{1}{\Delta t}(\kappa^{n+1} - \kappa^n), \quad \delta\psi_h^{n+\frac{1}{2}}(\hat{w}_h) := \frac{1}{\Delta t}(\psi_h^{n+1}(\hat{w}_h) - \psi_h^n(\hat{w}_h)), \quad \forall \hat{w}_h \in \widehat{\mathcal{U}}_{h,0}^{l,k}. \quad (3.50)$$

Let us consider the discrete energy (3.38) evaluated using the error. We set

$$\mathcal{E}_h^{n+\frac{1}{2}} := \frac{1}{2}\|\delta e_{\mathcal{T}}^{n+\frac{1}{2}}\|_{\Omega}^2 - \frac{\Delta t^2}{8}a_h(\delta\hat{e}_h^{n+\frac{1}{2}}, \delta\hat{e}_h^{n+\frac{1}{2}}) + \frac{1}{2}a_h(\hat{e}_h^{n+\frac{1}{2}}, \hat{e}_h^{n+\frac{1}{2}}), \quad \forall n \in \{0:N-1\}. \quad (3.51)$$

Unfortunately, even under a CFL condition,  $\mathcal{E}_h^{n+\frac{1}{2}}$  does not define a strongly convex functional on  $\delta e_{\mathcal{T}}^{n+\frac{1}{2}}$  and  $\hat{e}_h^{n+\frac{1}{2}}$ . The reason is that we no longer have  $a_h(\hat{e}_h^n, (0, w_{\mathcal{F}})) = 0$  for all  $w_{\mathcal{F}} \in \mathcal{U}_{\mathcal{F},0}^k$ , but only

$$a_h(\hat{e}_h^n, (0, w_{\mathcal{F}})) = \psi_h^n((0, w_{\mathcal{F}})), \quad (3.52)$$

as a consequence of (3.48). This leads to the following definition of discrete energy error: For all  $n \in \{0:N-1\}$ ,

$$\check{\mathcal{E}}_h^{n+\frac{1}{2}} := \mathcal{E}_h^{n+\frac{1}{2}} + \frac{\Delta t^2}{4}\delta\psi_h^{n+\frac{1}{2}}((0, \delta e_{\mathcal{F}}^{n+\frac{1}{2}})). \quad (3.53)$$

**Lemma 3.2.10** (Discrete energy error). *For all  $n \in \{0:N-1\}$ , we have*

$$\begin{aligned} \check{\mathcal{E}}_h^{n+\frac{1}{2}} &= \frac{1}{2}\|\delta e_{\mathcal{T}}^{n+\frac{1}{2}}\|_{\Omega}^2 + \frac{\Delta t^2}{8}a_h((0, \delta e_{\mathcal{F}}^{n+\frac{1}{2}}), (0, \delta e_{\mathcal{F}}^{n+\frac{1}{2}})) \\ &\quad - \frac{\Delta t^2}{8}a_h((\delta e_{\mathcal{T}}^{n+\frac{1}{2}}, 0), (\delta e_{\mathcal{T}}^{n+\frac{1}{2}}, 0)) + \frac{1}{2}a_h(\hat{e}_h^{n+\frac{1}{2}}, \hat{e}_h^{n+\frac{1}{2}}). \end{aligned} \quad (3.54)$$

Moreover, under the CFL condition (3.40),  $\check{\mathcal{E}}_h^{n+\frac{1}{2}}$  defines a strongly convex functional on  $\delta e_{\mathcal{T}}^{n+\frac{1}{2}}$  and  $\hat{e}_h^{n+\frac{1}{2}}$ .

*Proof.* Proceeding as in Lemma 3.2.4, we obtain

$$\begin{aligned} \Delta t^2 a_h(\delta\hat{e}_h^{n+\frac{1}{2}}, \delta\hat{e}_h^{n+\frac{1}{2}}) &= a_h(\hat{e}_h^{n+1} - \hat{e}_h^n, \hat{e}_h^{n+1} - \hat{e}_h^n) \\ &= a_h(\hat{e}_h^{n+1} - \hat{e}_h^n, (e_{\mathcal{T}}^{n+1} - e_{\mathcal{T}}^n, 0)) + a_h(\hat{e}_h^{n+1} - \hat{e}_h^n, (0, e_{\mathcal{F}}^{n+1} - e_{\mathcal{F}}^n)) \\ &= a_h(\hat{e}_h^{n+1} - \hat{e}_h^n, (e_{\mathcal{T}}^{n+1} - e_{\mathcal{T}}^n, 0)) + \Delta t \left\{ \psi_h^{n+1}((0, \delta e_{\mathcal{F}}^{n+\frac{1}{2}})) - \psi_h^n((0, \delta e_{\mathcal{F}}^{n+\frac{1}{2}})) \right\} \\ &= a_h((e_{\mathcal{T}}^{n+1} - e_{\mathcal{T}}^n, 0), (e_{\mathcal{T}}^{n+1} - e_{\mathcal{T}}^n, 0)) + a_h((0, e_{\mathcal{F}}^{n+1} - e_{\mathcal{F}}^n), (e_{\mathcal{T}}^{n+1} - e_{\mathcal{T}}^n, 0)) \\ &\quad + \Delta t^2 \delta\psi_h^{n+\frac{1}{2}}((0, \delta e_{\mathcal{F}}^{n+\frac{1}{2}})), \end{aligned}$$

where we used (3.52) and the notation (3.49) for  $\delta\psi_h^{n+\frac{1}{2}}$ . Owing to the symmetry of  $a_h$  and again (3.52), we infer that

$$\begin{aligned} a_h((0, e_{\mathcal{F}}^{n+1} - e_{\mathcal{F}}^n), (e_{\mathcal{T}}^{n+1} - e_{\mathcal{T}}^n, 0)) &= a_h((e_{\mathcal{T}}^{n+1} - e_{\mathcal{T}}^n, 0), (0, e_{\mathcal{F}}^{n+1} - e_{\mathcal{F}}^n)) \\ &= -a_h((0, e_{\mathcal{F}}^{n+1} - e_{\mathcal{F}}^n), (0, e_{\mathcal{F}}^{n+1} - e_{\mathcal{F}}^n)) + \Delta t^2 \delta\psi_h^{n+\frac{1}{2}}((0, \delta e_{\mathcal{F}}^{n+\frac{1}{2}})). \end{aligned}$$

This gives

$$\begin{aligned} \Delta t^2 a_h(\delta\hat{e}_h^{n+\frac{1}{2}}, \delta\hat{e}_h^{n+\frac{1}{2}}) &= a_h((e_{\mathcal{T}}^{n+1} - e_{\mathcal{T}}^n, 0), (e_{\mathcal{T}}^{n+1} - e_{\mathcal{T}}^n, 0)) - a_h((0, e_{\mathcal{F}}^{n+1} - e_{\mathcal{F}}^n), (0, e_{\mathcal{F}}^{n+1} - e_{\mathcal{F}}^n)) + 2\Delta t^2 \delta\psi_h^{n+\frac{1}{2}}((0, \delta e_{\mathcal{F}}^{n+\frac{1}{2}})) \\ &= \Delta t^2 a_h((\delta e_{\mathcal{T}}^{n+\frac{1}{2}}, 0), (\delta e_{\mathcal{T}}^{n+\frac{1}{2}}, 0)) - \Delta t^2 a_h((0, \delta e_{\mathcal{F}}^{n+\frac{1}{2}}), (0, \delta e_{\mathcal{F}}^{n+\frac{1}{2}})) + 2\Delta t^2 \delta\psi_h^{n+\frac{1}{2}}((0, \delta e_{\mathcal{F}}^{n+\frac{1}{2}})). \end{aligned}$$

This proves the identity (3.54). Finally, the strong convexity of  $\check{\mathcal{E}}_h^{n+\frac{1}{2}}$  under the CFL condition (3.40) is established as in the proof of Lemma 3.2.4.  $\square$



The next step is to write an energy identity mimicking the space semi-discrete case.

**Lemma 3.2.11** (Energy identity). *For all  $n \in \{1:N-1\}$ , the discrete energy error  $\check{\mathcal{E}}_h^{n+\frac{1}{2}}$  verifies*

$$\check{\mathcal{E}}_h^{n+\frac{1}{2}} = \mathcal{Z}_\psi^n + \mathcal{Z}_\kappa^n + \mathcal{Z}_{\text{IC}}^0, \quad (3.55)$$

with the space consistency error

$$\mathcal{Z}_\psi^n := \psi_h^n(\hat{e}_h^{n+\frac{1}{2}}) - \frac{\Delta t^2}{4} \delta \psi_h^{n+\frac{1}{2}}((0, \delta e_{\mathcal{F}}^{n+\frac{1}{2}})) + \Delta t \sum_{m=1}^{n-1} \delta \psi_h^{m+\frac{1}{2}}(\hat{e}_h^{m+\frac{1}{2}}), \quad (3.56a)$$

the time consistency error

$$\mathcal{Z}_\kappa^n := -(\kappa^n, e_{\mathcal{T}}^{n+\frac{1}{2}})_\Omega + \Delta t \sum_{m=1}^{n-1} (\delta \kappa^{m+\frac{1}{2}}, e_{\mathcal{T}}^{m+\frac{1}{2}})_\Omega, \quad (3.56b)$$

and the initial error

$$\mathcal{Z}_{\text{IC}}^0 := \mathcal{E}_h^{\frac{1}{2}} - \psi_h^1(\hat{e}_h^{\frac{1}{2}}) + (\kappa^1, e_{\mathcal{T}}^{\frac{1}{2}})_\Omega. \quad (3.56c)$$

*Proof.* Let us start by remarking that, similarly to the computations in the proof of Lemma 3.2.3, we have, for all  $m \in \{1:N-1\}$ ,

$$\mathcal{E}_h^{m+\frac{1}{2}} - \mathcal{E}_h^{m-\frac{1}{2}} = \frac{1}{\Delta t^2} (e_{\mathcal{T}}^{m+1} - 2e_{\mathcal{T}}^m + e_{\mathcal{T}}^{m-1}, e_{\mathcal{T}}^{m+1} - e_{\mathcal{T}}^{m-1})_\Omega + a_h(\hat{e}_h^m, \hat{e}_h^{m+1} - \hat{e}_h^{m-1}).$$

Using Lemma 3.2.9 with the test function  $\hat{w}_h := \frac{1}{2}(\hat{e}_h^{m+1} - \hat{e}_h^{m-1}) = \hat{e}_h^{m+\frac{1}{2}} - \hat{e}_h^{m-\frac{1}{2}}$  yields

$$\mathcal{E}_h^{m+\frac{1}{2}} - \mathcal{E}_h^{m-\frac{1}{2}} = \psi_h^m(\hat{e}_h^{m+\frac{1}{2}} - \hat{e}_h^{m-\frac{1}{2}}) - (\kappa^m, e_{\mathcal{T}}^{m+\frac{1}{2}} - e_{\mathcal{T}}^{m-\frac{1}{2}})_\Omega. \quad (3.57)$$

Summing (3.57) from  $m = 1$  to  $m = n$  yields

$$\mathcal{E}_h^{n+\frac{1}{2}} = \mathcal{E}_h^{\frac{1}{2}} + \sum_{m=1}^n \left\{ \psi_h^m(\hat{e}_h^{m+\frac{1}{2}} - \hat{e}_h^{m-\frac{1}{2}}) - (\kappa^m, e_{\mathcal{T}}^{m+\frac{1}{2}} - e_{\mathcal{T}}^{m-\frac{1}{2}})_\Omega \right\}. \quad (3.58)$$

The sums on the right-hand side of (3.58) can be reordered using the notation (3.49) to obtain

$$\begin{aligned} \sum_{m=1}^n \psi_h^m(\hat{e}_h^{m+\frac{1}{2}} - \hat{e}_h^{m-\frac{1}{2}}) &= \psi_h^n(\hat{e}_h^{n+\frac{1}{2}}) - \psi_h^1(\hat{e}_h^{\frac{1}{2}}) - \Delta t \sum_{m=1}^{n-1} \delta \psi_h^{m+\frac{1}{2}}(\hat{e}_h^{m+\frac{1}{2}}), \\ \sum_{m=1}^n (\kappa^m, e_{\mathcal{T}}^{m+\frac{1}{2}} - e_{\mathcal{T}}^{m-\frac{1}{2}})_\Omega &= (\kappa^n, e_{\mathcal{T}}^{n+\frac{1}{2}})_\Omega - (\kappa^1, e_{\mathcal{T}}^{\frac{1}{2}})_\Omega - \Delta t \sum_{m=1}^{n-1} (\delta \kappa^{m+\frac{1}{2}}, e_{\mathcal{T}}^{m+\frac{1}{2}})_\Omega. \end{aligned}$$

This gives

$$\begin{aligned} \mathcal{E}_h^{n+\frac{1}{2}} &= \mathcal{E}_h^{\frac{1}{2}} - \psi_h^1(\hat{e}_h^{\frac{1}{2}}) + (\kappa^1, e_{\mathcal{T}}^{\frac{1}{2}})_\Omega + \left\{ \psi_h^n(\hat{e}_h^{n+\frac{1}{2}}) - \Delta t \sum_{m=1}^{n-1} \delta \psi_h^{m+\frac{1}{2}}(\hat{e}_h^{m+\frac{1}{2}}) \right\} \\ &\quad - (\kappa^n, e_{\mathcal{T}}^{n+\frac{1}{2}})_\Omega + \Delta t \sum_{m=1}^{n-1} (\delta \kappa^{m+\frac{1}{2}}, e_{\mathcal{T}}^{m+\frac{1}{2}})_\Omega \\ &= \psi_h^n(\hat{e}_h^{n+\frac{1}{2}}) - \Delta t \sum_{m=1}^{n-1} \delta \psi_h^{m+\frac{1}{2}}(\hat{e}_h^{m+\frac{1}{2}}) + \mathcal{Z}_\kappa^n + \mathcal{Z}_{\text{IC}}^0, \end{aligned}$$

owing to the definitions (3.56b) and (3.56c). It remains to go from  $\mathcal{E}_h^{n+\frac{1}{2}}$  to  $\check{\mathcal{E}}_h^{n+\frac{1}{2}}$ . Using (3.53), we have

$$\check{\mathcal{E}}_h^{n+\frac{1}{2}} = \psi_h^n(\hat{e}_h^{n+\frac{1}{2}}) - \Delta t \sum_{m=1}^{n-1} \delta \psi_h^{m+\frac{1}{2}}(\hat{e}_h^{m+\frac{1}{2}}) + \mathcal{Z}_\kappa^n + \mathcal{Z}_{\text{IC}}^0 + \frac{\Delta t^2}{4} \delta \psi_h^{n+\frac{1}{2}}((0, \delta e_{\mathcal{F}}^{n+\frac{1}{2}})).$$

Recalling the definition of  $\mathcal{Z}_\psi^n$  proves the claim.  $\square$

### 3.2.2.4 Bound on consistency and initial errors

We now bound the three terms on the right-hand side of (3.55). Each estimate is stated as a separate lemma.

**Lemma 3.2.12** (Estimate on the space consistency error). *Let  $\mathcal{Z}_\psi^n$  be defined in (3.56a). For all  $n \in \{1:N-1\}$ , the following holds:*

$$\begin{aligned} |\mathcal{Z}_\psi^n| &\leq c_* (|\partial_t u|_{L^1([t^1, t^n]; *, h)} + |u(t^n)|_{*, h}) \max_{m \in \{1:n\}} \|\hat{e}_h^{m+\frac{1}{2}}\|_{\text{HHO}} \\ &\quad + \frac{1}{4} c_* \alpha^{-\frac{1}{2}} |\partial_t u|_{L^1([t^n, t^{n+1}]; *, h)} \left[ \Delta t^2 a_h((0, \delta e_{\mathcal{F}}^{m+\frac{1}{2}}), (0, \delta e_{\mathcal{F}}^{m+\frac{1}{2}})) \right]^{\frac{1}{2}}, \end{aligned}$$

where  $c_*$  results from (3.20) and  $\alpha$  is the coercivity constant of the discrete bilinear form  $a_h$ .

*Proof.* Let  $n \in \{1:N-1\}$  and let  $m \in \{1:n\}$ . We observe that

$$\Delta t \delta \psi_h^{m+\frac{1}{2}}(\hat{e}_h^{m+\frac{1}{2}}) = \psi_h(t^{m+1}, \hat{e}_h^{m+\frac{1}{2}}) - \psi_h(t^m, \hat{e}_h^{m+\frac{1}{2}}) = \int_{t^m}^{t^{m+1}} \dot{\psi}_h(s; \hat{e}_h^{m+\frac{1}{2}}) ds.$$

Since  $\partial_t u(t) \in Y \cap H^{1+\nu}(\Omega)$  for all  $t \in \bar{J}$  by assumption (indeed,  $B(\partial_t u(t)) = \partial_t f(t) - \partial_t^3 u(t)$ ,  $f \in C^1(\bar{J}; L^2(\Omega))$ , and  $u \in C^4(\bar{J}; L^2(\Omega))$ ), we can invoke (3.20) which gives

$$\|\dot{\psi}_h(t; \cdot)\|_{(\text{HHO})'} \leq c_* |\partial_t u(t)|_{*, h}.$$

Thus, we obtain

$$\left| \Delta t \sum_{m=1}^{n-1} \delta \psi_h^{m+\frac{1}{2}}(\hat{e}_h^{m+\frac{1}{2}}) \right| \leq c_* \sum_{m=1}^{n-1} |\partial_t u|_{L^1([t^m, t^{m+1}]; *, h)} \|\hat{e}_h^{m+\frac{1}{2}}\|_{\text{HHO}}.$$

The same reasoning is used for the second term on the right-hand-side of (3.56a) which is bounded as follows:

$$\left| \frac{\Delta t^2}{4} \delta \psi_h^{n+\frac{1}{2}}((0, \delta e_{\mathcal{F}}^{n+\frac{1}{2}})) \right| \leq \frac{1}{4} c_* \alpha^{-\frac{1}{2}} |\partial_t u|_{L^1([t^n, t^{n+1}]; *, h)} \left[ \Delta t^2 a_h((0, \delta e_{\mathcal{F}}^{n+\frac{1}{2}}), (0, \delta e_{\mathcal{F}}^{n+\frac{1}{2}})) \right]^{\frac{1}{2}},$$

where we used the coercivity of the discrete bilinear form  $a_h$ . Finally, the same argument for the first term on the right-hand side of (3.56a) yields

$$\left| \psi_h^n(\hat{e}_h^{n+\frac{1}{2}}) \right| \leq c_* |u(t^n)|_{*, h} \|\hat{e}_h^{n+\frac{1}{2}}\|_{\text{HHO}}. \quad (3.59)$$

This yields the claim.  $\square$

**Lemma 3.2.13** (Estimate on the time consistency error). *Let  $\mathcal{Z}_\kappa^n$  be defined in (3.56b). For all  $n \in \{1:N-1\}$ , the following holds:*

$$|\mathcal{Z}_\kappa^n| \leq C \Theta \Delta t^2 \left( \|\partial_t^5 u\|_{L^1([0, t^{n+1}]; L^2(\Omega))} + \|\partial_t^4 u\|_{C^0([t^{n-1}, t^{n+1}]; L^2(\Omega))} \right) \max_{m \in \{1:n\}} \|\hat{e}_h^{m+\frac{1}{2}}\|_{\text{HHO}}.$$

*Proof.* Let  $n \in \{1:N-1\}$ . A straightforward calculation using fourth-order Taylor expansions with integral remainder shows that

$$\kappa^n = \frac{1}{6\Delta t^2} \left\{ \int_{t^n}^{t^{n+1}} (t^{n+1} - t)^3 \partial_t^4 u(t) dt - \int_{t^{n-1}}^{t^n} (t^{n-1} - t)^3 \partial_t^4 u(t) dt \right\}.$$

Rewriting  $\delta\kappa^{m+\frac{1}{2}}$ , for all  $m \in \{1:n-1\}$ , using this identity gives

$$\begin{aligned}
\delta\kappa^{m+\frac{1}{2}} &= \frac{1}{\Delta t} (\kappa^{m+1} - \kappa^m) \\
&= \frac{1}{6\Delta t^3} \left( \int_{t^{m+1}}^{t^{m+2}} (t^{m+2} - t)^3 \partial_t^4 u(t) dt - \int_{t^m}^{t^{m+1}} (t^m - t)^3 \partial_t^4 u(t) dt \right. \\
&\quad \left. - \int_{t^m}^{t^{m+1}} (t^{m+1} - t)^3 \partial_t^4 u(t) dt + \int_{t^{m-1}}^{t^m} (t^{m-1} - t)^3 \partial_t^4 u(t) dt \right) \\
&= \frac{1}{6\Delta t^3} \left( \int_{t^m}^{t^{m+1}} (t^{m+1} - t)^3 \int_t^{t+\Delta t} \partial_t^5 u(s) ds dt - \int_{t^{m-1}}^{t^m} (t^{m-1} - t)^3 \int_t^{t+\Delta t} \partial_t^5 u(s) ds dt \right).
\end{aligned}$$

Let us now use the following estimate on the first integral:

$$\left\| \int_{t^m}^{t^{m+1}} (t^{m+1} - t)^3 \int_t^{t+\Delta t} \partial_t^5 u(s) ds dt \right\|_{\Omega} \leq \Delta t^4 \|\partial_t^5 u\|_{L^1([t^m; t^{m+2}]; L^2(\Omega))}, \quad (3.60)$$

and a similar estimate on the second integral. Then, we invoke the discrete Poincaré inequality (2.21) to obtain

$$\begin{aligned}
\left| \Delta t \sum_{m=1}^{n-1} (\delta\kappa^{m+\frac{1}{2}}, e_{\mathcal{T}}^{m+\frac{1}{2}})_{\Omega} \right| &\leq C\mu_{\sharp}^{-1} \ell_{\Omega} \Delta t^2 \sum_{m=1}^{n-1} \|\partial_t^5 u\|_{L^1([t^{m-1}; t^{m+2}]; L^2(\Omega))} \|\hat{e}_h^{m+\frac{1}{2}}\|_{\text{HHO}} \\
&= C\Theta \Delta t^2 \sum_{m=1}^{n-1} \|\partial_t^5 u\|_{L^1([t^{m-1}; t^{m+2}]; L^2(\Omega))} \|\hat{e}_h^{m+\frac{1}{2}}\|_{\text{HHO}}.
\end{aligned} \quad (3.61)$$

The first term on the right-hand-side of (3.56b) can be bounded in the same way, yielding

$$\left| (\kappa^n, e_{\mathcal{T}}^{n+\frac{1}{2}})_{\Omega} \right| \leq C\Theta \Delta t^2 \|\partial_t^4 u\|_{C^0([t^{n-1}; t^{n+1}]; L^2(\Omega))} \|\hat{e}_h^{n+\frac{1}{2}}\|_{\text{HHO}}. \quad (3.62)$$

Taking the maximum over  $m \in \{1:n-1\}$  on the right-hand-side of (3.61) and summing to (3.62) yields the expected result.  $\square$

**Lemma 3.2.14** (Estimate on the initial error). *Let  $\mathcal{Z}_{\text{IC}}^0$  be defined in (3.56c). Assume the CFL condition (3.40). There is a constant  $C_{\text{IC}}$  such that*

$$|\mathcal{Z}_{\text{IC}}^0| \leq C_{\text{IC}} \left\{ |u_0|_{*,h}^2 + |u(t^1)|_{*,h}^2 + \Delta t^4 (\|\partial_t^3 u\|_{C^0([0; t^1]; L^2(\Omega))}^2 + \Theta^2 \|\partial_t^4 u\|_{C^0([t^0; t^2]; L^2(\Omega))}^2) \right\}.$$

*Proof.* Recall that

$$\begin{aligned}
\mathcal{Z}_{\text{IC}}^0 &= \mathcal{E}_h^{\frac{1}{2}} - \psi_h^1(\hat{e}_h^{\frac{1}{2}}) + (\kappa^1, e_{\mathcal{T}}^{\frac{1}{2}})_{\Omega} \\
&= \frac{1}{2\Delta t^2} \|e_{\mathcal{T}}^1 - e_{\mathcal{T}}^0\|_{\Omega}^2 - \frac{1}{8} a_h(\hat{e}_h^1 - \hat{e}_h^0, \hat{e}_h^1 - \hat{e}_h^0) + \frac{1}{2} a_h(\hat{e}_h^{\frac{1}{2}}, \hat{e}_h^{\frac{1}{2}}) - \psi_h^1(\hat{e}_h^{\frac{1}{2}}) + (\kappa^1, e_{\mathcal{T}}^{\frac{1}{2}})_{\Omega}.
\end{aligned} \quad (3.63)$$

Since the second term is negative, we infer that  $\mathcal{Z}_{\text{IC}}^0 \leq A_1 + A_2 + A_3$  with

$$A_1 := \frac{1}{2\Delta t^2} \|e_{\mathcal{T}}^1 - e_{\mathcal{T}}^0\|_{\Omega}^2, \quad A_2 := \frac{1}{2} a_h(\hat{e}_h^{\frac{1}{2}}, \hat{e}_h^{\frac{1}{2}}), \quad A_3 := (\kappa^1, e_{\mathcal{T}}^{\frac{1}{2}})_{\Omega} - \psi_h^1(\hat{e}_h^{\frac{1}{2}}).$$

*Step 1: Bound on  $A_1$ .* We use a first-order Taylor expansion of  $u(t^1)$  with integral remainder to obtain

$$u(t^1) = u_0 + \Delta t v_0 + \frac{\Delta t^2}{2} \partial_t^2 u(0) + \frac{1}{2} \int_0^{t^1} (t^1 - s)^2 \partial_t^3 u(s) ds.$$

The definition of  $u_{\mathcal{T}}^1$  via the initial condition (6.43c) then gives, for all  $w_{\mathcal{T}} \in \mathcal{U}_{\mathcal{T}}^1$ ,

$$\begin{aligned} \frac{1}{\Delta t}(e_{\mathcal{T}}^1, w_{\mathcal{T}})_{\Omega} &= \frac{1}{\Delta t}(u_{\mathcal{T}}^1 - u(t^1), w_{\mathcal{T}})_{\Omega} \\ &= \frac{\Delta t}{2}(f^0 - \partial_t^2 u(0), w_{\mathcal{T}})_{\Omega} - \frac{1}{2\Delta t} \int_0^{t^1} (t^1 - s)^2 (\partial_t^3 u(s), w_{\mathcal{T}})_{\Omega} ds - \frac{\Delta t}{2} a_h(\hat{u}_h^0, (w_{\mathcal{T}}, 0)) \\ &= \frac{\Delta t}{2} \left\{ (B(u^0), w_{\mathcal{T}})_{\Omega} - a_h(\hat{u}_h^0, (w_{\mathcal{T}}, 0)) \right\} - \frac{1}{2\Delta t} \int_0^{t^1} (t^1 - s)^2 (\partial_t^3 u(s), w_{\mathcal{T}})_{\Omega} ds \\ &= \frac{\Delta t}{2} \psi_h^0((w_{\mathcal{T}}, 0)) - \frac{1}{2\Delta t} \int_0^{t^1} (t^1 - s)^2 (\partial_t^3 u(s), w_{\mathcal{T}})_{\Omega} ds. \end{aligned}$$

We use both discrete inverse inequalities from Lemma 2.2.1 to obtain

$$\begin{aligned} \|(w_{\mathcal{T}}, 0)\|_{\text{HHO}} &= \left( \sum_{T \in \mathcal{T}_h} \mu_T^2 \{ \|\nabla w_T\|_T^2 + h_T^{-1} \|w_T\|_{\partial T}^2 \} \right)^{\frac{1}{2}} \\ &\leq C \mu_{\sharp} h^{-1} \|w_{\mathcal{T}}\|_{\Omega}. \end{aligned}$$

Using this inequality, the CFL condition (3.40) and the inequality from (3.20) to bound  $\|\psi_h^0\|_{(\text{HHO})'}$ , we obtain

$$\begin{aligned} \frac{1}{\Delta t} |(e_{\mathcal{T}}^1, w_{\mathcal{T}})_{\Omega}| &\leq \frac{\Delta t}{2} \|\psi_h^0\|_{(\text{HHO})'} \|(w_{\mathcal{T}}, 0)\|_{\text{HHO}} + \frac{\Delta t^2}{2} \|\partial_t^3 u\|_{C^0([0; t^1]; L^2(\Omega))} \|w_{\mathcal{T}}\|_{\Omega} \\ &\leq C(|u_0|_{*,h} + \Delta t^2 \|\partial_t^3 u\|_{C^0([0; t^1]; L^2(\Omega))}) \|w_{\mathcal{T}}\|_{\Omega}. \end{aligned}$$

Since  $\|e_{\mathcal{T}}^1\|_{\Omega} \leq \sup_{w_{\mathcal{T}} \in \mathcal{U}_{\mathcal{T}}^1} \frac{(e_{\mathcal{T}}^1, w_{\mathcal{T}})_{\Omega}}{\|w_{\mathcal{T}}\|_{\Omega}}$ , we conclude that  $\frac{1}{\Delta t} \|e_{\mathcal{T}}^1\|_{\Omega} \leq C(|u_0|_{*,h} + \Delta t^2 \|\partial_t^3 u\|_{C^0([0; t^1]; L^2(\Omega))})$ . Since  $e_{\mathcal{T}}^0 = 0$ , this finally gives

$$A_1 = \frac{1}{2\Delta t^2} \|e_{\mathcal{T}}^1 - e_{\mathcal{T}}^0\|_{\Omega}^2 \leq C(|u_0|_{*,h}^2 + \Delta t^4 \|\partial_t^3 u\|_{C^0([0; t^1]; L^2(\Omega))}^2). \quad (3.64)$$

*Step 2: Bound on  $\|(0, e_{\mathcal{F}}^0)\|_{\text{HHO}}$  and  $\|(0, e_{\mathcal{F}}^1)\|_{\text{HHO}}$ .* Recalling (3.21), we have

$$\|(0, e_{\mathcal{F}}^0)\|_{\text{HHO}} \leq C|u_0|_{*,h}.$$

Moreover, to bound  $e_{\mathcal{F}}^1$ , we use the following equality (see (3.52) for  $n = 1$  and  $w_{\mathcal{F}} = e_{\mathcal{F}}^1$ ):

$$a_h((0, e_{\mathcal{F}}^1), (0, e_{\mathcal{F}}^1)) = \psi_h^1((0, e_{\mathcal{F}}^1)) - a_h((e_{\mathcal{T}}^1, 0), (0, e_{\mathcal{F}}^1)).$$

Using the coercivity and boundedness of the bilinear form  $a_h$ , we obtain

$$\|(0, e_{\mathcal{F}}^1)\|_{\text{HHO}}^2 \leq \frac{1}{\alpha} a_h((0, e_{\mathcal{F}}^1), (0, e_{\mathcal{F}}^1)) \leq \frac{1}{\alpha} \left( \|\psi_h^1\|_{(\text{HHO})'} + \varpi \|(e_{\mathcal{T}}^1, 0)\|_{\text{HHO}} \right) \|(0, e_{\mathcal{F}}^1)\|_{\text{HHO}},$$

so that

$$\|(0, e_{\mathcal{F}}^1)\|_{\text{HHO}} \leq C(\|\psi_h^1\|_{(\text{HHO})'} + \|(e_{\mathcal{T}}^1, 0)\|_{\text{HHO}}). \quad (3.65)$$

To bound the second term on the right-hand-side, we invoke again the coercivity and boundedness of  $a_h$  together with the discrete inverse inequality from Lemma 2.2.1 and the CFL condition (3.40), to obtain (since  $e_{\mathcal{F}}^0 = 0$ )

$$\begin{aligned} \alpha \|(e_{\mathcal{T}}^1, 0)\|_{\text{HHO}}^2 &\leq a_h((e_{\mathcal{T}}^1 - e_{\mathcal{T}}^0, 0), (e_{\mathcal{T}}^1 - e_{\mathcal{T}}^0, 0)) \leq \varpi \|(e_{\mathcal{T}}^1 - e_{\mathcal{T}}^0, 0)\|_{\text{HHO}}^2 \\ &\leq C \mu_{\sharp}^2 h^{-2} \|e_{\mathcal{T}}^1 - e_{\mathcal{T}}^0\|_{\Omega}^2 \\ &\leq C' \frac{1}{\Delta t^2} \|e_{\mathcal{T}}^1 - e_{\mathcal{T}}^0\|_{\Omega}^2, \end{aligned} \quad (3.66)$$

which can be bounded using (3.64). Since  $\|\psi_h^1\|_{\text{HHO}'} \leq c_*|u(t^1)|_{*,h}$ , this gives altogether

$$\|(0, e_{\mathcal{F}}^1)\|_{\text{HHO}} \leq C\left(|u_0|_{*,h} + |u(t^1)|_{*,h} + \Delta t^2 \|\partial_t^3 u\|_{C^0([0;t^1];L^2(\Omega))}\right). \quad (3.67)$$

*Step 3: Bound on  $\|\hat{e}_h^{\frac{1}{2}}\|_{\text{HHO}}$ .* The symmetry of  $a_h$ , the fact that  $e_{\mathcal{T}}^0 = 0$  and the above arguments give

$$\begin{aligned} a_h(\hat{e}_h^1 + \hat{e}_h^0, \hat{e}_h^1 + \hat{e}_h^0) &\leq 2a_h((e_{\mathcal{T}}^1 + e_{\mathcal{T}}^0, 0), (e_{\mathcal{T}}^1 + e_{\mathcal{T}}^0, 0)) + 2a_h((0, e_{\mathcal{F}}^1 + e_{\mathcal{F}}^0), (0, e_{\mathcal{F}}^1 + e_{\mathcal{F}}^0)) \\ &\leq 2a_h((e_{\mathcal{T}}^1, 0), (e_{\mathcal{T}}^1, 0)) + 4a_h((0, e_{\mathcal{F}}^1), (0, e_{\mathcal{F}}^1)) + 4a_h((0, e_{\mathcal{F}}^0), (0, e_{\mathcal{F}}^0)) \\ &= 2a_h((e_{\mathcal{T}}^1 - e_{\mathcal{T}}^0, 0), (e_{\mathcal{T}}^1 - e_{\mathcal{T}}^0, 0)) + 4a_h((0, e_{\mathcal{F}}^1), (0, e_{\mathcal{F}}^1)) + 4a_h((0, e_{\mathcal{F}}^0), (0, e_{\mathcal{F}}^0)) \\ &\leq C(\|(e_{\mathcal{T}}^1 - e_{\mathcal{T}}^0, 0)\|_{\text{HHO}}^2 + \|(0, e_{\mathcal{F}}^1)\|_{\text{HHO}}^2 + \|(0, e_{\mathcal{F}}^0)\|_{\text{HHO}}^2) \\ &\leq C'\left(\frac{1}{\Delta t^2} \|(e_{\mathcal{T}}^1 - e_{\mathcal{T}}^0, 0)\|_{\Omega}^2 + \|(0, e_{\mathcal{F}}^1)\|_{\text{HHO}}^2 + \|(0, e_{\mathcal{F}}^0)\|_{\text{HHO}}^2\right) \\ &\leq C''(|u_0|_{*,h}^2 + |u(t^1)|_{*,h}^2 + \Delta t^4 \|\partial_t^3 u\|_{C^0([0;t^1];L^2(\Omega))}^2), \end{aligned}$$

owing to (3.64), (3.65) and (3.67). Since  $\|\hat{e}_h^{\frac{1}{2}}\|_{\text{HHO}}^2 \leq \frac{1}{\alpha} a_h(\hat{e}_h^{\frac{1}{2}}, \hat{e}_h^{\frac{1}{2}})$ , we conclude that

$$\|\hat{e}_h^{\frac{1}{2}}\|_{\text{HHO}} \leq C\left(|u_0|_{*,h} + |u(t^1)|_{*,h} + \Delta t^2 \|\partial_t^3 u\|_{C^0([0;t^1];L^2(\Omega))}\right). \quad (3.68)$$

*Step 4: Bound on  $A_2$  and  $A_3$ .* The bound (3.68) directly yields

$$|A_2| \leq C\left(|u_0|_{*,h}^2 + |u(t^1)|_{*,h}^2 + \Delta t^4 \|\partial_t^3 u\|_{C^0([0;t^1];L^2(\Omega))}^2\right).$$

Finally, recalling the arguments from the proofs of Lemmas 3.2.12 and 3.2.13, we have

$$\begin{aligned} \left|(\kappa^1, e_{\mathcal{T}}^{\frac{1}{2}})_{\Omega}\right| &\leq C\Theta\Delta t^2 \|\partial_t^4 u\|_{C^0([t^0;t^2];L^2(\Omega))} \|\hat{e}_h^{\frac{1}{2}}\|_{\text{HHO}}, \\ \left|\psi_h^1(\hat{e}_h^{\frac{1}{2}})\right| &\leq C|u(t^1)|_{*,h} \|\hat{e}_h^{\frac{1}{2}}\|_{\text{HHO}}. \end{aligned}$$

The above estimate of  $\|\hat{e}_h^{\frac{1}{2}}\|_{\text{HHO}}$  then gives

$$\begin{aligned} |A_3| &\leq \left|(\kappa^1, e_{\mathcal{T}}^{\frac{1}{2}})_{\Omega}\right| + \left|\psi_h^1(\hat{e}_h^{\frac{1}{2}})\right| \\ &\leq C(|u(t^1)|_{*,h} + \Theta\Delta t^2 \|\partial_t^4 u\|_{C^0([t^0;t^2];L^2(\Omega))})(|u_0|_{*,h} + |u(t^1)|_{*,h} + \Delta t^2 \|\partial_t^3 u\|_{C^0([0;t^1];L^2(\Omega))}) \\ &\leq C'\left\{|u_0|_{*,h}^2 + |u(t^1)|_{*,h}^2 + \Delta t^4 (\Theta^2 \|\partial_t^4 u\|_{C^0([t^0;t^2];L^2(\Omega))}^2 + \|\partial_t^3 u\|_{C^0([0;t^1];L^2(\Omega))}^2)\right\}, \end{aligned}$$

where the last bound follows from Young's inequality. Gathering the two estimates from Step 4 with the one on  $A_1$  established in Step 1 proves the claim.  $\square$

### 3.2.2.5 Proof of Theorem 3.2.6

We are now ready to complete the proof of Theorem 3.2.6. On the one hand, using the convexity result of Lemma 3.2.10, we have, for all  $n \in \{1:N-1\}$ ,

$$\frac{1}{4} \|\delta e_{\mathcal{T}}^{n+\frac{1}{2}}\|_{\Omega}^2 + \alpha \|\hat{e}_h^{n+\frac{1}{2}}\|_{\text{HHO}}^2 + \frac{\Delta t^2}{8} a_h((0, \delta e_{\mathcal{F}}^{n+\frac{1}{2}}), (0, \delta e_{\mathcal{F}}^{n+\frac{1}{2}})) \leq \tilde{\mathcal{E}}_h^{n+\frac{1}{2}}. \quad (3.69)$$

On the other hand, Lemma 3.2.11 gives, for all  $n \in \{1:N-1\}$ ,

$$\begin{aligned} \tilde{\mathcal{E}}_h^{n+\frac{1}{2}} &\leq |\mathcal{Z}_{\psi}^n| + |\mathcal{Z}_{\kappa}^n| + |\mathcal{Z}_{\text{IC}}^0| \\ &\leq \max_{n \in \{1:N-1\}} |\mathcal{Z}_{\psi}^n| + \max_{n \in \{1:N-1\}} |\mathcal{Z}_{\kappa}^n| + |\mathcal{Z}_{\text{IC}}^0|. \end{aligned}$$

Combining the above two inequalities yields, for all  $n \in \{1:N-1\}$ ,

$$\frac{1}{4} \|\delta e_{\mathcal{T}}^{n+\frac{1}{2}}\|_{\Omega}^2 + \alpha \|\hat{e}_h^{n+\frac{1}{2}}\|_{\text{HHO}}^2 + \frac{\Delta t^2}{8} a_h((0, \delta e_{\mathcal{F}}^{n+\frac{1}{2}}), (0, \delta e_{\mathcal{F}}^{n+\frac{1}{2}})) \leq \max_{n \in \{1:N-1\}} |\mathcal{Z}_{\psi}^n| + \max_{n \in \{1:N-1\}} |\mathcal{Z}_{\kappa}^n| + |\mathcal{Z}_{\text{IC}}^0|.$$

Since the right-hand-side of this inequality does not depend on  $n$  and since each term on the left-hand-side is positive, we infer that

$$\begin{aligned} \frac{1}{4} \max_{n \in \{1:N-1\}} \|\delta e_{\mathcal{T}}^{n+\frac{1}{2}}\|_{\Omega}^2 &\leq \max_{n \in \{1:N-1\}} |\mathcal{Z}_{\psi}^n| + \max_{n \in \{1:N-1\}} |\mathcal{Z}_{\kappa}^n| + |\mathcal{Z}_{\text{IC}}^0|, \\ \alpha \max_{n \in \{1:N-1\}} \|\hat{e}_h^{n+\frac{1}{2}}\|_{\text{HHO}}^2 &\leq \max_{n \in \{1:N-1\}} |\mathcal{Z}_{\psi}^n| + \max_{n \in \{1:N-1\}} |\mathcal{Z}_{\kappa}^n| + |\mathcal{Z}_{\text{IC}}^0|, \\ \frac{\Delta t^2}{8} \max_{n \in \{1:N-1\}} a_h((0, \delta e_{\mathcal{F}}^{n+\frac{1}{2}}), (0, \delta e_{\mathcal{F}}^{n+\frac{1}{2}})) &\leq \max_{n \in \{1:N-1\}} |\mathcal{Z}_{\psi}^n| + \max_{n \in \{1:N-1\}} |\mathcal{Z}_{\kappa}^n| + |\mathcal{Z}_{\text{IC}}^0|. \end{aligned}$$

This gives

$$\begin{aligned} \frac{1}{4} \max_{n \in \{1:N-1\}} \|\delta e_{\mathcal{T}}^{n+\frac{1}{2}}\|_{\Omega}^2 + \alpha \max_{n \in \{1:N-1\}} \|\hat{e}_h^{n+\frac{1}{2}}\|_{\text{HHO}}^2 + \frac{\Delta t^2}{8} \max_{n \in \{1:N-1\}} a_h((0, \delta e_{\mathcal{F}}^{n+\frac{1}{2}}), (0, \delta e_{\mathcal{F}}^{n+\frac{1}{2}})) \\ \leq 3 \left\{ \max_{n \in \{1:N-1\}} |\mathcal{Z}_{\psi}^n| + \max_{n \in \{1:N-1\}} |\mathcal{Z}_{\kappa}^n| + |\mathcal{Z}_{\text{IC}}^0| \right\}. \end{aligned} \quad (3.70)$$

Let us collect the results of Lemmas 3.2.12 and 3.2.13 and apply Young's inequality to the upper bounds on  $\mathcal{Z}_{\psi}^n$  and  $\mathcal{Z}_{\kappa}^n$ . We obtain

$$\begin{aligned} \max_{n \in \{1:N-1\}} |\mathcal{Z}_{\psi}^n| &\leq C_{\psi} \left( |\partial_t u|_{L^1([t^1, t^N]; *, h)}^2 + |u|_{C^0([t^1, t^N-1]; *, h)}^2 \right) + \frac{\alpha}{8} \max_{n \in \{1:N-1\}} \|\hat{e}_h^{n+\frac{1}{2}}\|_{\text{HHO}}^2 \\ &\quad + \frac{\Delta t^2}{24} \max_{n \in \{1:N-1\}} a_h((0, \delta e_{\mathcal{F}}^{n+\frac{1}{2}}), (0, \delta e_{\mathcal{F}}^{n+\frac{1}{2}})), \end{aligned}$$

for the space consistency error, and

$$\max_{n \in \{1:N-1\}} |\mathcal{Z}_{\kappa}^n| \leq C_{\kappa} \Theta^2 \Delta t^4 \left( \|\partial_t^5 u\|_{L^1([0; t^N]; L^2(\Omega))}^2 + \|\partial_t^4 u\|_{C^0([0; t^N]; L^2(\Omega))}^2 \right) + \frac{\alpha}{8} \max_{n \in \{1:N-1\}} \|\hat{e}_h^{n+\frac{1}{2}}\|_{\text{HHO}}^2,$$

for the time consistency error, with constants  $C_{\psi}$  and  $C_{\kappa}$ . Using these two bounds and the bound on  $\mathcal{Z}_{\text{IC}}^0$  from Lemma 3.2.14 in (3.70) and rearranging the terms involving  $\max_{n \in \{1:N-1\}} \|\hat{e}_h^{n+\frac{1}{2}}\|_{\text{HHO}}^2$ , we obtain

$$\begin{aligned} \frac{1}{4} \max_{n \in \{1:N-1\}} \|\delta e_{\mathcal{T}}^{n+\frac{1}{2}}\|_{\Omega}^2 + \frac{1}{4} \alpha \max_{n \in \{1:N-1\}} \|\hat{e}_h^{n+\frac{1}{2}}\|_{\text{HHO}}^2 &\leq C_{\psi} \left( |\partial_t u|_{L^1([t^1, t^N]; *, h)}^2 + |u|_{C^0([t^1, t^N-1]; *, h)}^2 \right) \\ &\quad + C_{\kappa} \Theta^2 \Delta t^4 \left( \|\partial_t^5 u\|_{L^1([0; t^N]; L^2(\Omega))}^2 + \|\partial_t^4 u\|_{C^0([0; t^N]; L^2(\Omega))}^2 \right) \\ &\quad + C_{\text{IC}} \left( |u_0|_{*, h}^2 + |u(t^1)|_{*, h}^2 + \Delta t^4 \left( \|\partial_t^3 u\|_{C^0([0; t^1]; L^2(\Omega))}^2 + \Theta^2 \|\partial_t^4 u\|_{C^0([t^0; t^2]; L^2(\Omega))}^2 \right) \right) \end{aligned}$$

Taking the square root and reorganizing the terms gives

$$\begin{aligned} \max_{n \in \{1:N-1\}} \|\delta e_{\mathcal{T}}^{n+\frac{1}{2}}\|_{\Omega} + \max_{n \in \{1:N-1\}} \|\hat{e}_h^{n+\frac{1}{2}}\|_{\text{HHO}} &\leq C_1 \left\{ |\partial_t u|_{L^1([t^1, t^N]; *, h)} + |u|_{C^0([0; t^N-1]; *, h)} \right\} \\ &\quad + C_2 \Delta t^2 \left\{ \|\partial_t^3 u\|_{C^0([0; t^1]; L^2(\Omega))} + \Theta \|\partial_t^4 u\|_{C^0([0; t^N]; L^2(\Omega))} + \Theta \|\partial_t^5 u\|_{L^1([0; t^N]; L^2(\Omega))} \right\}, \end{aligned}$$

which yields the claim.

### 3.3 Numerical experiments

The theoretical convergence orders established in Sections 3.1 and 3.2 can be verified on an analytical test case. Let us separate the verification of the space convergence order from that of the time convergence order, by considering first a polynomial-in-time test case, and then a more general test case.

**Verification of the space convergence order** Let us set  $\Omega := (0, 2)^2$ , the source term  $f(x, y, t) := 2(\pi^2 t^2 + 1) \sin(\pi x) \sin(\pi y)$ , the initial conditions  $u_0(x, y) := 0$  and  $v_0(x, y) := 0$ , and a constant speed of sound  $\mu := 1$ . The solution to (3.1) is

$$u(x, y, t) = t^2 \sin(\pi x) \sin(\pi y). \quad (3.71)$$

Notice that the homogeneous Dirichlet boundary conditions are verified at all times. The solution (3.71) satisfies  $\partial_t^3 u = 0$ , so that the leapfrog scheme does not produce any discretization error. The only error term is due to the space semi-discretization by the HHO method and the expected order of convergence of the energy error is  $h^{k+1}$  and that of the  $L^2$ -error is  $h^{k+2}$ .

The shape of the solution is displayed in Figure 3.1a at  $t = 0.3$ . Let  $\hat{U}(t) := (\mathbf{U}_{\mathcal{T}}(t), \mathbf{U}_{\mathcal{F}}(t))$  collect the degrees of freedom of the exact solution  $u(t)$  projected onto the HHO space, i.e. the degrees of freedom of  $\hat{I}_h^{k,l}(u(t))$  at each time step  $t \in \bar{\mathcal{J}}$ , with  $\mathbf{U}_{\mathcal{T}}(t)$  collecting the degrees of freedom of  $\Pi_{\mathcal{T}}^l(u(t))$  and  $\mathbf{U}_{\mathcal{F}}(t)$  those of  $\Pi_{\mathcal{F}}^k(u(t))$ . The discrete errors are computed globally, using the mass and stiffness matrices  $\mathcal{M}$  and  $\mathcal{A}$  and the following definitions:

$$\|e_{\mathcal{T}}^n\|_{\Omega}^2 = (\mathbf{U}_{\mathcal{T}}(t^n) - \mathbf{U}_{\mathcal{T}}^n)^\top \mathcal{M} (\mathbf{U}_{\mathcal{T}}(t^n) - \mathbf{U}_{\mathcal{T}}^n), \quad a_h(\hat{e}_h^n, \hat{e}_h^n) = (\hat{U}(t^n) - \hat{U}^n)^\top \mathcal{A} (\hat{U}(t^n) - \hat{U}^n), \quad (3.72)$$

which correspond to the  $L^2$ - and  $H^1$ -errors evaluated at the discrete time nodes.

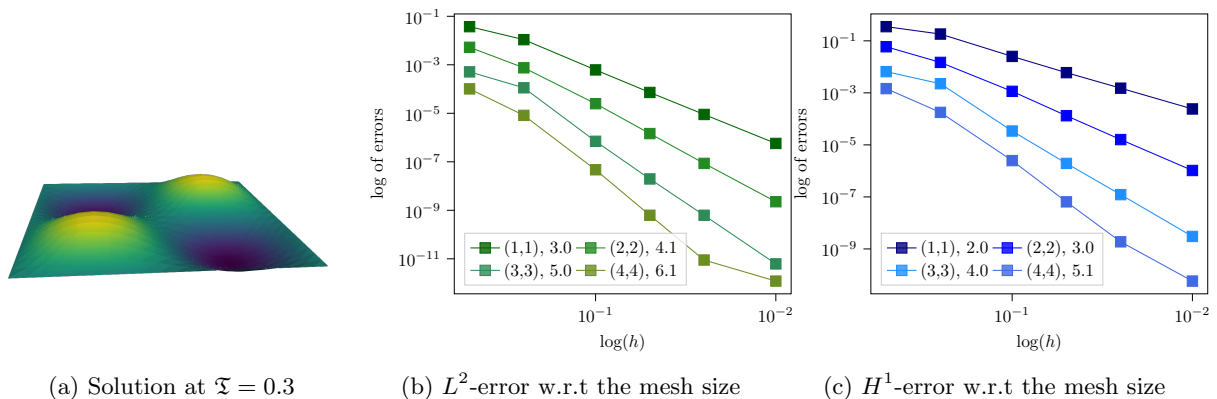


Figure 3.1: Linear wave problem, analytical test case with polynomial time dependence - Convergence of the  $L^2$ - and  $H^1$ -errors as a function of the mesh size, polynomial orders  $k \in \{1, 2, 3, 4\}$ , equal-order setting,  $\mathfrak{T} = 0.3$ .

Figures 3.1b and 3.1c display these errors on a series of refined triangular meshes for  $k \in \{1, 2, 3, 4\}$  in the equal-order setting, at the final time  $\mathfrak{T} = 0.3$ . The convergence order (obtained by linear regression) is also reported in the legends. The time step is taken so that the CFL condition is verified with a margin of 20%. As expected, the convergence orders are  $(k + 1)$  for the  $H^1$ -error and  $(k + 2)$  for the  $L^2$ -error. The last point of the  $L^2$ -error curve for  $k = 4$  seems to reach a limit, which is probably due to floating point precision. The mixed-order setting delivers the same results which are not displayed for brevity.

**Verification of the space-time convergence order** Let us set  $\Omega := (0, 1)^2$ , the source term  $f(x, y, t) := 0$ , the initial conditions  $u_0(x, y) := 0$  and  $v_0(x, y) := \sin(\pi x) \sin(\pi y)$ , and a constant speed

of sound  $\mu := 1$ . The solution to (3.1) is now

$$u(x, y, t) = \frac{1}{\sqrt{2\pi}} \sin(\sqrt{2}\pi t) \sin(\pi x) \sin(\pi y). \quad (3.73)$$

This solution is non-polynomial in space and time so that both time and space discretization errors are present, with expected orders  $\Delta t^2$  in time,  $h^{k+1}$  in space for the  $H^1$ -error and  $h^{k+2}$  in space for the  $L^2$ -error. If the time step is taken so that it respects the CFL condition, and the ratio  $\frac{\Delta t}{h}$  remains constant, the time discretization error will dominate the space discretization error on the finest meshes, except for the  $H^1$ -error in the lowest-order settings ( $k \in \{0, 1\}$ ) and the  $L^2$ -error for  $k = 0$ .

This phenomenon is illustrated in Figure 3.2, where the time step  $\Delta t$  is proportional to the mesh size  $h$ , and verifies  $\Delta t := 0.8\eta\mu^{-1}h$  (see the CFL condition (3.40)). The convergence rate of the  $L^2$ -error for polynomial orders  $(k, l) = (0, 0)$  and the  $H^1$ -error for the polynomial orders  $(k, l) = (0, 0)$  and  $(k, l) = (1, 1)$  are unaffected. But for higher orders, the optimal space convergence rate is observed on the coarse meshes, whereas a second-order decay rate dictated by the time discretization error is observed on the finest meshes. The transition between the two regimes is smooth.

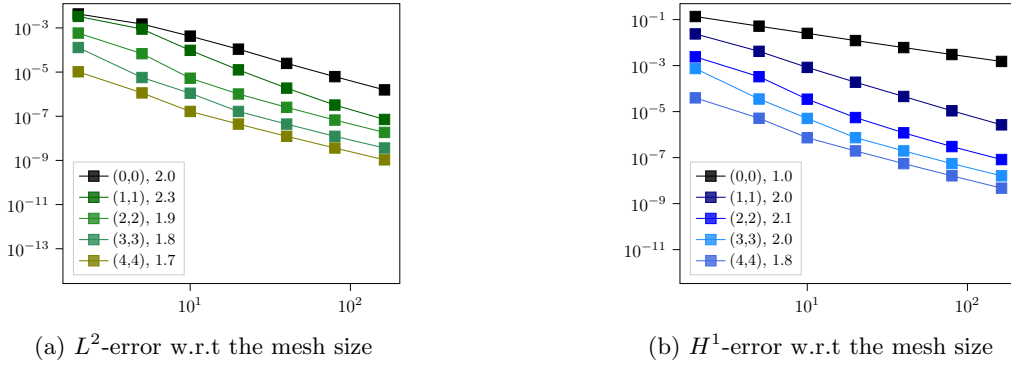


Figure 3.2: Linear wave problem, analytical test case with nonpolynomial time dependence - Convergence of the  $L^2$ - and  $H^1$ -errors as a function of the mesh size, polynomial orders  $k \in \{0, 1, 2, 3, 4\}$ , equal-order setting,  $\mathfrak{T} = 0.2$ ,  $\Delta t = 0.8\eta\mu^{-1}h$

If the time step is taken to be constant for all the meshes in the sequence, but much smaller than the critical time step, the time discretization error becomes negligible w.r.t the space discretization error and higher convergence rates in space are recovered, see Figure 3.3.

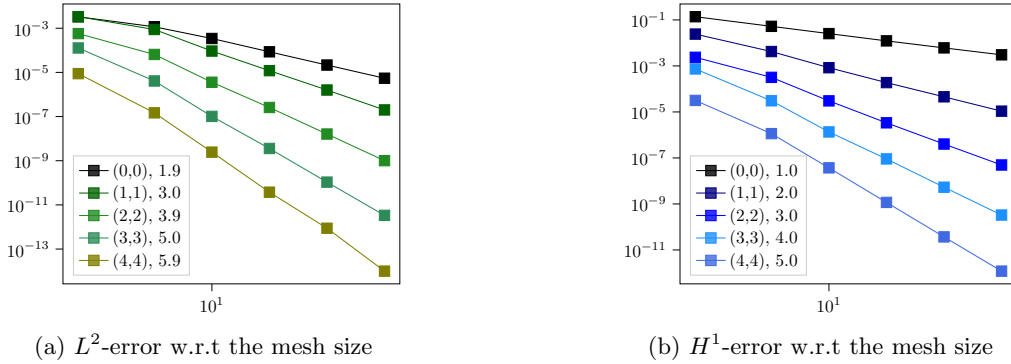


Figure 3.3: Linear wave problem, analytical test case with nonpolynomial time dependence - Convergence of the  $L^2$ - and  $H^1$ -errors as a function of the mesh size, polynomial orders  $k \in \{0, 1, 2, 3, 4\}$ , equal-order setting,  $\mathfrak{T} = 0.2$ ,  $\Delta t \ll \eta\mu^{-1}h$ .



### 3.4 Other possible schemes

The leapfrog scheme is one of the most frequently used time-schemes for the second-order formulation in time of the wave equation. But other schemes are available. An alternative is also to write the wave equation as a first-order system. Although in this Thesis we focus on the leapfrog scheme, we briefly compare here the leapfrog scheme with other time schemes, in particular regarding the possibility of a fully explicit time-stepping.

**Second-order formulation.** In the second-order formulation in time, [34] also considers the space semi-discretization of the acoustic wave equation with the HHO method and a Newmark scheme in time with real parameters  $\beta$  and  $\gamma$ . This reference reports numerical examples with optimal convergence orders. The Newmark scheme is an implicit scheme, unless  $\beta = 0$ , in which case the scheme reduces to the leapfrog scheme. The Newmark scheme is unconditionally stable if  $\frac{1}{2} < \gamma < 2\beta$  (classically one takes  $\beta = \frac{1}{4}$  and  $\gamma = \frac{1}{2}$ ). In the context of the HHO space semi-discretization, the Newmark scheme works as follows: Knowing  $\hat{u}_h^n, \hat{v}_h^n, \hat{a}_h^n$  from previous time steps (notice that here the speed and acceleration are hybrid quantities all living in  $\widehat{\mathcal{U}}_{h,0}^{l,k}$ ), we perform the following steps:

1. Predictor step: Compute

$$\hat{u}_h^{*,n} := \hat{u}_h^n + \Delta t \hat{v}_h^n + \frac{1}{2} \Delta t^2 (1 - 2\beta) \hat{a}_h^n, \quad (3.74a)$$

$$\hat{v}_h^{*,n} := \hat{v}_h^n + \Delta t (1 - \gamma) \hat{a}_h^n. \quad (3.74b)$$

2. Find the acceleration  $\hat{a}_h^{n+1}$  at time  $t^{n+1}$  such that

$$(a_{\mathcal{T}}^{n+1}, w_{\mathcal{T}})_{\Omega} + \beta \Delta t^2 a_h(\hat{a}_h^{n+1}, \hat{w}_h) = (f^{n+1}, w_{\mathcal{T}})_{\Omega} - a_h(\hat{u}_h^{*,n}, \hat{w}_h), \quad \forall \hat{w}_h \in \widehat{\mathcal{U}}_{h,0}^{l,k}. \quad (3.75)$$

3. Corrector step: Update the potential and the speed as follows:

$$\hat{u}_h^{n+1} := \hat{u}_h^{*,n} + \beta \Delta t^2 \hat{a}_h^{n+1}, \quad (3.76a)$$

$$\hat{v}_h^{n+1} := \hat{v}_h^{*,n} + \gamma \Delta t \hat{a}_h^{n+1}. \quad (3.76b)$$

The scheme is initialized with  $\hat{u}_h^0 = \hat{I}_h^{k,l}(u_0)$ ,  $\hat{v}_h^0 = \hat{I}_h^{k,l}(v_0)$ , the initial cell acceleration  $a_{\mathcal{T}}^0 \in \mathcal{U}_{\mathcal{T}}^l$  is found by solving

$$(a_{\mathcal{T}}^0, w_{\mathcal{T}})_{\Omega} = (f^0, w_{\mathcal{T}})_{\Omega} - a_h(\hat{u}_h^0, (w_{\mathcal{T}}, 0)), \quad \forall w_{\mathcal{T}} \in \mathcal{U}_{\mathcal{T}}^l,$$

and then the initial face acceleration  $a_{\mathcal{F}}^0 \in \mathcal{U}_{\mathcal{F},0}^k$  is found by solving  $a_h((a_{\mathcal{T}}^0, a_{\mathcal{F}}^0), (0, w_{\mathcal{F}})) = 0$  for all  $w_{\mathcal{F}} \in \mathcal{U}_{\mathcal{F},0}^k$ . The second step above, see (3.75), contains the implicit part of the scheme, since it requires to solve the algebraic problem

$$\left( \begin{bmatrix} \mathcal{M} & 0 \\ 0 & 0 \end{bmatrix} + \beta \Delta t^2 \begin{bmatrix} \mathcal{A}_{\mathcal{T}\mathcal{T}} & \mathcal{A}_{\mathcal{T}\mathcal{F}} \\ \mathcal{A}_{\mathcal{F}\mathcal{T}} & \mathcal{A}_{\mathcal{F}\mathcal{F}} \end{bmatrix} \right) \begin{pmatrix} \mathbf{A}_{\mathcal{T}}^{n+1} \\ \mathbf{A}_{\mathcal{F}}^{n+1} \end{pmatrix} = \begin{bmatrix} \mathbf{F}_{\mathcal{T}}^{n+1} - \mathcal{A}_{\mathcal{T}\mathcal{T}} \mathbf{U}_{\mathcal{T}}^{*,n} - \mathcal{A}_{\mathcal{T}\mathcal{F}} \mathbf{U}_{\mathcal{F}}^{*,n} \\ -\mathcal{A}_{\mathcal{F}\mathcal{T}} \mathbf{U}_{\mathcal{T}}^{*,n} - \mathcal{A}_{\mathcal{F}\mathcal{F}} \mathbf{U}_{\mathcal{F}}^{*,n} \end{bmatrix}. \quad (3.77)$$

This problem requires the inversion of a matrix that is not block-diagonal. Since  $\mathcal{M} + \beta \Delta t^2 \mathcal{A}_{\mathcal{T}\mathcal{T}}$  is block diagonal, the cell unknowns can be eliminated via static condensation, in the same manner as in the static case. This, however, still requires the inversion of a face-face matrix that is not block-diagonal.

**First-order system.** The linear acoustic wave equation can also be formulated as a space-time first-order system. The unknowns are the velocity  $v := \partial_t u$  and the dual vector-valued variable  $\sigma := \mu \nabla u$ .

The system reads as follows:

$$\partial_t \boldsymbol{\sigma} - \mu \nabla v = \mathbf{0}, \quad \text{in } \Omega \times J, \quad (3.78a)$$

$$\partial_t v - \nabla \cdot (\mu \boldsymbol{\sigma}) = f, \quad \text{in } \Omega \times J, \quad (3.78b)$$

$$v|_{t=0} = v_0, \quad \text{in } \Omega, \quad (3.78c)$$

$$\boldsymbol{\sigma}|_{t=0} = \mu \nabla u_0, \quad \text{in } \Omega, \quad (3.78d)$$

$$v = 0, \quad \text{on } \partial\Omega \times J. \quad (3.78e)$$

For a.e.  $t \in J$ , a weak form for this problem is to look for  $v \in H^1(J; L^2(\Omega)) \cap L^2(J; H_0^1(\Omega))$  and  $\boldsymbol{\sigma} \in H^1(J; \mathbf{L}^2(\Omega))$  such that

$$\begin{aligned} (\partial_t \boldsymbol{\sigma}(t), \mathbf{q})_\Omega - (\mu \nabla v(t), \mathbf{q})_\Omega &= 0, & \forall \mathbf{q} \in \mathbf{L}^2(\Omega), \\ (\partial_t v(t), w)_\Omega + (\mu \boldsymbol{\sigma}(t), \nabla w)_\Omega &= (f(t), w)_\Omega, & \forall w \in H_0^1(\Omega). \end{aligned} \quad (3.79)$$

The HHO unknowns are in this case the hybrid velocity  $\hat{v}_h \in C^1(\bar{J}; \hat{\mathcal{U}}_{h,0}^{l,k})$  and the cellwise dual variable  $\boldsymbol{\sigma}_\mathcal{T} \in C^1(\bar{J}; \mathcal{U}_\mathcal{T}^l)$ . The space semi-discrete system reads

$$\begin{aligned} (\partial_t \boldsymbol{\sigma}_\mathcal{T}(t), \mathbf{q})_\Omega - (\mu \mathbf{G}_\mathcal{T}^k(\hat{v}_h(t)), \mathbf{q})_\Omega &= 0, & \forall \mathbf{q} \in \mathcal{U}_\mathcal{T}^l, \\ (\partial_t v_\mathcal{T}(t), w_\mathcal{T})_\Omega + (\mu \boldsymbol{\sigma}_\mathcal{T}(t), \mathbf{G}_\mathcal{T}^k(\hat{w}_h))_\Omega + s_h(\hat{v}_h(t), \hat{w}_h) &= (f(t), w_\mathcal{T})_\Omega, & \forall \hat{w}_h \in \hat{\mathcal{U}}_{h,0}^{l,k}, \end{aligned} \quad (3.80)$$

with the initial conditions  $\boldsymbol{\sigma}_\mathcal{T}(0) \in C^1(\bar{J}; \mathcal{U}_\mathcal{T}^l) = \mu \mathbf{G}_\mathcal{T}^k(\hat{I}_h^{k,l}(u_0))$  and  $v_\mathcal{T}(0) = \Pi_\mathcal{T}^l(v_0)$ .

Let  $\mathcal{M}$  be the mass matrix associated with the inner product of  $\mathbf{L}^2(\Omega)$ ,  $\mathcal{G}$  the rectangular matrix associated with the gradient reconstruction operator, decomposed into the blocks  $\mathcal{G}_\mathcal{T}$  and  $\mathcal{G}_\mathcal{F}$  related to the cell and face unknowns, respectively. Let  $\mathbf{V}_\mathcal{T}(t), \mathbf{V}_\mathcal{F}(t)$  respectively be the cell and face components vectors of the scalar velocity  $\hat{v}_h(t)$  and let  $\mathbf{S}_\mathcal{T}(t)$  be the component vector of  $\boldsymbol{\sigma}_\mathcal{T}(t)$ . The algebraic formulation of (3.80) reads

$$\begin{bmatrix} \mathcal{M} & 0 & 0 \\ 0 & \mathcal{M} & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} \partial_t \mathbf{S}_\mathcal{T}(t) \\ \partial_t \mathbf{V}_\mathcal{T}(t) \\ \cdot \end{bmatrix} + \begin{bmatrix} 0 & -\mathcal{G}_\mathcal{T} & -\mathcal{G}_\mathcal{F} \\ \mathcal{G}_\mathcal{T}^\dagger & \mathcal{S}_\mathcal{T}\mathcal{T} & \mathcal{S}_\mathcal{T}\mathcal{F} \\ \mathcal{G}_\mathcal{F}^\dagger & \mathcal{S}_\mathcal{F}\mathcal{T} & \mathcal{S}_\mathcal{F}\mathcal{F} \end{bmatrix} \begin{bmatrix} \mathbf{S}_\mathcal{T}(t) \\ \mathbf{V}_\mathcal{T}(t) \\ \mathbf{V}_\mathcal{F}(t) \end{bmatrix} = \begin{bmatrix} 0 \\ \mathbf{F}_\mathcal{T}(t) \\ 0 \end{bmatrix}, \quad (3.81)$$

with the same stabilization matrix  $\mathcal{S}$  as for the second-order formulation in time. In the two first lines, the mass matrices are block-diagonal and the block equations are easily solved. In the mixed-order setting,  $\mathcal{S}_\mathcal{F}\mathcal{F}$  is block-diagonal, so that the static coupling resulting from the last equation is easily solved in this case. Hence, any time explicit scheme will lead to a fully explicit problem. For instance, an explicit Runge–Kutta (RK) scheme can be implemented. The implementation of RK schemes is detailed in [34], and the proof of convergence of the space semi-discrete formulation can be found in [36]. In the equal-order setting, the problem coupling the face unknowns is not block-diagonal, so that implicit RK schemes can be preferred.

### 3.5 Conclusion

The HHO space semi-discretization of the linear acoustic wave equation has been presented in this chapter. In the second-order formulation, this leads to a static coupling between the cell and the face unknowns. The proof of the rates of convergence, established in [36], have been detailed. The time discretization with the leapfrog scheme, along with the proof of convergence in fully discrete energy norm, has also been developed. Moreover, numerical illustrations have been given. The high orders of convergence in space ( $h^{k+2}$  for the  $L^2$ -norm and  $h^{k+1}$  for the energy norm) observed in the static case are retrieved both in the space semi-discrete and the fully discrete cases. A consequence of the static coupling is that any classical explicit time-scheme leads to a semi-implicit time stepping when combined with the HHO space

semi-discretization. In the mixed-order setting, one can use the first-order formulation, briefly presented in Section 3.4, with for instance an explicit Runge–Kutta scheme. In the equal-order setting, however, the first-order formulation does not provide a straightforward explicit time-scheme.

Considering our final goal of a robust explicit time-scheme for nonlinear mechanics, an alternative to the direct resolution of the semi-implicit coupling in the second-order formulation is developed in the next chapter.

## Chapter 4

# Explicit HHO method for the acoustic wave equation

This chapter is dedicated to the conception, mathematical analysis and numerical evaluation of a splitting procedure to solve the static coupling between the cell and face unknowns arising in the space discretization of the wave equation via the HHO method. Indeed, as observed in the previous chapter, this method leads to a semi-implicit scheme, even if the time-scheme is explicit. The example of the leapfrog scheme was considered in the previous chapter. Using the notation introduced previously, we must solve, at each time step  $t^n$  of the time stepping, the following static coupling between the cell unknowns  $\mathbf{U}_{\mathcal{T}}^n$  and the face unknowns  $\mathbf{U}_{\mathcal{F}}^n$ :

$$\mathcal{A}_{\mathcal{F}\mathcal{F}}\mathbf{U}_{\mathcal{F}}^n = -\mathcal{A}_{\mathcal{F}\mathcal{T}}\mathbf{U}_{\mathcal{T}}^n, \quad (4.1)$$

where  $\mathbf{U}_{\mathcal{T}}^n$  is the data,  $\mathbf{U}_{\mathcal{F}}^n$  the unknown and  $\mathcal{A}$  the complete stiffness matrix defined by (2.31).

This static coupling requires to solve the linear system involving the sparse but non block-diagonal matrix  $\mathcal{A}_{\mathcal{F}\mathcal{F}}$  with either a direct or an iterative solver. The burden may not be excessive in the linear case, but it will become significant in the nonlinear case, such as the realistic mechanical cases that are the target of this Thesis. Indeed, in a nonlinear context, this static coupling requires to solve a nonlinear equation at each time step, using for instance a Newton algorithm. Let us also observe that, when using an explicit time-scheme, the time step is bounded by a critical time step via a CFL condition in order to preserve stability. Thus, the expensive nonlinear solver will be called on a very large number of time steps. Hence, using an explicit time-scheme with the HHO space semi-discretization will be very costly.

In order to reduce the computational cost of solving (4.1), a splitting procedure is presented (Section 4.1). It is introduced in the linear setting, since it allows for a mathematical analysis of the convergence of the splitting procedure. Some numerical experiments are also presented in order to illustrate the impact of the splitting on the quality of the solution and on the critical time step (Section 4.2). These experiments also give optimal values of the splitting parameters regarding the execution time and the error for different settings, as well as a procedure for determining these values. Finally, a wave propagation numerical test case on a contrasted material is simulated (Section 4.3). The results illustrate the efficiency of the HHO method in this setting, compared to lowest-order standard finite elements. A short study of the computational efficiency is also made, but it will be more relevant in the nonlinear case and more thoroughly described in the next chapter.

### 4.1 Study of the operator splitting

This section is dedicated to the conception and the mathematical analysis of a splitting procedure to solve equation (4.1). The splitting procedure depends on the setting for the polynomial order in the

HHO method (mixed- or equal-order). In the mixed-order setting, the procedure relies on the face-face stabilization matrix  $\mathcal{S}_{\mathcal{F}\mathcal{F}}$  which is block-diagonal. In the equal-order setting, this is not the case, but a block-diagonal matrix  $\mathcal{S}_{\mathcal{F}\mathcal{F}}^*$  can be extracted from  $\mathcal{S}_{\mathcal{F}\mathcal{F}}$  to formulate the splitting procedure.

For both cases, a convergence criterion is given along with a proof. This criterion does not depend on the mesh size, but only on the polynomial order and on the shape regularity of the mesh.

#### 4.1.1 Design of the operator splitting

Let us first recall some definitions of Section 3.2.1:  $N$  is the number of discrete time intervals such that  $(t^n)_{n \in \{0:N\}}$  are the evenly spaced discrete time nodes with  $t^0 = 0$  and  $t^N := \mathfrak{T}$ . The leapfrog scheme consists of solving, for all  $n \in \{1:N-1\}$  and all  $\hat{w}_h \in \widehat{\mathcal{U}}_{h,0}^{l,k}$ ,

$$\frac{1}{\Delta t^2}(u_{\mathcal{T}}^{n+1} - 2u_{\mathcal{T}}^n + u_{\mathcal{T}}^{n-1}, w_{\mathcal{T}})_{\Omega} + a_h(\hat{u}_h^n, \hat{w}_h) = (f^n, w_{\mathcal{T}})_{\Omega}, \quad (4.2)$$

where the unknowns are  $u_{\mathcal{T}}^{n+1}$  and  $u_{\mathcal{F}}^n$ , whereas  $u_{\mathcal{T}}^n$  and  $u_{\mathcal{T}}^{n-1}$  are known from prior time steps or the initial conditions (3.25). Equation (4.2) can be rewritten in two steps:

1. For all  $w_{\mathcal{F}} \in \mathcal{U}_{\mathcal{F},0}^k$ , solve  $a_h(\hat{u}_h^n, (0, w_{\mathcal{F}})) = 0$ , to find  $u_{\mathcal{F}}^n$  as a function of  $u_{\mathcal{T}}^n$ ,
2. For all  $w_{\mathcal{T}} \in \mathcal{U}_{\mathcal{T}}^l$ , solve  $\frac{1}{\Delta t^2}(u_{\mathcal{T}}^{n+1} - 2u_{\mathcal{T}}^n + u_{\mathcal{T}}^{n-1}, w_{\mathcal{T}})_{\Omega} + a_h(\hat{u}_h^n, (w_{\mathcal{T}}, 0)) = (f^n, w_{\mathcal{T}})_{\Omega}$ , to find  $u_{\mathcal{T}}^{n+1}$ .

The first step is the functional writing of (4.1) and also rewrites

$$a_h((0, u_{\mathcal{F}}^n), (0, w_{\mathcal{F}})) = -a_h((u_{\mathcal{T}}^n, 0), (0, w_{\mathcal{F}})), \quad \forall w_{\mathcal{F}} \in \mathcal{U}_{\mathcal{F},0}^k. \quad (4.3)$$

**Mixed-order case.** We consider first the mixed-order setting since the definition of the stabilization is simpler and offers a more direct expression of the splitting procedure. Let us remark that the coupling equation (4.3) writes in a more detailed manner

$$s_h((0, u_{\mathcal{F}}^n), (0, w_{\mathcal{F}})) + b_h((0, u_{\mathcal{F}}^n), (0, w_{\mathcal{F}})) = -b_h((u_{\mathcal{T}}^n, 0), (0, w_{\mathcal{F}})) - s_h((u_{\mathcal{T}}^n, 0), (0, w_{\mathcal{F}})), \quad \forall w_{\mathcal{F}} \in \widehat{\mathcal{U}}_{h,0}^{l,k},$$

where  $b_h$  and  $s_h$  are defined in Chapter 2 by (2.27).

The splitting proceeds as follows: for all  $n \in \{1:N\}$ , set  $u_{\mathcal{F}}^{n,0} = u_{\mathcal{F}}^{n-1}$  and iterate on  $m \geq 0$  by finding  $u_{\mathcal{F}}^{n,m+1} \in \mathcal{U}_{\mathcal{F},0}^k$  such that

$$s_h((0, u_{\mathcal{F}}^{n,m+1}), (0, w_{\mathcal{F}})) = -b_h((u_{\mathcal{T}}^n, u_{\mathcal{F}}^{n,m}), (0, w_{\mathcal{F}})) - s_h((u_{\mathcal{T}}^n, 0), (0, w_{\mathcal{F}})), \quad (4.4)$$

for all  $w_{\mathcal{F}} \in \mathcal{U}_{\mathcal{F},0}^k$ . Its algebraic form is

$$\mathcal{S}_{\mathcal{F}\mathcal{F}} \mathbf{U}_{\mathcal{F}}^{n,m+1} = -\mathcal{B}_{\mathcal{F}\mathcal{F}} \mathbf{U}_{\mathcal{F}}^{n,m} - \mathcal{A}_{\mathcal{F}\mathcal{T}} \mathbf{U}_{\mathcal{T}}^n,$$

with  $\mathcal{B}$  and  $\mathcal{S}$  defined by (2.31). This splitting procedure is a fixed-point algorithm using the affine function  $g : X \rightarrow \mathcal{S}_{\mathcal{F}\mathcal{F}}^{-1} \mathcal{B}_{\mathcal{F}\mathcal{F}} X - \mathcal{A}_{\mathcal{F}\mathcal{T}} \mathbf{U}_{\mathcal{T}}^n$  for  $X \in \mathcal{U}_{\mathcal{F},0}^k$ . Its convergence is thus ensured for  $\delta < 1$ , where  $\delta$  is the Lipschitz constant of  $g$ . Since  $g$  is an affine function, this Lipschitz constant is found by computing the spectral radius of  $\mathcal{S}_{\mathcal{F}\mathcal{F}}^{-1} \mathcal{B}_{\mathcal{F}\mathcal{F}}$  and the convergence criterion is

$$\rho(\mathcal{S}_{\mathcal{F}\mathcal{F}}^{-1} \mathcal{B}_{\mathcal{F}\mathcal{F}}) < 1, \quad (4.5)$$

where  $\rho(\mathcal{Q})$  designates the spectral radius of the matrix  $\mathcal{Q}$ . This condition does not depend on the time index  $n$ . The splitting ends when a convergence condition  $\epsilon > 0$  is reached on the increment, i.e.

$$\|u_{\mathcal{F}}^{n,m+1} - u_{\mathcal{F}}^{n,m}\|_{\mathcal{F}} < \epsilon \|u_{\mathcal{F}}^{n,0}\|_{\mathcal{F}}, \quad (4.6)$$

and we set  $u_{\mathcal{F}}^n := u_{\mathcal{F}}^{n,m+1}$ . This splitting procedure is computationally effective since, as mentioned above, the face-face stabilization submatrix  $\mathcal{S}_{\mathcal{F}\mathcal{F}}$  is block-diagonal, the size of each block being  $\binom{k+d-1}{d-1}$ , i.e.,  $(k+1)$  for  $d=2$  and  $\frac{1}{2}(k+1)(k+2)$  for  $d=3$ .

**Equal-order case** The equal-order setting does not offer the same simplicity since the stabilization couples together the degrees of freedom of all the faces of a cell. One can, however, draw on the mixed-order setting and split the stabilization form into the mixed-order stabilization form, which leads to a block-diagonal matrix, and the remainder. Specifically, let  $\zeta_T$  be the local form such that for all  $T \in \mathcal{T}_h$ , and all  $v_{\partial T}, w_{\partial T} \in \mathcal{U}_{\partial T}^k$ ,

$$\begin{aligned} \zeta_T((0, v_{\partial T}), (0, w_{\partial T})) &= \gamma \mu_T^2 \sum_{F \in \mathcal{F}_T} \eta_{TF}^{-1} \left\{ ((I - \Pi_T^k) R_T^{k+1}(0, v_{\partial T})|_F, w_F)_F \right. \\ &\quad + (v_F, (I - \Pi_T^k) R_T^{k+1}(0, w_{\partial T})|_F)_F \\ &\quad \left. + (\Pi_F^k (I - \Pi_T^k) R_T^{k+1}(0, v_{\partial T})|_F, \Pi_F^k (I - \Pi_T^k) R_T^{k+1}(0, w_{\partial T})|_F)_F \right\}, \end{aligned} \quad (4.7)$$

and let  $s_T^*$  be the local form defined by

$$s_T^*((0, v_{\partial T}), (0, w_{\partial T})) := \gamma \mu_T^2 \sum_{F \in \mathcal{F}_T} \eta_{TF}^{-1} (v_F, w_F)_F. \quad (4.8)$$

Then the equal-order stabilization form writes

$$s_T((0, v_{\partial T}), (0, w_{\partial T})) = s_T^*((0, v_{\partial T}), (0, w_{\partial T})) + \zeta_T((0, v_{\partial T}), (0, w_{\partial T})). \quad (4.9)$$

Let us introduce the global forms  $s_h^*((0, v_{\mathcal{F}}), (0, w_{\mathcal{F}})) := \sum_{T \in \mathcal{T}_h} s_T^*((0, v_{\partial T}), (0, w_{\partial T}))$  and  $\zeta_h((0, v_{\mathcal{F}}), (0, w_{\mathcal{F}})) := \sum_{T \in \mathcal{T}_h} \zeta_T((0, v_{\partial T}), (0, w_{\partial T}))$ , so that  $s_h = s_h^* + \zeta_h$ . This leads to the following iterative procedure, with the same initial condition as for the mixed-order setting: For all  $m \geq 0$ , find  $u_{\mathcal{F}}^{n, m+1} \in \mathcal{U}_{\mathcal{F}, 0}^k$  such that

$$\begin{aligned} s_h^*((0, u_{\mathcal{F}}^{n, m+1}), (0, w_{\mathcal{F}})) &= -b_h((u_{\mathcal{T}}^n, u_{\mathcal{F}}^{n, m}), (0, w_{\mathcal{F}})) \\ &\quad - \zeta_h((0, u_{\mathcal{F}}^{n, m}), (0, w_{\mathcal{F}})) - s_h((u_{\mathcal{T}}^n, 0), (0, w_{\mathcal{F}})), \end{aligned} \quad (4.10)$$

for all  $w_{\mathcal{F}} \in \mathcal{U}_{\mathcal{F}, 0}^k$ . At the algebraic level, we define two matrices  $\mathcal{S}_{\mathcal{F}\mathcal{F}}^*$  and  $\mathcal{Z}_{\mathcal{F}\mathcal{F}}$  such that  $\mathcal{S}_{\mathcal{F}\mathcal{F}} = \mathcal{S}_{\mathcal{F}\mathcal{F}}^* + \mathcal{Z}_{\mathcal{F}\mathcal{F}}$ ,  $\mathcal{Z}_{\mathcal{F}\mathcal{F}}$  corresponds to the bilinear form  $\zeta_h(\cdot, \cdot)$  and  $\mathcal{S}_{\mathcal{F}\mathcal{F}}^*$  to  $s_h^*(\cdot, \cdot)$ . Then the splitting procedure translates into the following iterative algorithm: For all  $m \geq 0$ , find  $U_{\mathcal{F}}^{n, m+1} \in \mathbb{R}^{N_{\mathcal{F}}}$  such that

$$\mathcal{S}_{\mathcal{F}\mathcal{F}}^* U_{\mathcal{F}}^{n, m+1} = -\mathcal{B}_{\mathcal{F}\mathcal{F}} U_{\mathcal{F}}^{n, m} - \mathcal{Z}_{\mathcal{F}\mathcal{F}} U_{\mathcal{F}}^{n, m} - \mathcal{A}_{\mathcal{F}\mathcal{T}} U_{\mathcal{T}}^n. \quad (4.11)$$

As for the mixed-order setting, the convergence condition is that  $\delta < 1$ , where  $\delta$  is the Lipschitz constant of the vector-valued function  $g : X \rightarrow (\mathcal{S}_{\mathcal{F}\mathcal{F}}^*)^{-1} (\mathcal{Z}_{\mathcal{F}\mathcal{F}} + \mathcal{B}_{\mathcal{F}\mathcal{F}}) X$ , which translates in terms of a spectral radius as

$$\rho((\mathcal{S}_{\mathcal{F}\mathcal{F}}^*)^{-1} (\mathcal{B}_{\mathcal{F}\mathcal{F}} + \mathcal{Z}_{\mathcal{F}\mathcal{F}})) < 1. \quad (4.12)$$

**Remark 4.1.1** (Rewriting as a Neumann series). The splitting procedure can be interpreted as a Neumann series to invert the matrix  $\mathbb{A}_{\mathcal{F}\mathcal{F}}$ . In the mixed-order setting, the algebraic form writes as follows, with  $N$  the number of splitting iterations:

$$\frac{1}{\Delta t^2} \mathcal{M} U_{\mathcal{T}}^{n+1} + \left[ \mathcal{A}_{\mathcal{T}\mathcal{T}} - \mathcal{A}_{\mathcal{T}\mathcal{F}} \left( \sum_{n=0}^N (-1)^n (\mathcal{S}_{\mathcal{F}\mathcal{F}}^{-1} \mathcal{B}_{\mathcal{F}\mathcal{F}})^n \mathcal{S}_{\mathcal{F}\mathcal{F}}^{-1} \right) \mathcal{A}_{\mathcal{F}\mathcal{T}} \right] U_{\mathcal{T}}^n = F_{\mathcal{T}}^n + \frac{1}{\Delta t^2} \mathcal{M} (2U_{\mathcal{T}}^n - U_{\mathcal{T}}^{n-1}). \quad (4.13)$$

In the equal-order setting, the above formulation must be adapted by replacing  $\mathcal{S}_{\mathcal{F}\mathcal{F}}^{-1}$  by  $(\mathcal{S}_{\mathcal{F}\mathcal{F}}^*)^{-1}$  and  $\mathcal{B}_{\mathcal{F}\mathcal{F}}$  by  $\mathcal{B}_{\mathcal{F}\mathcal{F}} + \mathcal{Z}_{\mathcal{F}\mathcal{F}}$ .

## 4.1.2 Convergence

In this section, we show that the convergence criteria (4.5) and (4.12) can be achieved by choosing  $\gamma$ , the coefficient scaling the stabilization, large enough. Moreover, we derive explicit lower bounds on the coefficient  $\gamma$  that are uniform in the mesh size.

Invoking a discrete trace inequality (see, e.g., [97, Lem. 12.8] on simplicial meshes and [81, Lem. 1.46] on more general meshes), we infer that there exists a trace constant  $C_{\text{tr}}$  independent of  $h$  and which only depends on the polynomial order  $k$  and the mesh regularity parameter  $\rho$  such that

$$\max_{T \in \mathcal{T}_h} \max_{F \in \mathcal{F}_T} \max_{v \in \mathbb{P}_d^k(T)} \frac{|T|^{1/2} \|v\|_F}{|F|^{1/2} \|v\|_T} \leq C_{\text{tr}}. \quad (4.14)$$

Moreover, we introduce a nondimensional geometric constant  $\rho_{\sharp}$  such that

$$\max_{T \in \mathcal{T}_h} \left( \sum_{F \in \mathcal{F}_T} \frac{\eta_{TF} |F|}{|T|} \right)^{\frac{1}{2}} \leq \rho_{\sharp} \quad (4.15)$$

recalling that  $\eta_{TF}$  is the length scale used in the stabilization. The value of  $\rho_{\sharp}$  can be chosen to be independent of  $h$  owing to the regularity of the mesh sequence.

### Mixed-order case

**Lemma 4.1.2** (Convergence in the mixed-order setting). *A sufficient condition ensuring the convergence of the iterative algorithm (4.4) in the mixed-order setting is*

$$\gamma > (C_{\text{tr}} \rho_{\sharp})^2. \quad (4.16)$$

*Proof.* The main idea of the proof is the following equivalence regarding the convergence condition (4.5):

$$\rho(\mathcal{S}_{\mathcal{F}\mathcal{F}}^{-1} \mathcal{B}_{\mathcal{F}\mathcal{F}}) < 1 \iff s_h((0, w_{\mathcal{F}}), (0, w_{\mathcal{F}})) > b_h((0, w_{\mathcal{F}}), (0, w_{\mathcal{F}})), \quad \forall w_{\mathcal{F}} \in \mathcal{U}_{\mathcal{F},0}^k. \quad (4.17)$$

Let  $w_{\mathcal{F}} \in \mathcal{U}_{\mathcal{F},0}^k$  and let  $T \in \mathcal{T}_h$ . Recalling that  $\mu_T$  denotes the constant value taken by  $\mu$  in the cell  $T \in \mathcal{T}_h$ , we have

$$\begin{aligned} b_T((0, w_{\partial T}), (0, w_{\partial T})) &= \mu_T^2 \|\mathbf{G}_T^k(0, w_{\partial T})\|_T^2 = \mu_T^2 \sum_{F \in \mathcal{F}_T} (\mathbf{G}_T^k(0, w_{\partial T}) \cdot \mathbf{n}_{TF}, w_F) \\ &\leq \sum_{F \in \mathcal{F}_T} \mu_T \eta_{TF}^{\frac{1}{2}} \|\mathbf{G}_T^k(0, w_{\partial T}) \cdot \mathbf{n}_{TF}\|_F \mu_T \eta_{TF}^{-\frac{1}{2}} \|w_F\|_F \\ &\leq \left( \sum_{F \in \mathcal{F}_T} \mu_T^2 \eta_{TF} \|\mathbf{G}_T^k(0, w_{\partial T}) \cdot \mathbf{n}_{TF}\|_F^2 \right)^{\frac{1}{2}} \left( \sum_{F \in \mathcal{F}_T} \mu_T^2 \eta_{TF}^{-1} \|w_F\|_F^2 \right)^{\frac{1}{2}}, \end{aligned} \quad (4.18)$$

where we used the definition of the gradient reconstruction operator and the Cauchy–Schwarz inequality. Since  $\mathbf{G}_T^k(0, w_{\partial T}) \cdot \mathbf{n}_{TF} \in \mathbb{P}_d^k(T)$  because  $\mathbf{n}_{TF}$  is a constant vector for all  $F \in \mathcal{F}_h$ , we can use the discrete inverse trace inequality (4.14) on each face  $F \in \mathcal{F}_T$  to infer that

$$\eta_{TF}^{\frac{1}{2}} \|\mathbf{G}_T^k(0, w_{\partial T}) \cdot \mathbf{n}_{TF}\|_F \leq C_{\text{tr}} \|\mathbf{G}_T^k(0, w_{\partial T})\|_T \left( \frac{\eta_{TF} |F|}{|T|} \right)^{\frac{1}{2}}, \quad \forall F \in \mathcal{F}_h$$

Recognizing the definition of the mixed-order stabilization form without  $\gamma$  in the rightmost sum of (4.18), we obtain

$$b_T((0, w_{\partial T}), (0, w_{\partial T})) \leq C_{\text{tr}} \mu_T \|\mathbf{G}_T^k(0, w_{\partial T})\|_T \left( \sum_{F \in \mathcal{F}_T} \frac{\eta_{TF} |F|}{|T|} \right)^{\frac{1}{2}} \left( \frac{1}{\gamma} s_T((0, w_{\partial T}), (0, w_{\partial T})) \right)^{\frac{1}{2}}.$$

Using the definition of  $\rho_{\sharp}$  and that  $b_T((0, w_{\partial T}), (0, w_{\partial T})) = \mu_T^2 \|\mathbf{G}_T^k(0, w_{\partial T})\|_T^2$  gives

$$b_T((0, w_{\partial T}), (0, w_{\partial T})) \leq \frac{C_{\text{tr}}^2 \rho_{\sharp}^2}{\gamma} s_T((0, w_{\partial T}), (0, w_{\partial T})).$$

Summing over all the mesh cells, we obtain

$$b_h((0, w_{\mathcal{F}}), (0, w_{\mathcal{F}})) \leq \frac{C_{\text{tr}}^2 \rho_{\sharp}^2}{\gamma} s_h((0, w_{\mathcal{F}}), (0, w_{\mathcal{F}})).$$

Thus, the splitting converges if  $\frac{C_{\text{tr}}^2 \rho_{\sharp}^2}{\gamma} < 1$ , which proves the claim.  $\square$

**Equal-order case** The equal-order setting does not provide the same simplicity and requires a hypothesis on the spectrum of the matrix  $(\mathcal{S}_{\mathcal{F}\mathcal{F}}^*)^{-1} \mathcal{S}_{\mathcal{F}\mathcal{F}}$ . Let us define

$$\alpha := \rho((\mathcal{S}_{\mathcal{F}\mathcal{F}}^*)^{-1} \mathcal{Z}_{\mathcal{F}\mathcal{F}}) = \rho((\mathcal{S}_{\mathcal{F}\mathcal{F}}^*)^{-1} \mathcal{S}_{\mathcal{F}\mathcal{F}}) - 1, \quad (4.19)$$

Notice that the value of  $\alpha$  does not depend on  $\gamma$  (since all the involved matrices are proportional to  $\gamma$ ) and that the matrices considered in (4.19) are symmetric.

**Lemma 4.1.3** (Convergence in the equal-order setting). *Under the assumption  $\alpha < 1$ , a sufficient condition for the convergence of the iterative algorithm (5.18) in the equal-order setting is*

$$\gamma > \frac{(C_{\text{tr}} \rho_{\sharp})^2}{1 - \alpha}. \quad (4.20)$$

*Proof.* The following equivalence holds regarding the convergence condition (4.12):

$$\begin{aligned} \rho((\mathcal{S}_{\mathcal{F}\mathcal{F}}^*)^{-1} (\mathcal{B}_{\mathcal{F}\mathcal{F}} + \mathcal{Z}_{\mathcal{F}\mathcal{F}})) &< 1 \\ \iff s_h^*((0, w_{\mathcal{F}}), (0, w_{\mathcal{F}})) &> b_h((0, w_{\mathcal{F}}), (0, w_{\mathcal{F}})) + \zeta_h((0, w_{\mathcal{F}}), (0, w_{\mathcal{F}})), \quad \forall w_{\mathcal{F}} \in \mathcal{U}_{\mathcal{F},0}^k. \end{aligned} \quad (4.21)$$

Since  $\zeta_h$  and  $s_h^*$  are symmetric and  $s_h^*$  is nonnegative, the definition of  $\alpha$  means that

$$\zeta_h((0, w_{\mathcal{F}}), (0, w_{\mathcal{F}})) \leq \alpha s_h^*((0, w_{\mathcal{F}}), (0, w_{\mathcal{F}})), \quad \forall w_{\mathcal{F}} \in \mathcal{U}_{\mathcal{F},0}^k. \quad (4.22)$$

Considering this, a sufficient condition to obtain the bound announced in (4.21) is

$$(1 - \alpha) s_h^*((0, w_{\mathcal{F}}), (0, w_{\mathcal{F}})) > b_h((0, w_{\mathcal{F}}), (0, w_{\mathcal{F}})), \quad \forall w_{\mathcal{F}} \in \mathcal{U}_{\mathcal{F},0}^k. \quad (4.23)$$

Indeed, if (4.23) holds, we infer that

$$b_h((0, w_{\mathcal{F}}), (0, w_{\mathcal{F}})) < s_h^*((0, w_{\mathcal{F}}), (0, w_{\mathcal{F}})) - \alpha s_h^*((0, w_{\mathcal{F}}), (0, w_{\mathcal{F}})) \leq s_h^*((0, w_{\mathcal{F}}), (0, w_{\mathcal{F}})) - \zeta_h((0, w_{\mathcal{F}}), (0, w_{\mathcal{F}})).$$

Lemma 4.1.2 established that, if  $\gamma > (C_{\text{tr}} \rho_{\sharp})^2$ , then,

$$s_h^*((0, w_{\mathcal{F}}), (0, w_{\mathcal{F}})) > b_h((0, w_{\mathcal{F}}), (0, w_{\mathcal{F}})), \quad \forall w_{\mathcal{F}} \in \mathcal{U}_{\mathcal{F},0}^k. \quad (4.24)$$

Therefore, taking  $\gamma > \frac{(C_{\text{tr}} \rho_{\sharp})^2}{1 - \alpha}$  ensures the condition (4.23). This concludes the proof.  $\square$

**Remark 4.1.4** (Assumption  $\alpha < 1$ ). This assumption is verified numerically in Section 4.1.3, see in particular Table 4.2. Although a proof of this assumption is not available, our numerical experiments indicate that it is reasonable to expect that it holds for polynomial degrees from 0 to 4 and relatively simple mesh shapes.

**Remark 4.1.5** (Local vs. global spectral radius). We notice that the minimal value of  $\gamma$  given in Lemmas 4.1.2 and 4.1.3 is derived by reasoning locally on a single mesh cell  $T \in \mathcal{T}_h$ . Slightly sharper values can be derived by reasoning globally on the mesh and taking into account the homogeneous Dirichlet boundary conditions. This point is further quantified in the following.



### 4.1.3 Numerical study of stability parameter

In this section, we evaluate numerically the influence of the stability parameter  $\gamma$  on the splitting procedure for various polynomial orders. Both equal-order and mixed-order settings are tested, with polynomial orders  $k \in \{0, 1, 2, 3, 4\}$ , as well as various mesh types: Cartesian in two and three dimensions, unstructured quadrangles, structured and unstructured triangles and unstructured tetrahedra. The polynomial order for the HHO method is indicated via the pair of integers  $(l, k)$  with  $l$  for the cell unknowns and  $k$  for the face unknowns.

**2D case.** We first verify that the minimum value of the parameter  $\gamma$  ensuring the convergence of the splitting procedure is bounded, for all the polynomial orders, uniformly in the mesh size. To this purpose, we first consider a series of refined right-triangular meshes of  $\Omega := (0, 1)^2$  and compute the value of  $\gamma$  on each mesh via the spectral radius of  $\mathcal{S}_{\mathcal{FF}}^{-1}\mathcal{B}_{\mathcal{FF}}$  in the mixed-order setting and the spectral radius of  $\mathcal{S}_{\mathcal{FF}}^{-1}(\mathcal{B}_{\mathcal{FF}} + \mathcal{Z}_{\mathcal{FF}})$  in the equal-order setting, with homogeneous Dirichlet boundary conditions. We compare the resulting values to the value  $\gamma^*$ , which is the value of  $\gamma$  obtained in a single cell without boundary conditions. Figure 4.1 illustrates, as expected, that  $\gamma^*$  gives a reliable upper bound on  $\gamma$ . This

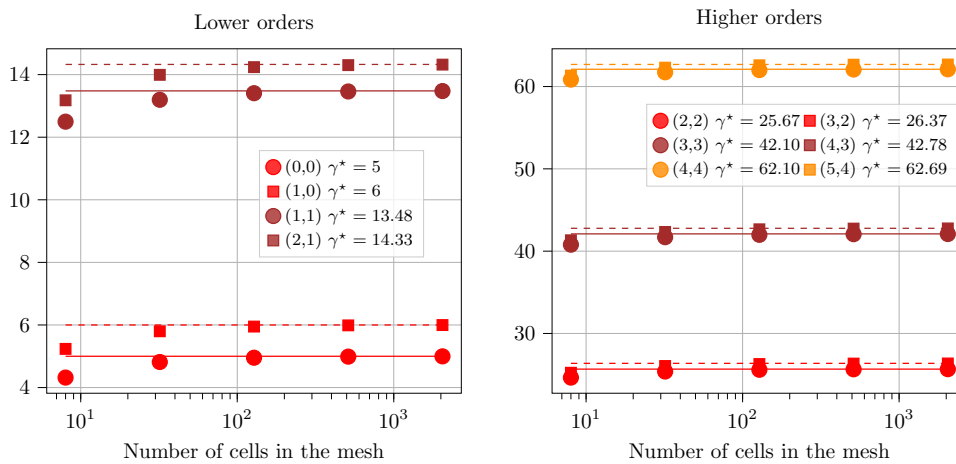


Figure 4.1: Spectral radius as a function of the mesh size in right-triangular cells with homogeneous Dirichlet conditions in 2D, compared to the reference value on a single cell without boundary conditions (horizontal lines).

upper bound turns out to be quite sharp even on moderately refined meshes. Indeed, for all polynomial orders, on meshes with 100 cells or more, the value of  $\gamma$  reported in Figure 4.1 coincides with the value of  $\gamma^*$ , materialized by the horizontal lines. Therefore,  $\gamma^*$  is a very good minimal value for  $\gamma$  to be used in practice.

Table 4.1 reports the value of  $\gamma^*$  for a square cell, a right-isocles triangular cell as well as estimates for general quadrangular and triangular cells belonging to shape-regular sequences of unstructured meshes. The triangular meshes are generated using `gmsh`, and the quadrangular meshes are also created using `gmsh`, which creates the quadrangular meshes from the triangular meshes by merging pairs of adjacent triangles. For these two meshes, the reported value of  $\gamma^*$  is the largest observed value, rounded to the above integer on all the mesh cells and for all the meshes in the sequence. One notices that, in all cases, the value of  $\gamma^*$  increases with the polynomial order. Moreover, for a given polynomial order, the smallest value of  $\gamma^*$  is obtained on square cells, and the largest value on unstructured triangular meshes.

As a second verification, we check the condition  $\alpha < 1$  in the equal-order setting on all the previous cases. The value of  $\alpha$  appears not to depend on the mesh size when homogeneous Dirichlet boundary conditions are enforced. As before, when unstructured meshes are used, the largest value observed on

Order (cell, face)	(0,0)	(1,1)	(2,2)	(3,3)	(4,4)	(1,0)	(2,1)	(3,2)	(4,3)	(5,4)
Squares	1	5	11	19	29	2	6	12	20	30
Right triangles	5	13.48	25.67	42.10	62.10	6	14.33	26.37	42.78	62.69
Unstructured quadrangles	3	9	19	32	48	4	9	19	32	48
Unstructured triangles	6	15	28	45	65	7	15	28	45	65

Table 4.1:  $\gamma^*$  computed via the spectral radius of  $\mathcal{S}_{\mathcal{F}\mathcal{F}}^{-1}\mathcal{B}_{\mathcal{F}\mathcal{F}}$  or  $(\mathcal{S}_{\mathcal{F}\mathcal{F}}^*)^{-1}(\mathcal{B}_{\mathcal{F}\mathcal{F}} + \mathcal{Z}_{\mathcal{F}\mathcal{F}})$  on a single cell (first two lines) and on a shape-regular mesh sequence (last two lines),  $k \in \{0, 1, 2, 3, 4\}$ , equal- and mixed-order settings.

the mesh sequence is reported. As seen in Table 4.2, the value of  $\alpha$  is in all cases smaller than 1, thereby confirming the assumption made in Lemma 4.1.3. The lowest-order setting even yields  $\alpha$  close to zero. The value of  $\alpha$  increases with the polynomial order before it stabilizes for higher orders. For unstructured mesh sequences, the finer meshes are more uniform (i.e. the different cells of the mesh have a more similar shape) leading to a smaller value of  $\alpha$ . Therefore, the values reported in Table 4.2 are pessimistic on finer meshes.

Order (cell, face)	(0,0)	(1,1)	(2,2)	(3,3)	(4,4)
Squares	0.04	0.11	0.19	0.26	0.23
Right triangles	0.05	0.28	0.38	0.36	0.39
Unstructured quadrangles	0.13	0.21	0.44	0.52	0.61
Unstructured triangles	0.05	0.17	0.25	0.24	0.32

Table 4.2:  $\alpha$ , the spectral radius  $(\mathcal{S}_{\mathcal{F}\mathcal{F}}^*)^{-1}\mathcal{Z}_{\mathcal{F}\mathcal{F}}$ ,  $k \in \{0, 1, 2, 3, 4\}$ , equal- and mixed-order settings.

A classical result on the leapfrog scheme (3.24) gives a stability condition on the time step depending on the spectral radius of a certain matrix (see, e.g., [117]). Here, the matrix to be considered is

$$\mathcal{D}(\gamma) := \mathcal{M}^{-1}(\mathcal{A}_{\mathcal{T}\mathcal{T}}(\gamma) - \mathcal{A}_{\mathcal{T}\mathcal{F}}(\gamma)\mathcal{A}_{\mathcal{F}\mathcal{F}}(\gamma)^{-1}\mathcal{A}_{\mathcal{F}\mathcal{T}}(\gamma)), \quad (4.25)$$

where  $\mathcal{A}(\gamma)$  is the stiffness matrix defined in (2.31) with the dependence on  $\gamma$  made explicit. The stability condition on  $\Delta t$  then reads

$$\Delta t(\gamma) \leq \Delta t^{\text{opt}}(\gamma) := \frac{2}{\sqrt{\rho(\mathcal{D}(\gamma))}}. \quad (4.26)$$

This allows us to compare stability conditions for different values of  $\gamma$ . In Table 4.3, we compare values of  $\Delta t^{\text{opt}}$  for  $\gamma = 1$ , i.e. without splitting, to those for  $\gamma = \gamma^*$  from Table 4.1. The results show that using the splitting procedure leads to a reduction of the time step by at most a factor of two on squares and by at most three on right-isocles triangles. This reduction of the time step illustrates the fact that the modification of the value of  $\gamma$  impacts the CFL condition (4.26) by changing the spectrum of  $\mathcal{D}$ . The fact that the time step is impacted negatively (i.e. reduced) when  $\gamma$  increases is further illustrated in Section 4.2.3. Moreover, increasing the polynomial degree generally alleviates the tightening of the stability condition except on triangular meshes. The value of the CFL does not change with the mesh size  $h$ , so that the values given in Table 4.3 should be independent on  $h$ .

**3D case.** The same studies can be performed on 3D meshes. Only Cartesian hexahedral and unstructured tetrahedral meshes with polynomial orders up to  $k \leq 2$  are presented. Table 4.4 summarizes the results. Comparing square and hexahedral meshes, going to 3D does not deteriorate too much the stability condition. On tetrahedral meshes, however,  $\gamma^*$  grows faster with the polynomial degree, and the splitting procedure impacts more strongly the stability condition, by at most a factor of five.

Order (cell, face)	(0,0)	(1,1)	(2,2)	(3,3)	(4,4)	(1,0)	(2,1)	(3,2)	(4,3)	(5,4)
Squares	1.00	0.75	0.63	0.68	0.60	0.66	0.54	0.52	0.54	0.52
Right triangles	0.37	0.58	0.65	0.54	0.70	0.40	0.36	0.41	0.53	0.63
Unstructured quadrangles	0.60	0.43	0.32	0.37	0.43	0.45	0.41	0.39	0.47	0.51
Unstructured triangles	0.38	0.30	0.42	0.47	0.56	0.25	0.32	0.39	0.44	0.51

Table 4.3:  $\Delta t^{\text{opt}}(\gamma^*)/\Delta t^{\text{opt}}(1)$ , showing the tightening of the stability condition induced by  $\gamma^*$  on 2D meshes,  $h = 0.1$ ,  $k \in \{0, 1, 2, 3, 4\}$ , equal- and mixed-order settings.

Order (cell, face)	(0,0)	(1,1)	(2,2)	(1,0)	(2,1)	(3,2)
$\gamma^*$ hexahedra	1.83	8.5	17	2.83	9.5	18
$\Delta t^{\text{opt}}(\gamma^*)/\Delta t^{\text{opt}}(1)$ hexahedra	0.57	0.58	0.58	0.5	0.42	0.48
$\gamma^*$ tetrahedra	33	60	90	34	60	90.5
$\Delta t^{\text{opt}}(\gamma^*)/\Delta t^{\text{opt}}(1)$ tetrahedra	0.19	0.17	0.51	0.16	0.26	0.57

Table 4.4: Value of  $\gamma^*$  and tightening of the stability condition induced by  $\gamma^*$  on 3D meshes (Cartesian hexahedra and unstructured tetrahedra),  $h = 0.55$ ,  $k \in \{0, 1, 2\}$ , equal-order setting.

## 4.2 Parametric study of the splitting procedure

This section is dedicated to the study of the impact of the splitting procedure on the quality of the solution and the execution time, depending on the value of its different parameters. The goal is to identify values or range of values for the scaling parameter, the number of iterations and the convergence criterion that optimize the execution time, without deteriorating the solution compared to the semi-implicit setting.

Firstly, a case with a fixed number of iterations is considered. Then a fixed convergence criterion is imposed, as defined in equation (4.6). The number of iterations required at each time step is reported below. Secondly, computations with various values of  $\gamma$  are made to measure the impact on the quality of the solution, the number of iterations needed for the splitting procedure to satisfy the convergence criterion, and the CFL condition. This gives us a way to optimally choose  $\gamma$  in terms of error for a given execution time. Finally, a focus is made on the case  $\gamma \rightarrow \infty$ , and the impact of a very large  $\gamma$  on the quality of the solution and the CFL condition is explored.

### 4.2.1 Impact of the number of iterations

In this paragraph and the next one, we consider a manufactured solution in  $\Omega := (0, 1)^2$  with a non-polynomial behavior in space and a quadratic in time behavior. Specifically, we set  $\mu := 1$  in (3.1) and

$$u(x, y, t) := t^2 \sin(\pi x) \sin(\pi y), \quad (4.27)$$

leading to homogeneous Dirichlet conditions and zero initial conditions. The source term is  $f(x, y, t) := 2(\pi^2 t^2 + 1) \sin(\pi x) \sin(\pi y)$ . Here, there is no time discretization error. Indeed, since the time integration scheme is of order 2, it integrates  $t \mapsto t^2$  exactly.

**Error for a fixed number of iterations.** Figure 4.2 illustrates the convergence of the  $L^2$ -error at the final time  $t = 1.0$  with respect to the number of iterations in the splitting procedure on a series of spatially refined square meshes. We focus on the equal-order setting (the results for the mixed-order setting show the same behavior and are not displayed.) The time step is chosen so that the ratio  $\frac{h}{\Delta t}$  remains constant and smaller than the stability condition (4.26). The reference curve is the error without splitting. The value of  $\gamma$  is determined from Table 4.1 by setting  $\gamma := 1.5\gamma^*$ .

For a given number of iterations in the splitting procedure, the error ends up stagnating below a certain value for  $h$ . The stagnation value and the mesh size for which stagnation starts both decrease as the number of iterations increases. This value of  $h$  corresponds to the splitting error being larger than the discretization error. When the splitting error is sufficiently small, the space discretization error prevails. Moreover, the higher the polynomial order, the more iterations needed to converge. Figure 4.2 also shows that the convergence order of the  $L^2$ -error depending on the mesh size  $h$  is, as expected,  $h^{k+2}$ .

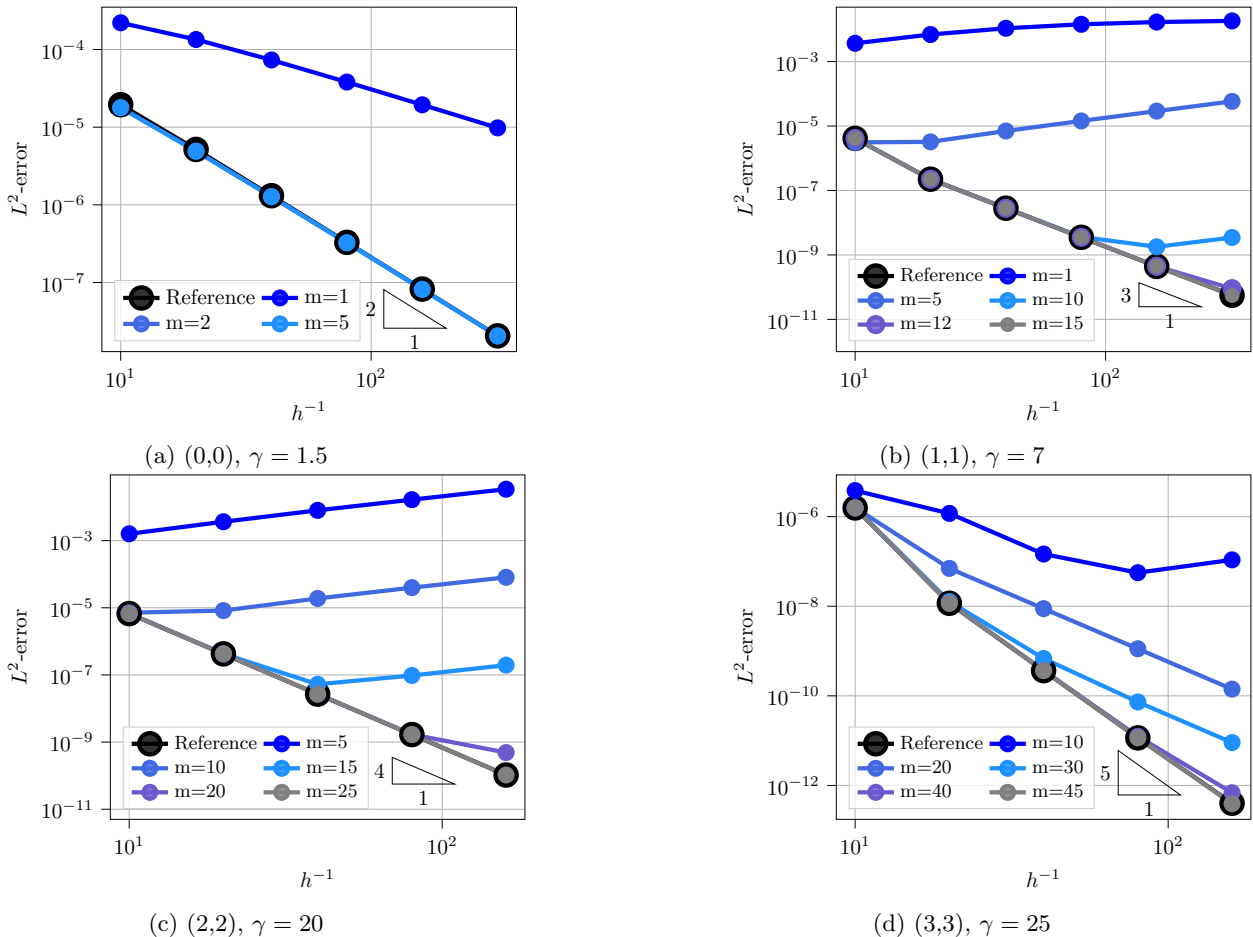


Figure 4.2: Linear wave problem - Convergence of the  $L^2$ -error at the final time  $t = 1.0$  as a function of the mesh size for varying numbers of splitting iterations  $m$ ,  $k \in \{0, 1, 2, 3\}$ , equal-order setting. The reference corresponds to the semi-implicit scheme ( $m \rightarrow \infty$ ).

**Number of iterations depending on the convergence criterion.** Fixing the number of iterations is not the standard way to ensure the convergence of the splitting algorithm. We rather choose a convergence criterion  $\epsilon$  and stop the splitting procedure as soon as the difference between two iterates is smaller than  $\epsilon$ ; see (4.6). Figure 4.3 illustrates the number of iterations needed to reach the convergence criterion (4.6) for different values of  $\epsilon$ . Figure 4.4 displays the associated  $L^2$ - and  $H^1$ -errors. The experiments displayed in both figures are run with the equal-order setting, with polynomial order  $k \in \{0, 1, 2, 3\}$  and on a Cartesian mesh with  $h = 0.05$ . The solution is still given by equation (4.27). The number of iterations increases as  $\epsilon$  decreases, which is expected. It also increases as the simulation runs. This latter phenomenon is more visible for the higher orders ( $k = 2$  and  $k = 3$ ). It can also be noticed that the convergence criterion must be small enough to ensure proper convergence of the splitting procedure. Indeed, for the largest

values of  $\epsilon$ , the errors shown on Figure 4.4 are much larger than for the smallest values of  $\epsilon$  (this is the expected behavior). These results illustrate the fact that  $\epsilon$  controls the convergence of the splitting error. Under a certain value of  $\epsilon$ , the splitting error is (much) smaller than the time- and space-discretization errors and does not contribute significantly to the total error. Under this value, diminishing  $\epsilon$  does not increase the quality of the solution. Taking  $\epsilon$  close to this threshold ( $10^{-4}, 10^{-8}, 10^{-8}, 10^{-8}$  respectively for  $k \in \{0, 1, 2, 3\}$ ) ensures the best trade-off between accuracy and computational efficiency (the fewer iterations, the less expensive the computation). This optimal value of  $\epsilon$  depends on the time- and space-discretization errors, which respectively depend on the time step  $\Delta t$  and the mesh size  $h$ . Since the time step  $\Delta t$  is constrained by the mesh size  $h$  via the CFL condition (4.26), the optimal value of  $\epsilon$  varies for all values of  $h$ . More precisely, under the CFL condition (4.26) the smaller  $h$ , the smaller the time- and space-discretization errors, and the smaller  $\epsilon$  needs to be for the splitting error to be negligible.

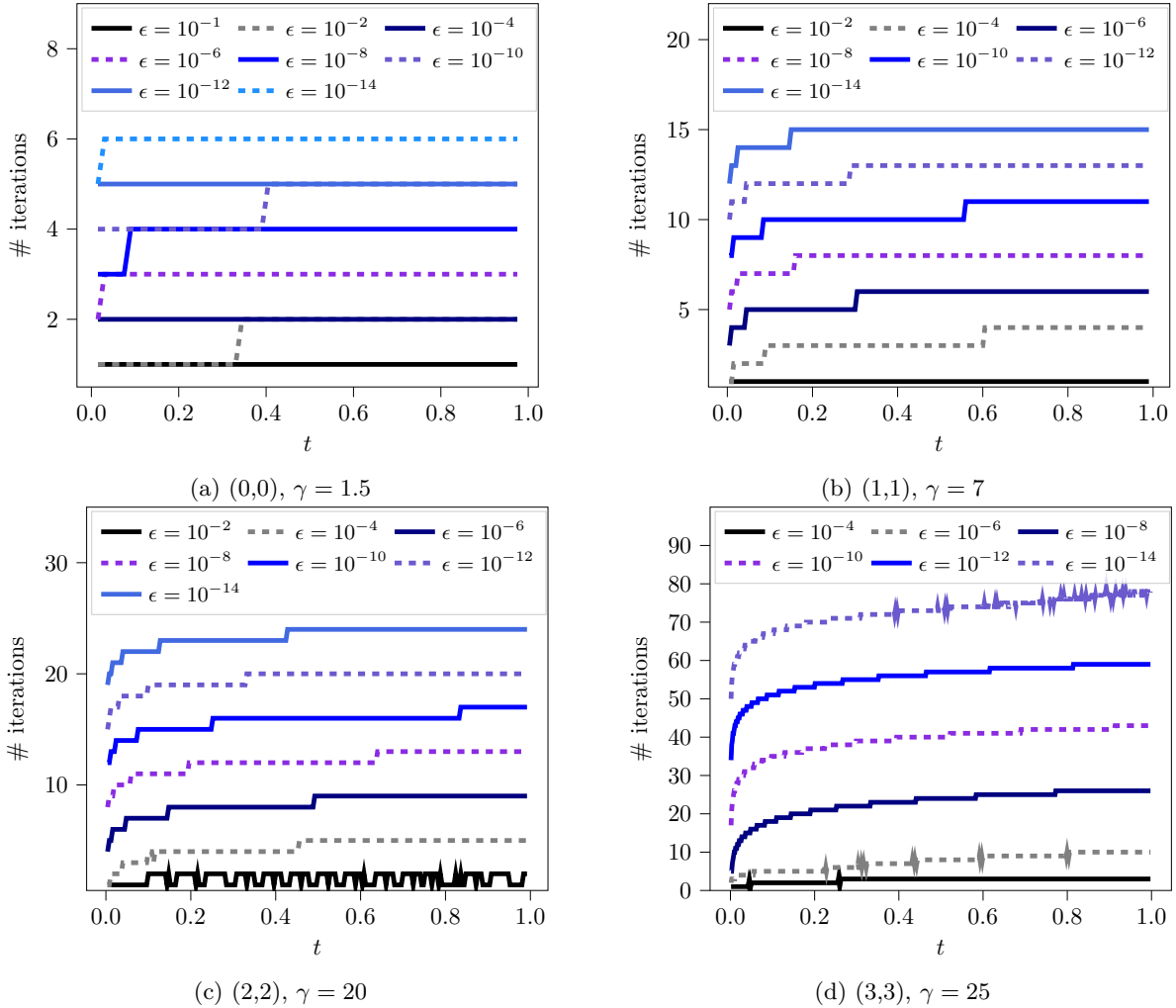
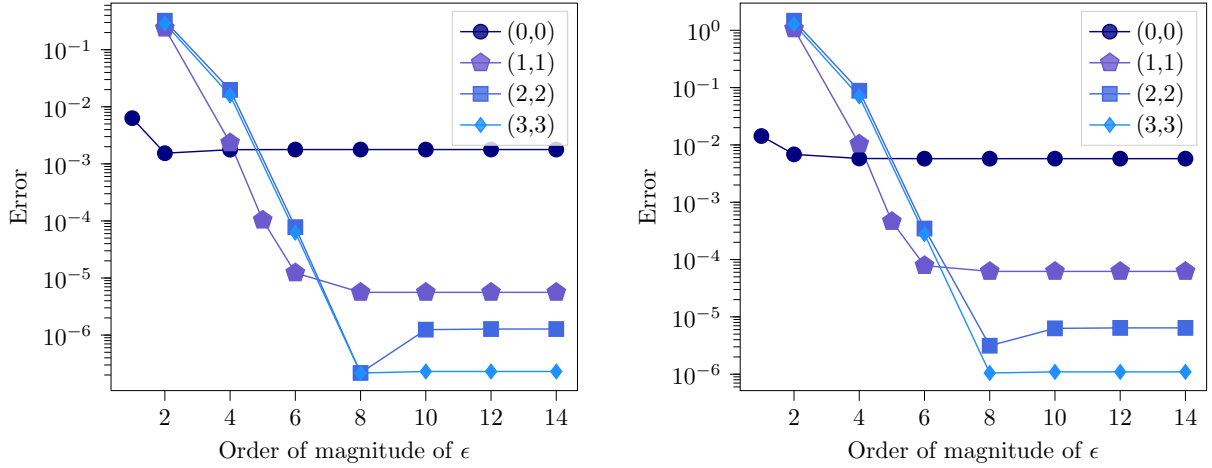


Figure 4.3: Linear wave problem, quadratic in time - Number of iterations over time ( $t \in [0, 1]$ ) in order to satisfy the convergence criterion (4.6) in the splitting procedure,  $k \in \{0, 1, 2, 3\}$ , equal-order setting.

It can be interesting to run the same experiments, but for a different solution. We now consider another manufactured solution that is non-polynomial in space and time, given by

$$u(x, y, t) := \frac{1}{\sqrt{2\pi}} \sin(\sqrt{2}\pi t) \sin(\pi x) \sin(\pi y), \quad (4.28)$$



(a)  $L^2$ -error as function of the order of magnitude of  $\epsilon$  (log-scale)      (b)  $H^1$ -error as function of the order of magnitude of  $\epsilon$  (log-scale)

Figure 4.4: Linear wave problem, quadratic in time.  $L^2$ - and  $H^1$ -errors at time  $t = 1$  for a range of values of  $\epsilon$  in the convergence criterion (4.6),  $k \in \{0, 1, 2, 3\}$ , equal-order setting.

leading to homogeneous Dirichlet conditions,  $f := 0$  and initial conditions

$$u_0(x, y) := 0, \quad v_0(x, y) := \sin(\pi x) \sin(\pi y). \quad (4.29)$$

Figure 4.5 shows the number of iterations needed to satisfy the convergence criterion (4.6) for this manufactured solution, in the time interval  $t \in [0, 2]$ . The manufactured solution is sinusoidal in time of period  $\sqrt{2}$ , so that the time interval covers more than one period. The behavior over time differs from the quadratic case. The number of iterations still increases as  $\epsilon$  decreases, but for most values of  $\epsilon$  and polynomial orders  $k \in \{1, 2, 3\}$ , the periodicity of the solution is passed on to the number of splitting iterations. More precisely, the number of iterations decreases at the times when the solution is maximal and its velocities are minimal ( $t = m\frac{\sqrt{2}}{4}, m \in \{1, 3, 5\}$ ). This can be explained by the fact that, since the velocities of the solution are small, starting the splitting with the value at the previous time step is a good approximation of the value at the current time step. Conversely, when the velocities are large, the solution changes more swiftly at each time step, and more iterations are needed to converge since the fixed-point procedure is started further away from the solution. This can also explain why the number of iterations increases steadily for the previous test case, where the solution is quadratic in time: as time progresses, the difference between the previous solution and the current solution increases. The error as a function of  $\epsilon$  is very similar to the quadratic case and is not displayed for brevity.

**Conclusion on the impact of the number of splitting iterations.** The number of splitting iterations must be such that the truncation error of the splitting is negligible with respect to the time and space discretization errors. To reach this precision, the convergence criterion of the splitting procedure must be small enough. Taking a smaller convergence criterion does not impact the quality of the solution, but more iterations are needed, hence making the computations more expensive. Finding this optimal convergence criterion requires an a priori knowledge on the time and space discretization errors, which is in practice not available.

## 4.2.2 Optimal choice of the scaling factor

We now focus on the dependence of the error and of the execution time on the value of  $\gamma$ . Since we saw in Section 4.1.3 that the value of  $\gamma$  impacts the CFL condition (see Tables 4.3 and 4.4) and the number

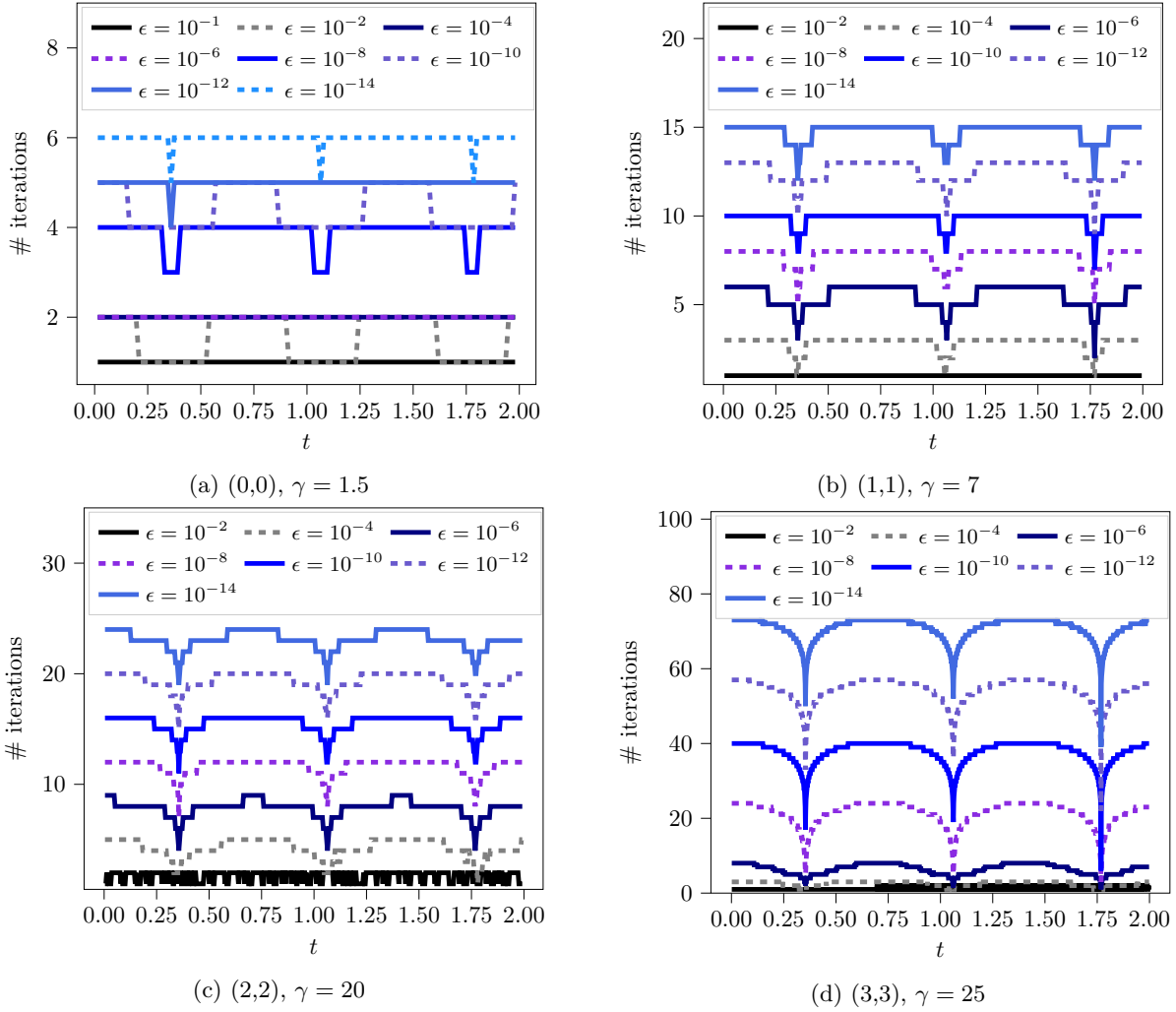


Figure 4.5: Linear wave problem, sinusoidal in time - Number of iterations over time ( $t \in [0, 2]$ ) in order to satisfy the convergence criterion (4.6),  $k \in \{0, 1, 2, 3\}$ , equal-order setting.

of iterations in the splitting procedure, we want to measure the impact on the total execution time. And since changing  $\gamma$  changes the solution of (4.2), we want to verify that this impact is limited.

We still consider the manufactured solution given by equation (4.27). Since there is no error in time, the error is the sum of two errors, coming from the space discretization on the one hand and from the splitting procedure on the other hand. The splitting error has two sources: the truncation error of the iterative procedure and the error caused by the variation of  $\gamma$ . The previous test case illustrated the convergence of the truncation error. The impact of the value of  $\gamma$  on the error is actually quite moderate (see Section 4.2.3 for further insight). Regarding the execution time, we expect that the critical time step decreases as  $\gamma$  increases, thus impacting the execution cost. Further insight into this impact is again provided in subsection 4.2.3.

Figure 4.6 shows the discrete energy error at  $t = 0.1$  as a function of the execution time for the equal- and mixed-order settings, with  $k \in \{0, 1, 2\}$  for different values of  $\gamma$ . Order 3 is not shown for brevity. The discrete energy error, computed using the projection of the exact solution onto the mesh cells and faces, is the sum of the  $L^2$ -error on the cell velocities and the hybrid  $H^1$ -error evaluated using the gradient reconstruction. Computations are performed on a sequence of unstructured triangular meshes,

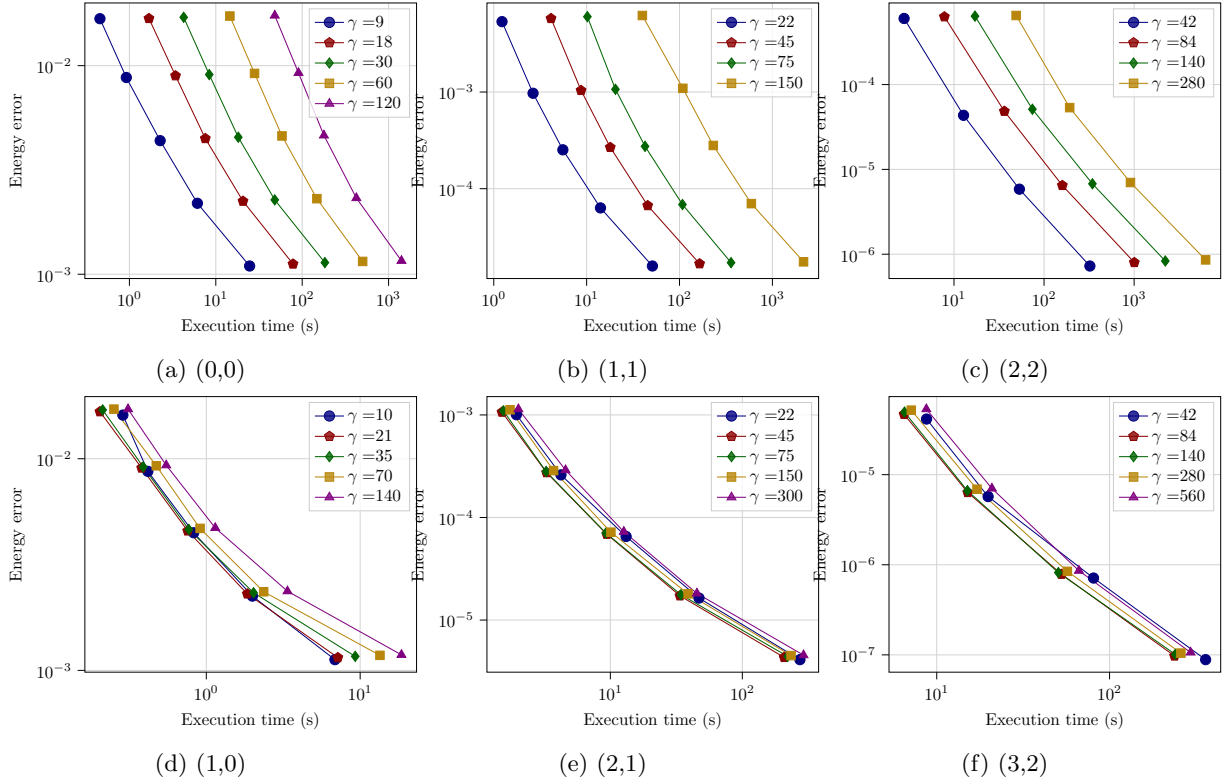


Figure 4.6: Linear wave problem - Energy error at  $t = 0.1$  as a function of the execution time,  $\gamma \in \{1.5\gamma^*, 3\gamma^*, 5\gamma^*, 10\gamma^*\}$ ,  $k \in \{0, 1, 2\}$ , equal- and mixed-order settings.

obtained by repeated refinements of a coarse initial mesh. This ensures that all the meshes have the same regularity. Computations are run on 16 cores with distributed memory, and the meshes are equally distributed on all the processors (i.e. each processor has approximately the same number of mesh cells). Each marker on the curves of Figure 4.6 corresponds to the error and execution time on a mesh from the sequence. The scaling parameter  $\gamma$  takes the values  $\{1.5\gamma^*, 3\gamma^*, 5\gamma^*, 10\gamma^*\}$  and, when the execution cost is reasonable,  $20\gamma^*$ . For each value of  $\gamma$ , the critical time step  $\Delta t^{\text{opt}}(\gamma)$  can be computed via equation (4.26), and the actual time step is then set to  $\Delta t := 0.8\Delta t^{\text{opt}}(\gamma)$ . Unlike the experiments in Figure 4.2, the number of iterations for the splitting procedure is not fixed here. We rather stop the splitting procedure when the relative norm of the increment is smaller than  $\epsilon := 10^{-11}$ .

The first salient point is that, for a given mesh, the error does not depend on  $\gamma$ . Indeed, all the corresponding markers are almost horizontally aligned. This means that a reasonable increase in  $\gamma$  does not deteriorate the quality of the solution. The second salient point concerns the execution time. Here, the behavior differs between equal- and mixed-order settings. In the equal-order setting, the larger  $\gamma$ , the more expensive the execution. This could have two origins: a smaller critical time step for larger  $\gamma$ , requiring more time steps to reach the final time, or a higher number of splitting iterations. Figure 4.7 and Table 4.5 answer this question. The mean number of splitting iterations for each time step is displayed in Figures 4.7a and 4.7b for polynomial orders (0,0) and (1,1), respectively. As  $\gamma$  increases, more iterations are needed for the splitting procedure to converge. This can be explained by the results in Figure 4.11 which show that the spectral radius of the iteration matrix converges to 1 when  $\gamma$  grows. Moreover, Table 4.5 reports the critical time step  $\Delta t^{\text{opt}}(\gamma)$  for the mesh size  $h = 0.1$  and all polynomial orders. We observe that the critical time step  $\Delta t^{\text{opt}}(\gamma)$  decreases when  $\gamma$  increases. Since these two factors act together, this leads to a swift increase of the execution cost with  $\gamma$ .

On the contrary, in the mixed-order setting, there seems to be an optimal value of  $\gamma$ , for which the



execution cost is the smallest. Indeed, the number of iterations per time step decreases as  $\gamma$  increases (see Figures 4.7c and 4.7d), but the critical time step decreases as well. For the considered polynomial orders, this optimum is in the interval  $\gamma \in [3\gamma^*, 5\gamma^*]$ . Moreover, comparing Figures 4.7e and 4.7f with Figures 4.7g and 4.7h highlights that, for a given polynomial order, the total number of iterations in the mixed-order setting is smaller than that in the equal-order setting, by a factor of 10 or more. At the same time, Table 4.5 shows that the critical time step is just slightly smaller in the mixed-order setting. Combining these two factors makes the mixed-order computations faster than the equal-order ones, on a given mesh and for a given polynomial order.

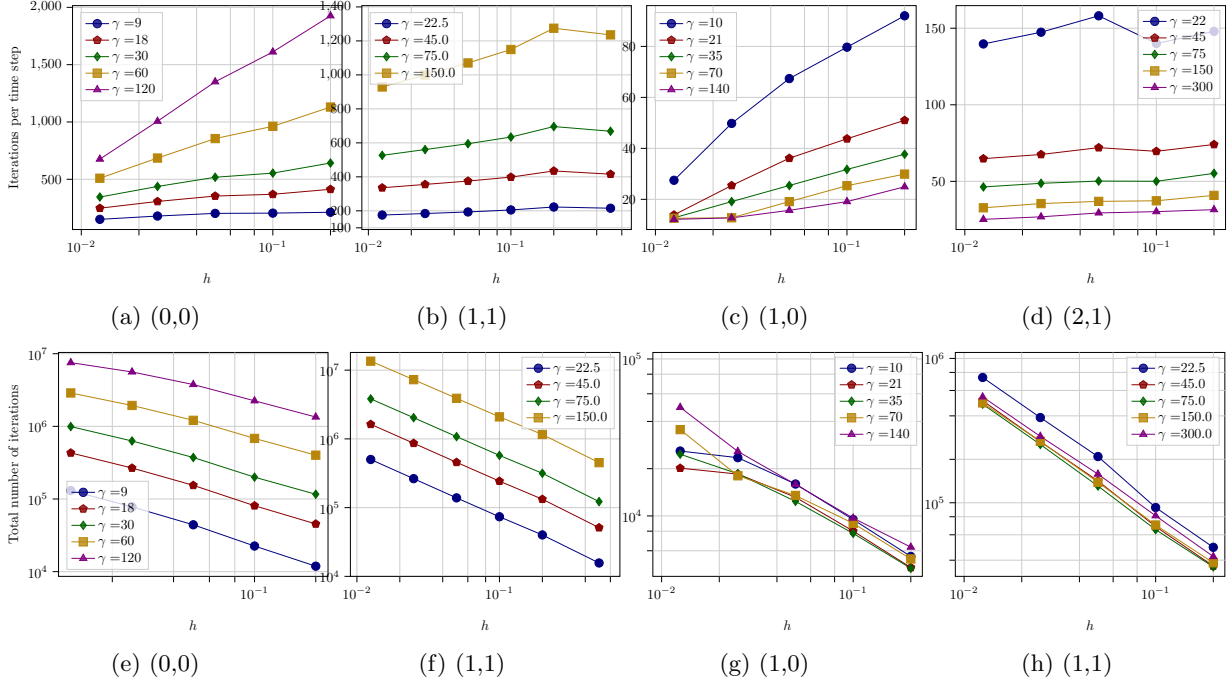


Figure 4.7: Linear wave problem - Mean number of splitting iterations in log scale per time step (top) and total number of splitting iterations (bottom) in log scale, as a function of the mesh size,  $k \in \{0, 1\}$ , equal- and mixed-order settings.

Order	(0,0)	(1,1)	(2,2)	(3,3)	(1,0)	(2,1)	(3,2)	(4,3)
$1.5\gamma^*$	4.61e-2	1.40e-2	3.75e-3	1.12e-3	4.25e-2	7.57e-3	1.75e-3	4.037e-4
$3\gamma^*$	2.30e-2	8.23e-3	2.33e-3	7.11e-4	2.75e-2	5.09e-3	1.23e-3	2.88e-4
$5\gamma^*$	1.39e-2	5.51e-3	1.67e-3	5.14e-4	2.06e-2	3.84e-3	9.46e-4	2.23e-4
$10\gamma^*$	7.07e-3	2.75e-3	1.09e-3	3.34e-4	1.42e-2	2.66e-3	6.58e-4	1.58e-4
$20\gamma^*$	3.59e-3	1.27e-3	5.56e-4	1.82e-4	9.89e-3	1.86e-3	4.61e-4	1.11e-4

Table 4.5: Linear wave problem - Critical time step  $\Delta t^{\text{opt}}(\gamma)$ ,  $h = 0.1$ ,  $\gamma \in \{1.5\gamma^*, 3\gamma^*, 5\gamma^*, 10\gamma^*, 20\gamma^*\}$ ,  $k \in \{0, 1, 2, 3\}$ , equal- and mixed-order settings.

The above results allow us to determine the optimal value of  $\gamma$ :  $\gamma = \gamma^*$  in the equal-order setting and  $\gamma \in [3\gamma^*, 5\gamma^*]$  in the mixed-order setting. Using these values of  $\gamma$ , we can now compare the error as a function of the execution cost for all polynomial orders. Figure 4.8 shows the energy error for polynomial degrees  $k \in \{0, 1, 2, 3\}$ . These results illustrate the fact that, for a given error or a given execution time, taking higher orders on a coarse mesh is more efficient than lower orders on a fine mesh. Moreover, the mixed- and equal-order settings essentially lead to similar curves, but with an offset. Indeed, on a given

mesh, the mixed-order setting is faster than the equal-order in order to reach the same error.

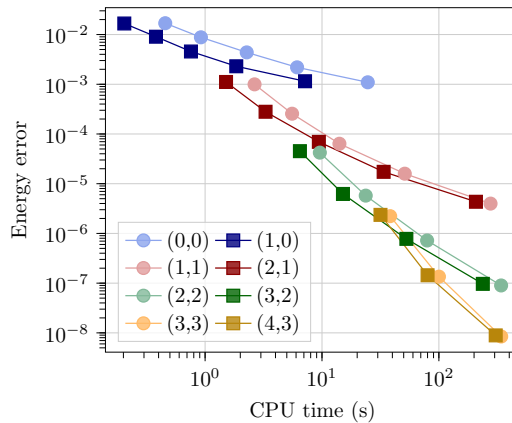


Figure 4.8: Linear wave problem - Energy error with optimal  $\gamma$  as a function of the execution time, with  $(k, \gamma) \in \{(0, 9), (1, 22.5), (2, 42), (3, 67.5)\}$  in the equal-order setting and  $(k, \gamma) \in \{(0, 1), (1, 45), (2, 84), (3, 135)\}$  in the mixed-order setting.

### 4.2.3 Consequences of a large scaling factor

The experiments of Section 4.2.2 showed that taking  $\gamma \in [1, 100]$  has a moderate to negligible impact on the  $L^2$ - and  $H^1$ -errors. It also showed that there exists an optimal range of values of  $\gamma$  in terms of execution time. This range is  $\gamma \in [3\gamma^*, 5\gamma^*]$  for the mixed-order setting and  $\gamma$  as close to  $\gamma^*$  as possible in the equal-order setting. We did not investigate, however, the impact of *very large* values of  $\gamma$ , say  $\gamma = 1000\gamma^*$ . This is the aim of this Section.

**Impact on the error in the static case** We consider the following problem: Find  $u \in H_0^1(\Omega)$  such that  $(\mu^2 \nabla u, \nabla w)_\Omega = (f, w)_\Omega$  for all  $w \in H_0^1(\Omega)$ . The HHO discretization consists of finding  $\hat{u}_h \in \hat{\mathcal{U}}_{h,0}^{l,k}$  such that

$$b_h(\hat{u}_h, \hat{w}_h) + s_h(\hat{u}_h, \hat{w}_h) = (f, w_\mathcal{T})_\Omega, \quad \forall \hat{w}_h \in \hat{\mathcal{U}}_{h,0}^{l,k}. \quad (4.30)$$

We consider an analytical test case with  $u(x, y) := \sin(\pi x) \sin(\pi y)$  obtained using  $\mu := 1$  in the domain  $\Omega := (0, 1)^2$  and the source term  $f(x, y) := 2\pi^2 \sin(\pi x) \sin(\pi y)$ . We study the  $L^2$ -error on the cells,  $\mathcal{E}_1 := (\Delta \mathbf{U}_\mathcal{T}^\top M_{\mathcal{T}\mathcal{T}} \Delta \mathbf{U}_\mathcal{T})^{\frac{1}{2}}$ , on the faces,  $\mathcal{E}_2 := h^{-\frac{1}{2}} (\Delta \mathbf{U}_\mathcal{F}^\top M_{\mathcal{F}\mathcal{F}} \Delta \mathbf{U}_\mathcal{F})^{\frac{1}{2}}$  (evaluated using the face mass matrix), as well as the (broken)  $H^1$ -error,  $\mathcal{E}_3 := (\Delta \hat{\mathbf{U}}^\top \mathcal{K}(\gamma) \Delta \hat{\mathbf{U}})^{\frac{1}{2}}$ , computed using the complete stiffness matrix (this error thus includes the stabilization and is therefore affected by the value of  $\gamma$ ). We test the values  $\gamma \in \{1, 10, 100, 1000, 10000\}$ . Figure 4.9 shows the three errors for the equal-order setting and  $k \in \{0, 1\}$ . Other polynomial orders are not displayed for brevity, since the results are similar. For the  $L^2$ - and  $H^1$ -errors, the convergence rates are, as expected,  $(k+2)$  and  $(k+1)$ , respectively. The main observation is that, for a fixed  $h$ , the error converges to some value as  $\gamma$  increases. In other words, the limit  $\gamma \rightarrow \infty$  does not yield a diverging solution, but a suitable solution with optimal convergence rates. Interestingly, this limit solution is reached for a rather small value of  $\gamma$  since the error curves for  $\gamma \geq 10$  almost overlap. Regarding the  $L^2$ -error on the faces, the behavior differs between  $k=0$  and the higher orders. For the higher orders, the behavior is the same as for the other errors, whereas for the lowest order, the error does not change when  $\gamma$  increases. This indicates that the value on the faces is not affected by the scaling of the stabilization when  $k=0$ .

The behavior as  $\gamma \rightarrow \infty$  can be partially explained by a series expansion in  $\gamma$ . Let us plug the Ansatz  $\hat{u}_h = \sum_{n=0}^{\infty} \hat{u}_h^n \gamma^{-n}$  into (4.30) and consider test functions having only nonzero face components. In order

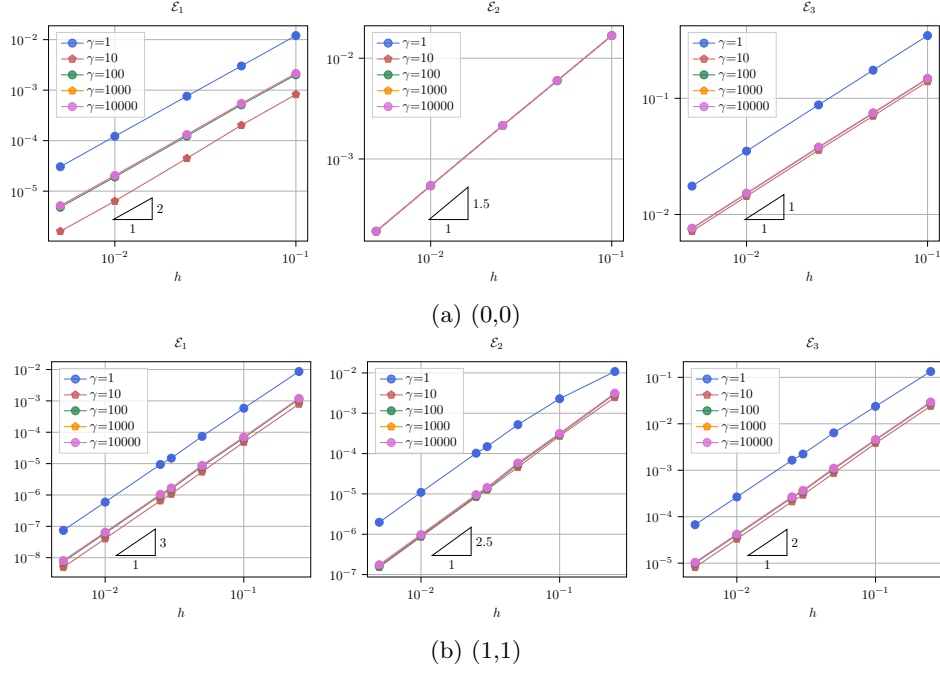


Figure 4.9: Large values of  $\gamma$  - Errors  $\mathcal{E}_1, \mathcal{E}_2, \mathcal{E}_3$  on the scalar Laplacian as a function of the mesh size for  $\gamma \in \{1, 10, 100, 1000, 10000\}$ , equal-order setting with  $k \in \{0, 1\}$ .

to analyze the dependance in  $\gamma$ , we consider in this paragraph that  $s_h$  does not contain the factor  $\gamma$ . This yields

$$\gamma s_h(\hat{u}_h^0, (0, w_{\mathcal{F}})) + b_h(\hat{u}_h^0, (0, w_{\mathcal{F}})) + s_h(\hat{u}_h^1, (0, w_{\mathcal{F}})) + \dots = 0, \quad \forall w_{\mathcal{F}} \in \mathcal{U}_{\mathcal{F},0}^k. \quad (4.31)$$

The dominant term in  $\gamma$  gives  $s_h(\hat{u}_h^0, w_{\mathcal{F}}) = 0$  for all  $w_{\mathcal{F}} \in \mathcal{U}_{\mathcal{F},0}^k$ . In the mixed-order setting, this condition reads

$$\sum_{T \in \mathcal{T}_h} \sum_{F \in \mathcal{F}_T} \eta_{TF}^{-1} (\Pi_F^k(u_T^0) - u_F^0, w_F) = 0, \quad \forall w_{\mathcal{F}} \in \mathcal{U}_{\mathcal{F},0}^k. \quad (4.32)$$

Letting  $\mathcal{T}_F$  collect the one or two cells adjacent to a face  $F$ , this translates into the following direct expression of the face unknowns:

$$u_F^0 = \frac{1}{\#\mathcal{T}_F} \sum_{T \in \mathcal{T}_F} \Pi_F^k(u_T^0). \quad (4.33)$$

This means that, in the mixed-order setting, when  $\gamma \rightarrow \infty$ , the face unknowns on  $F$  are equal to the projection onto  $\mathbb{P}_d^k(F)$  of the mean-value of the cell unknowns from the two adjacent cells if  $F$  is an interface and of the trace of the cell unknown if  $F$  is a boundary face. The equal-order setting does not lead to an explicit expression of the face unknowns in terms of the cell unknowns. Indeed, in this case, the condition  $s_h(\hat{u}_h, (0, w_{\mathcal{F}})) = 0$  for all  $w_{\mathcal{F}} \in \mathcal{U}_{\mathcal{F},0}^k$  is not sufficient to determine the face unknowns in terms of the cell unknowns, because the spectrum of  $\mathcal{S}_{\mathcal{F}\mathcal{F}}$  contains at least one zero eigenvalue as shown in Figure 4.10. The smallest eigenvalue is nonzero owing to round-off errors. Higher-order polynomials lead to the same behavior and are not displayed for brevity.

**Spectral radius of the iteration matrix as  $\gamma \rightarrow \infty$**  Let us now focus on the spectral radius of  $\mathcal{S}_{\mathcal{F}\mathcal{F}}^{-1} \mathcal{B}_{\mathcal{F}\mathcal{F}}$  in the mixed-order setting and of  $(\mathcal{S}_{\mathcal{F}\mathcal{F}}^*)^{-1} (\mathcal{B}_{\mathcal{F}\mathcal{F}} + \mathcal{Z}_{\mathcal{F}\mathcal{F}})$  in the equal-order setting. In Section 4.1.3, a study of the value of  $\gamma^*$  as a function of the polynomial order and the mesh regularity has been performed. Figure 4.11 reports the spectral radius of the iteration matrix as  $\gamma \rightarrow \infty$  on an unstructured triangular mesh with  $h = 0.1$ . In the mixed-order setting, the spectral radius of  $\mathcal{S}_{\mathcal{F}\mathcal{F}}^{-1} \mathcal{B}_{\mathcal{F}\mathcal{F}}$  goes to zero,

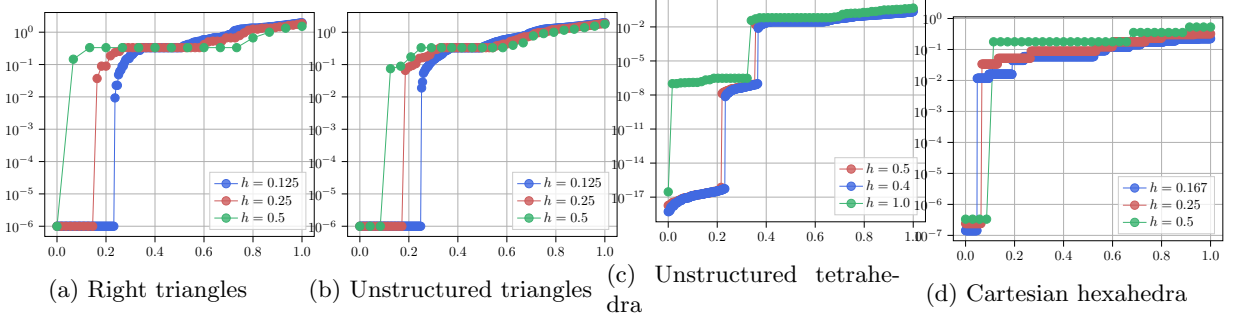


Figure 4.10: Large values of  $\gamma$  - Eigenvalues of  $\mathcal{S}_{\mathcal{F}\mathcal{F}}$  on different meshes for polynomial orders  $(0,0)$ .

whereas the spectral radius of  $(\mathcal{S}_{\mathcal{F}\mathcal{F}}^*)^{-1}(\mathcal{B}_{\mathcal{F}\mathcal{F}} + \mathcal{Z}_{\mathcal{F}\mathcal{F}})$  converges to one in the equal-order setting. This is the expected behavior and it pinpoints the existence of an optimal value of  $\gamma$  for the equal-order setting, since the smaller  $\gamma$ , the less iterations needed. The minimum of  $\gamma$  is very close to  $\gamma^*$ , which corroborates the experience of paragraph 4.2.2.

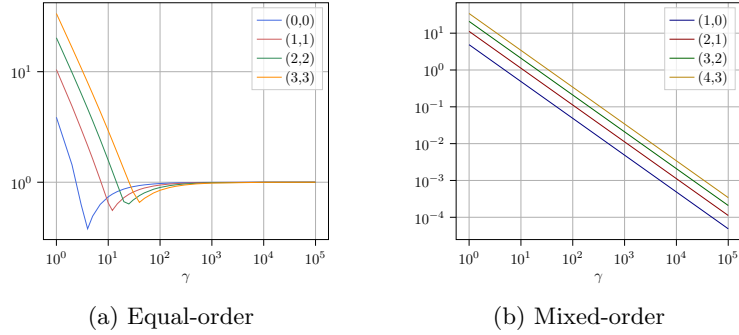


Figure 4.11: Large values of  $\gamma$  - Spectral radius of the iteration matrix as a function of  $\gamma$ , equal-order setting (left) and mixed-order setting (right) on an unstructured triangular mesh with  $h = 0.1$ .

**Stability condition as  $\gamma \rightarrow \infty$**  Equation (4.26) gives a formula to obtain the critical time step  $\Delta t^{\text{opt}}(\gamma)$  as a function of  $\gamma$ . Figure 4.12 reports the rate  $\frac{\Delta t^{\text{opt}}(\gamma)}{h}$  on 2D unstructured quadrangular and triangular meshes, and 3D Cartesian hexahedral and unstructured tetrahedral meshes. For each polynomial order, the value of  $\gamma^*$  is indicated by a bullet. As  $\gamma \rightarrow \infty$ ,  $\Delta t^{\text{opt}}(\gamma)$  go to zero as  $\frac{1}{\sqrt{\gamma}}$ . Thus, taking  $\gamma \gg \gamma^*$  is not optimal, since this entails very small time steps. Fortunately, Table 4.1 has already shown that the value of  $\gamma^*$  need not be too large. The fact that  $\Delta t^{\text{opt}}(\gamma) \rightarrow 0$  as  $\gamma \rightarrow \infty$  is easily proven using the min-max theorem for the largest eigenvalue. Indeed, since  $\mathcal{B}$  is a positive matrix, we have

$$\rho(\mathcal{D}(\gamma)) = \max_{X=(X_{\mathcal{T}}, X_{\mathcal{F}}) \in \tilde{\mathcal{U}}_{h,0}^{l,k}} \frac{X^{\top} \mathcal{A} X}{X^{\top} \mathcal{M} X_{\mathcal{T}}} \geq \max_{X_{\mathcal{T}} \in \mathcal{U}_{\mathcal{T}}^l} \frac{X_{\mathcal{T}}^{\top} \mathcal{A}_{\mathcal{T}\mathcal{T}} X_{\mathcal{T}}}{X_{\mathcal{T}}^{\top} \mathcal{M} X_{\mathcal{T}}} \geq \gamma \rho(\mathcal{M}^{-1} \mathcal{S}_{\mathcal{T}\mathcal{T}}).$$

This shows that  $\Delta t^{\text{opt}}(\gamma) \leq \frac{2}{\sqrt{\gamma \rho(\mathcal{M}^{-1} \mathcal{S}_{\mathcal{T}\mathcal{T}})}} \underset{\gamma \rightarrow \infty}{\sim} \frac{1}{\sqrt{\gamma}}$ . All the mixed-order curves in Figure 4.12 show the expected asymptotic behavior. For the equal-order curves, this asymptotic behavior is recovered for the larger values  $\gamma$ , but we also observe that there is a large range of values of  $\gamma$  for which  $\Delta t^{\text{opt}}(\gamma) \simeq C\gamma^{-1}$ , where  $C$  is some constant depending on the polynomial order and the mesh regularity. For triangular meshes with  $k = 3$  and for hexahedral meshes with  $k = 0$ , the asymptotic rate  $\gamma^{-\frac{1}{2}}$  is not even reached for  $\gamma = 10^{10}$ .

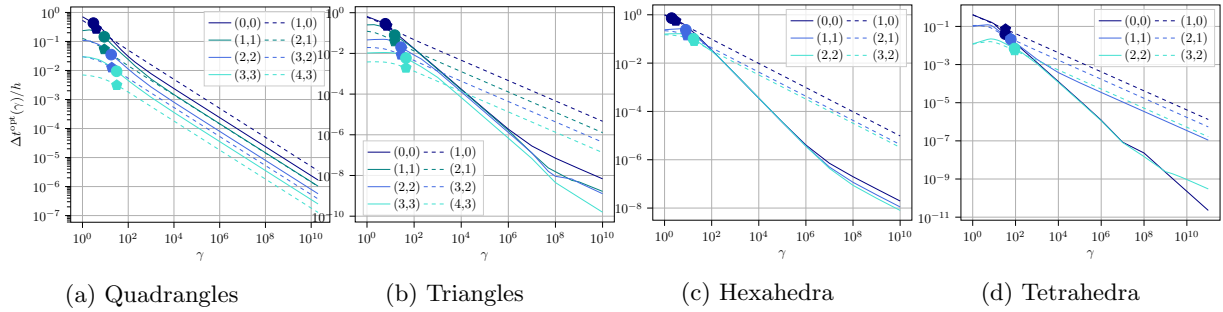


Figure 4.12: Large values of  $\gamma$  - Ratio  $\frac{\Delta t^{\text{opt}}(\gamma)}{h}$  as a function of  $\gamma$  on 2D and 3D unstructured meshes.

From these experiments with large  $\gamma$ , we conclude that large values of  $\gamma$  do not lead to large errors. In the static case, the errors with large values of  $\gamma$  are even smaller than with  $\gamma = 1$  (see Figure 4.9). However, large values of  $\gamma$  have a negative impact on the CFL condition, as the critical time step induced by this condition goes to 0 as  $\gamma \rightarrow \infty$ . For wave propagation problems, values of  $\gamma$  larger than  $1000\gamma^*$  should not be considered.

#### 4.2.4 Conclusion on the choice of the splitting parameter.

The above experiments allowed to grasp the impact of the value of  $\gamma$  on the number of iterations and the quality of the solution. In both equal- and mixed-order settings, there is an optimal value in terms of execution time:  $\gamma \simeq 1.5\gamma^*$  in the equal-order setting and  $\gamma \in [3\gamma^*, 5\gamma^*]$  in the mixed-order setting. As long as  $\gamma \in [\gamma^*, 10\gamma^*]$ , the impact on the quality of the solution is negligible.

As  $\gamma \rightarrow \infty$ , the obtained solution does not diverge, but the CFL condition gives  $\Delta t \rightarrow 0$  with rate  $\gamma^{-\frac{1}{2}}$ , thus increasing the number of time steps, and making the computations very expensive. Large values of  $\gamma^*$  could have been considered in the mixed-order setting since it makes the splitting converge in less iterations. The trade-off between splitting iterations and CFL condition leads to the range of optimal values of  $\gamma$  stated above.

Finally, we notice that the equal-order setting requires more iterations to converge, but when choosing  $\gamma = 1.5\gamma^*$ , it is competitive with respect to the mixed-order setting, as shown, for instance, in Figure 4.8.

### 4.3 Comparison to finite elements

The splitting procedure has been developed with nonlinear applications in mind, since, in this case, the cost of solving the static coupling at each time step becomes quite important. It is, however, interesting to measure the performance of the splitting in the linear case too, in particular by comparing its performance to that of the standard finite element method.

#### 4.3.1 Presentation of the test case

It is expected that the HHO method presents several advantages compared to finite elements for the simulation of acoustic waves, in particular a better performance with contrasted material coefficients. Let us consider the following test case (see Figure 4.13a):  $\Omega := (-1.5, 1.5)^2$  with two subdomains  $\Omega_1 := (-1.5, 1.5) \times (0, 1.5)$  and  $\Omega_2 := (-1.5, 1.5) \times (-1.5, 0)$  with respective speeds of sound  $c_1 := 3$  and  $c_2 := \sqrt{3}$ . A source is placed in the upper domain  $\Omega_1$  at coordinates  $S := (0, \frac{2}{3})$ , and two sensors are placed at symmetric positions in both domains at the points  $P_1 := (\frac{1}{3}, \frac{1}{3})$  and  $P_2 := (\frac{1}{3}, -\frac{1}{3})$ . The wave first propagates in  $\Omega_1$ , then is partly reflected on the boundary and partly transmitted to the second

subdomain. The sensor  $P_1$  should measure the initial wave and the reflected wave, while the sensor  $P_2$  should measure the transmitted wave.

The simulations are run from  $t = 0$  until  $t = 1$ . Reflecting boundary conditions are enforced (homogeneous Dirichlet boundary conditions). The source at  $t = 0$  is the derivative of a Ricker wavelet in time and a Dirac in space at point  $S$ . Figure 4.13b illustrates the one-dimensional Ricker wavelet centered at 0 and its expression is

$$f_R(t) := (1 - 2\pi^2 f_p^2 t^2) e^{-\pi^2 f_p^2 t^2},$$

with  $f_p$  the peak frequency. No analytical solution is known for this problem, but an open-source software **Gar6More2D** developed at INRIA computes the semi-analytical solution at given points under the assumption of an infinite domain of propagation. Therefore, the comparison with the present simulations is valid only until the waves reflecting on the walls reach the sensors. These maximal validity times can be computed by geometric arguments (see Table 4.6 below). For more details on the implementation of **Gar6More2D**, see the open Gitlab repository [89]. When running wave simulations with the HHO

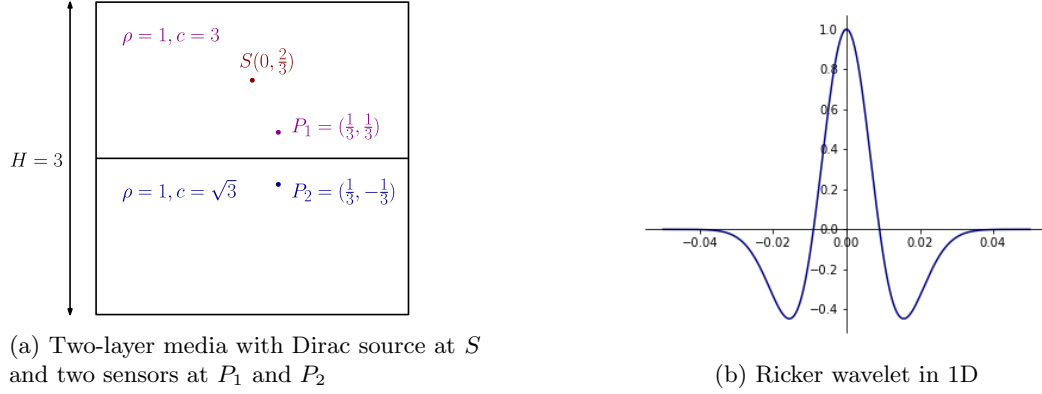


Figure 4.13: Linear wave propagation - Setup for the numerical experiment: scheme of the physical domain (left) and shape of the initial condition (right)

method, we need the initial conditions  $u_0$  and  $v_0$  to be used in (3.3). In order to simulate the same setting as in **Gar6More2D**, the initial condition is also a Ricker wavelet but, in space, it is expressed in terms of the variable  $r = \sqrt{(x - x_S)^2 + (y - y_S)^2}$ , the distance to the source. One can show that this initial condition is equivalent to the initial source used in **Gar6More2D**, with only a time shift. This time shift is a given data of the semi-analytical simulation with **Gar6More2D**. For a given time peak frequency  $f_p$  and a speed of sound  $c_p$  in the domain of the source (here  $c_1$ ), the corresponding wave number  $k_p$  is given by  $k_p = c_p^{-1} f_p$  and the initial condition reads

$$u_0(x, y) = 0, \quad v_0(x, y) := A e^{-\pi^2 r^2 k_p^2}, \quad (4.34)$$

with  $A$  an amplitude coefficient coming from the derivation of the Ricker wavelet.

The variables of interest are the pressure  $u$  and the components of its gradient  $v_x$  and  $v_y$ . Notice that these two quantities can be viewed as the velocity components. For both sensors and the numerical experiment setting illustrated in Figure 4.13a, the semi-analytical solution is displayed in Figure 4.14. The semi-analytical solution is valid for an infinite domain, whereas the simulations are computed with a finite domain  $\Omega$  and reflecting boundary conditions. Thus, for each sensor, there is a time  $t_r$  when the reflections on the boundaries reach the sensor. After this time, the analytical and numerical solutions differ. The times can be easily computed, as well as the time needed for the initial wave to reach the first sensor  $P_1$  and the time for the wave reflected on the interface or on the top and right boundaries to reach  $P_1$ . All the values of these times for the present setting are given in Table 4.6. The time at which the reflected wave reaches  $P_2$  is harder to compute analytically and numerical experiments show that this time is larger than  $t = 0.9$ . The first two values of Table 4.6 are in agreement with the peaks

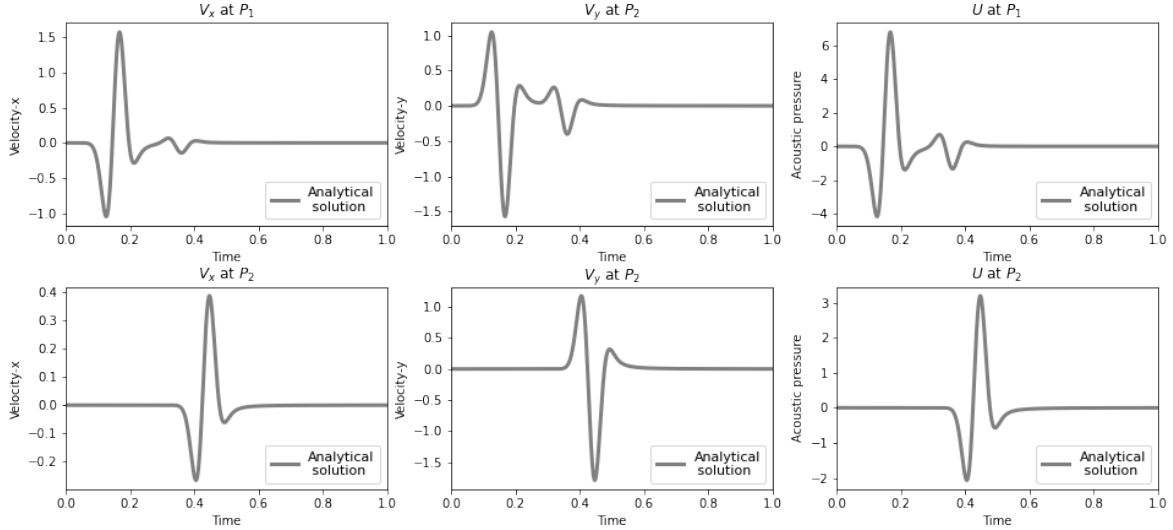


Figure 4.14: Linear wave propagation - Analytical solution:  $x$ -component of the gradient  $v_x$  (left),  $y$ -component  $v_y$  (center) and pressure  $u$  (right) at  $S_1$  (top row) and  $S_2$  (bottom row) on the time interval  $[0.0, 1.0]$ .

Wave path	Time
from $S$ to $P_1$	0.157
reflection on the interface to $P_1$	0.351
reflection on the top boundary to $P_1$	0.676
reflection on the right boundary to $P_1$	0.896

Table 4.6: Linear wave propagation - Times for different wave paths to reach the sensor  $P_1$

of the panels in Figure 4.14. The other two reflections can be observed in Figure 4.15 (first two panels). The reflections on the boundaries and the interaction with the wave reflected on the interface are visible in the two rightmost panels.

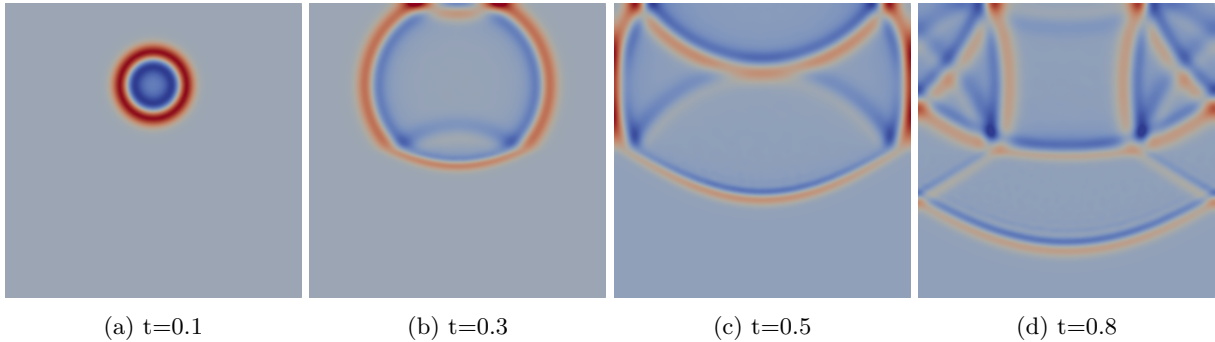


Figure 4.15: Linear wave propagation - Solution at different times, computed using  $P_1$  finite elements with MANTA

### 4.3.2 Error and execution time on a fixed mesh

The first experiment consists in comparing the different simulations on the same mesh and using the same time step, employing HHO with different orders and P1 finite elements. The simulations are run with the splitting procedure for HHO,  $\gamma = 1.5\gamma^*$ , since we verified in Section 4.2.2 that the solution is accurate, as long as the convergence criterion is small enough, and that, for the equal-order setting, the best value of  $\gamma$  is relatively close to  $\gamma^*$ . The scaling parameter and the time step for each polynomial order are given in Table 4.7.

Polynomial order	(0,0)	(1,1)	(2,2)	(3,3)	(1,0)	(2,1)	(3,2)	(4,3)	P1 FEM
$\gamma$	10	20	20	60	40	50	80	135	NA
$\Delta t, h = 0.05$	0.015	0.0025	0.002	0.001	0.004	0.0025	0.0015	0.0005	0.02
$\Delta t, h = 0.015$	0.0045	0.00075	0.0006	0.0003	0.0012	0.00075	0.00045	0.00015	0.006

Table 4.7: Linear wave propagation - Stabilization scaling and time step,  $h = 0.05$  and  $h = 0.015$ .

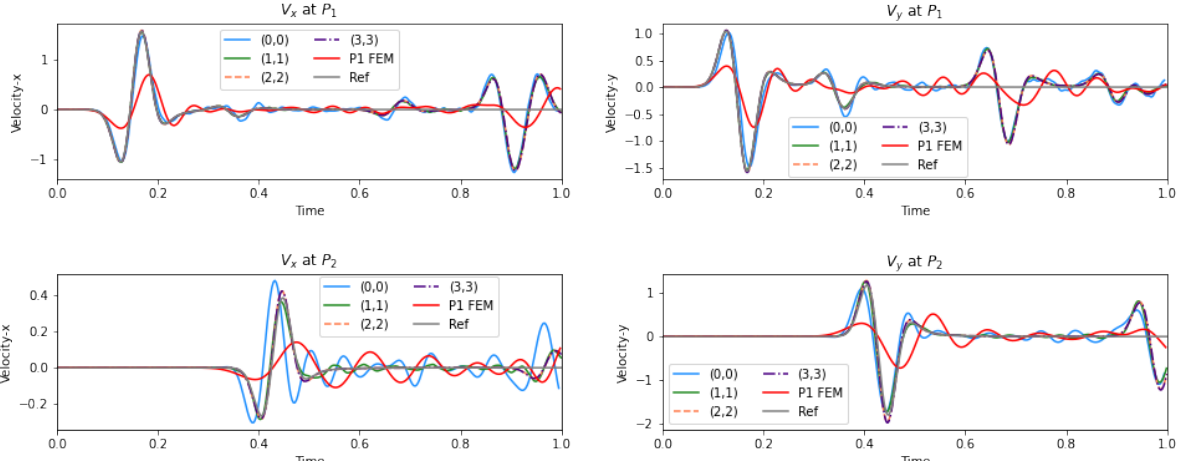
The results in Figure 4.16 illustrate a well-known phenomenon: as the polynomial order increases, the dispersion of the scheme is reduced. Figure 4.16a shows the solutions computed for  $h = 0.05$ . In this case, the mesh is very coarse. Visually, P1 finite elements behave the worst, even compared to HHO with order (0,0) (both methods are of order 1 and there are less degrees of freedom for P1 than for (0,0)). Figure 4.16b shows the solution on a finer mesh  $h = 0.015$ . In this case, the solutions in the upper domain (at  $P_1$ ) are all satisfactory, but in the lower domain (at  $P_2$ ), HHO with order (0,0) and finite elements are still dispersive. For instance, on the values at  $P_2$  on the finer mesh (Figure 4.16b) indicate that both components of the velocity are somewhat damped for finite elements, while they present some overshoots for HHO with order (0,0).

In terms of execution time, our current implementation of the splitting procedure is not conceived to perform at best in the linear case. Indeed, the coupling matrix does not change over the simulation, and providing that enough memory is available, this matrix can be factorized at the beginning of the simulation and reused at each time step. It can still be interesting to compare the execution time and the errors for HHO with splitting procedure and finite elements on a given mesh, in order to highlight some possible trends to be expected in the nonlinear case. The error and the execution time are reported in Table 4.8 for  $k \in \{0, 1, 2\}$  in both equal- and mixed-order settings. The solution is measured every time step and compared to the semi-analytical solution in the time interval  $[0, 0.6]$  to compute the  $\ell^2$ -error in time. Indeed, since the first reflection on the boundaries reaches the sensor  $P_1$  at  $t \simeq 0.676$ , the semi-analytical solution is a reference in the time interval  $[0, 0.6]$ . The measured errors are in agreement with the results of Figure 4.16: the higher the polynomial order and the finer the mesh, the smaller the error. The results of Table 4.8 also show that the HHO method with polynomial orders (1,1) or (2,2) is more expensive on a given mesh than P1 finite elements, but (much) more accurate, whereas HHO with polynomial order (0,0) yields comparable errors and smaller execution time than P1 finite elements. HHO in the mixed-order setting behaves very similarly to the equal-order setting in terms of error, but is quite faster, as already observed in Section 4.2. This first example shows the efficiency of the HHO method compared to P1 finite elements on the present heterogeneous wave propagation problem.

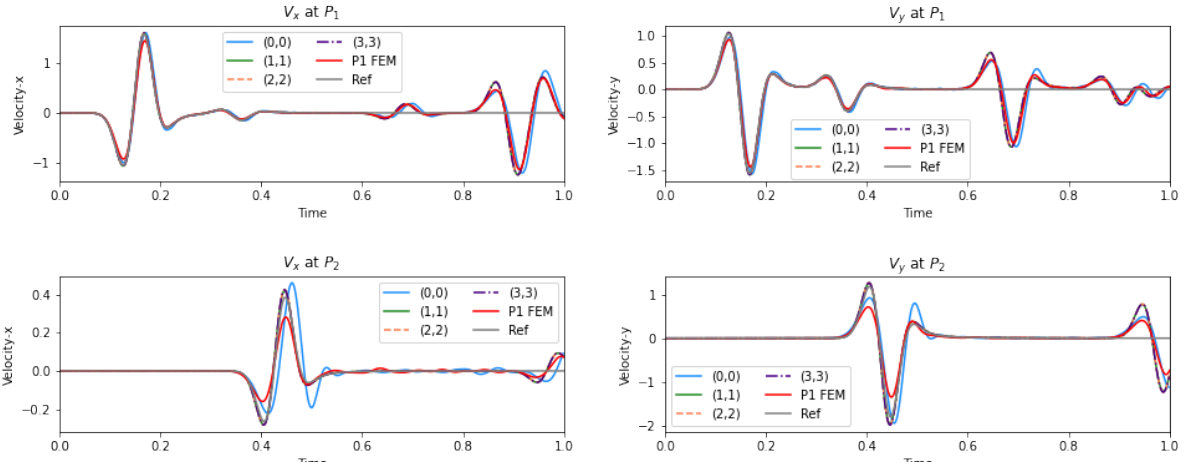
### 4.3.3 Error with a fixed number of degrees of freedom

Another way to compare the HHO method with the finite elements method is to perform numerical experiments with a fixed number of degrees of freedom, since the latter can be a measure of the computational time per time step, as well as a measure of the memory cost for each method. Figure 4.17 shows the solutions for computations with approximately the same number of degrees of freedom in space (up to 10% variation) and the same time step. This is achieved by taking different meshes for each method. This experiment is done on nonuniform triangular meshes and Table 4.9 gives the mesh size for each polynomial order. Table 4.10 reports the  $\ell^2$ -error in time on the two components of the velocity at both





(a)  $h = 0.05$

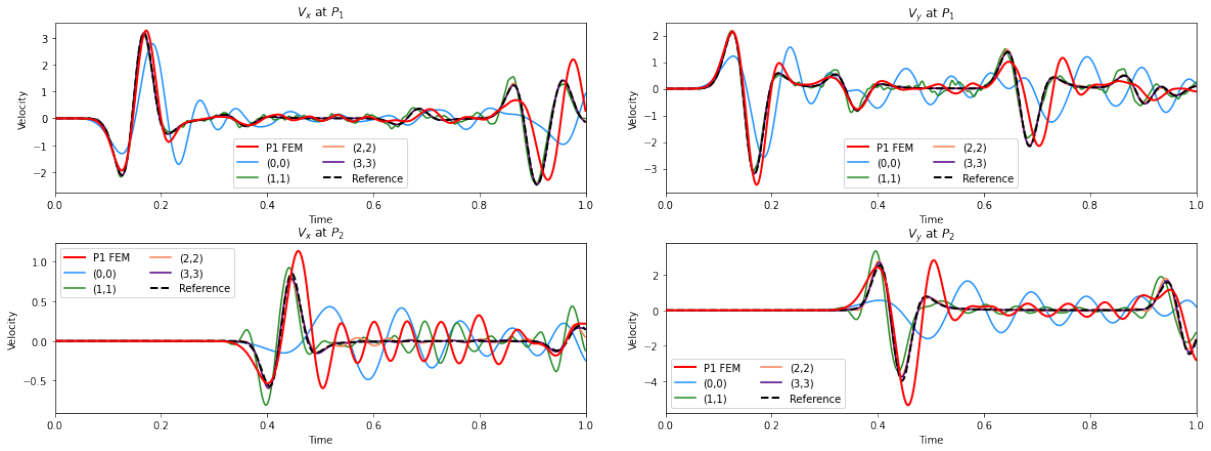


(b)  $h = 0.15$

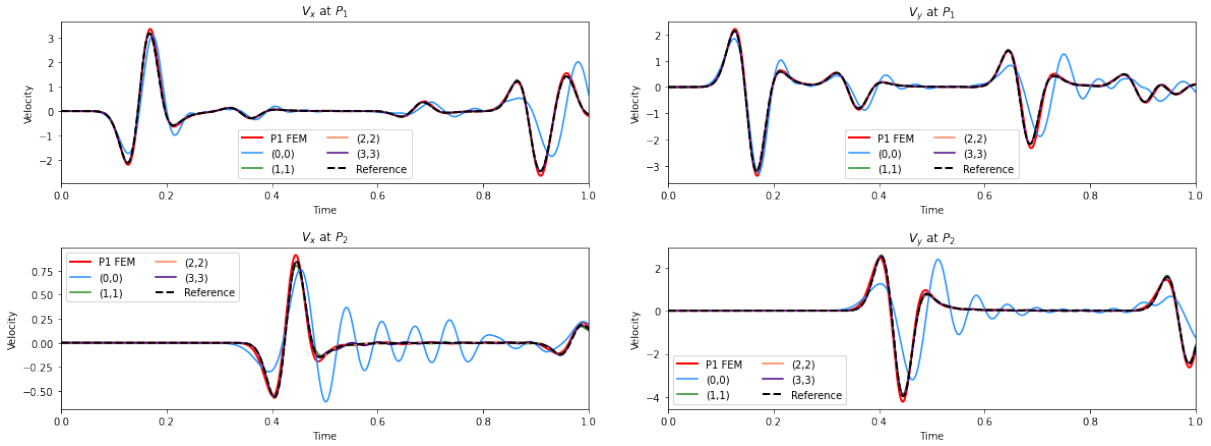
Figure 4.16: Simulations on the same mesh,  $h = 0.05$  and  $h = 0.15$ , HHO with polynomial order  $(k, k)$ ,  $k \in \{0, 1, 2, 3\}$ , P1 finite elements and semi-analytical reference solution.

Case	(0,0)	(1,1)	(2,2)	(1,0)	(2,1)	(3,2)	P1 FEM
$h = 0.05 - V_x$ at $S_1$	4.21e-2	1.67e-1	4.86e-4	1.22e-2	4.89e-4	3.86e-4	5.72e-2
$h = 0.05 - V_y$ at $S_1$	5.21e-2	4.96e-2	8.34e-4	1.82e-2	4.98e-4	8.15e-4	6.89e-2
$h = 0.05 - V_x$ at $S_2$	1.12e-2	4.86e-2	4.26e-4	2.62e-1	2.73e-4	3.24e-4	5.95e-2
$h = 0.05 - V_y$ at $S_2$	1.15e-1	1.44e-1	5.98e-3	2.32e-2	6.34e-3	2.62e-3	0.13
$h = 0.05$ - Exec. time (s)	30	83	194	25	38	97	26
$h = 0.015 - V_x$ at $S_1$	5.87e-4	2.26e-4	2.68e-4	1.10e-3	3.87e-4	3.58e-4	2.23e-3
$h = 0.015 - V_y$ at $S_1$	6.94e-4	1.87e-4	2.65e-4	7.45e-4	3.81e-4	3.22e-4	2.67e-3
$h = 0.015 - V_x$ at $S_2$	5.20e-4	2.98e-5	1.14e-5	1.06e-3	1.90e-5	1.33e-5	2.89e-3
$h = 0.015 - V_y$ at $S_2$	1.96e-3	4.71e-4	3.14e-4	1.27e-3	4.20e-4	3.99e-4	6.17e-2
$h = 0.015$ - Exec. time (s)	506	4417	10536	859	2465	5064	421

Table 4.8: Linear wave propagation -  $\ell^2([0, 0.6])$ -error on the gradient at the 2 sensor points and execution time on the mesh  $h = 0.05$  and  $h = 0.015$ ,  $k \in \{0, 1, 2\}$ .



(a) 13k dofs,  $\pm 10\%$



(b) 54k dofs,  $\pm 10\%$

Figure 4.17: Linear wave propagation - Simulations with approximately the same number of dofs, obtained with (0,0), (1,1), (2,2), (3,3) and P1 finite elements. The reference solution is obtained using (3,3) on 300k dofs ( $h=0.025$ ).

Number of dofs	(1,1)	(2,2)	(3,3)	P1 FEM
13k	0.074	0.091	0.11	0.03
54k	0.037	0.045	0.053	0.01

Table 4.9: Linear wave propagation - Size of the mesh for each method

Number of dofs	Quantity	(1,1)	(2,2)	(3,3)	P1 FEM
13k	$V_x$ at $S_1$	7.92e-3	3.48e-4	2.82e-4	4.86e-2
13k	$V_y$ at $S_1$	2.39e-2	7.58e-4	3.60e-4	7.27e-2
13k	$V_x$ at $S_2$	3.62e-1	3.85e-2	2.39e-2	8.59e-1
13k	$V_y$ at $S_2$	2.92e-1	3.61e-2	2.20e-2	8.30e-1
54k	$V_x$ at $S_1$	8.47e-4	3.00e-4	2.90e-4	1.07e-2
54k	$V_y$ at $S_1$	8.70e-4	3.90e-4	4.30e-4	1.01e-2
54k	$V_x$ at $S_2$	1.29e-2	2.14e-2	1.79e-2	9.67e-2
54k	$V_y$ at $S_2$	1.91e-2	1.87e-2	1.64e-2	8.65e-2

Table 4.10: Linear wave propagation -  $\ell^2$ -error for different methods and number of dofs

sensors. The error is computed using the semi-analytical solution in the same manner as in the previous experiment. We observe that, for a given number of dofs, the HHO method leads to a smaller error compared to P1 FEM. Errors for HHO with polynomial order (0,0) are not displayed because we saw in the previous experiment that this method does not capture well the contrasted properties on coarse meshes and leads to dispersive solutions (see Figure 4.17). For the coarser simulations (13k dofs), errors obtained with polynomial order (1,1) are about an order of magnitude larger than those obtained with polynomial orders (2,2) and (3,3). The two latter choices yield equivalent errors. The errors at the sensor  $P_2$  are larger because this sensor is located on the other side of the interface with respect to the source position. For the finer simulations (54k dofs), the errors obtained with polynomial order (1,1) in the upper domain are about 3 times larger than for polynomial orders (2,2) and (3,3), and are equivalent in the lower domain. All HHO orders yield errors at least four times smaller than the P1 FEM errors. This second experiment confirms the attractive performance of the HHO method on problems with contrasted properties.

## 4.4 Conclusion on the splitting procedure

In this chapter, we presented an iterative splitting procedure to solve the static coupling between the cell and face unknowns encountered at each time step. This splitting procedure is a fixed-point algorithm that requires to invert the block-diagonal matrix  $\mathcal{S}_{\mathcal{FF}}$  in the mixed-order setting and another block-diagonal matrix  $\mathcal{S}_{\mathcal{FF}}^*$  in the equal-order setting. In order for this splitting procedure to converge, the scaling parameter of the stabilization must be set larger than a constant depending only on the polynomial order and the mesh regularity. The minimal value for this scaling parameter has been given for a variety of mesh shapes and for polynomial orders  $k \in \{0, 1, 2, 3, 4\}$ .

Numerical experiments have been conducted in order to optimally set the value of all the splitting parameters: stabilization scaling parameter, convergence criterion and time step. These experiments showed that, provided that the stopping criterion is tight enough, the splitting error can be neglected in front of the discretization errors. The numerical experiments also showed that the optimal value of the scaling factor is  $\gamma = \gamma^*$  for the equal-order setting and  $\gamma \in [3\gamma^*, 5\gamma^*]$  in the mixed-order setting.

The numerical examples illustrated the efficiency of the HHO method compared to P1 FEM to solve problems with contrasted materials. In this case, the splitting procedure yields comparable execution time with smaller errors, on some coarse meshes. For the same number of degrees of freedom, the HHO

method with splitting yields smaller errors.

These observations emphasize the advantages of the HHO method with splitting, and motivate its use in the nonlinear case, where computational efficiency is of paramount importance. This topic is the subject of the next chapter.

## Chapter 5

# Explicit HHO method for the nonlinear acoustic wave equation

This chapter is dedicated to the adaptation of the splitting procedure devised in Chapter 4 to the nonlinear acoustic wave equation. We first expose the HHO space semi-discretization of the nonlinear acoustic wave equation, and then devise the splitting procedure. Two numerical experiments are presented where the differential operator in space corresponds either to a p-structure potential or to a vibrating membrane with nonquadratic energy. In both cases, the nonlinearity depends on a parameter, so that the impact of the nonlinearity on the efficiency of the splitting procedure can be studied.

The generic problem is posed in the domain  $\Omega \times J$  and consists in finding  $u : \Omega \times J \rightarrow \mathbb{R}$  such that

$$\partial_t^2 u + \nabla \cdot (\mu(u, \nabla u)^2 \nabla u) = f, \quad \text{in } \Omega \times J, \quad (5.1a)$$

$$u|_{t=0} = u_0, \quad \text{in } \Omega, \quad (5.1b)$$

$$\partial_t u|_{t=0} = v_0, \quad \text{in } \Omega, \quad (5.1c)$$

$$u = 0, \quad \text{on } \partial\Omega \times J, \quad (5.1d)$$

with  $f : \Omega \times J \rightarrow \mathbb{R}$  the source term,  $\mu : \mathbb{R} \times \mathbb{R}^d \rightarrow \mathbb{R}_+$  a nonlinear function representing the speed of sound, and  $u_0, v_0$  the initial data for the potential,  $u$ , and the velocity,  $\partial_t u$ . For simplicity, homogeneous Dirichlet boundary conditions are enforced on  $u$ , and we assume that the initial data  $u_0, v_0$  satisfy these conditions.

### 5.1 Splitting procedure for the nonlinear wave equation

In this section, the space discretization of the nonlinear equation (5.1) with the HHO method is presented, along with the adaptation of the splitting procedure introduced in Chapter 3.

Assuming  $f \in L^2(J; L^2(\Omega))$  and that  $\mu(u, \nabla u)$  stays bounded at all times, we seek the solution of (5.1) in  $U := H^2(J; L^2(\Omega)) \cap L^2(J; H_0^1(\Omega))$  such that

$$(\partial_t^2 u, w)_\Omega + (\mu(u, \nabla u)^2 \nabla u, \nabla w)_\Omega = (f, w)_\Omega, \quad \forall t \in J, \forall w \in H_0^1(\Omega). \quad (5.2)$$

#### 5.1.1 Complete discretization of the nonlinear wave equation

**HHO space discretization.** We use the discrete setting from Section 2.2.1 as well as the definition of the gradient reconstruction operator from Section 2.2.2. However, the definition of the stabilization operator differs from the linear case. Indeed, the local stabilization form depends on the nonlinear speed

of sound. Namely, for all  $\hat{y}_T, \hat{v}_T, \hat{w}_T \in \widehat{\mathcal{U}}_T^{l,k}$ , the local stabilization form is now defined as

$$\sigma_T(\hat{y}_T; \hat{v}_T, \hat{w}_T) = \gamma \bar{\mu}_T(\hat{y}_T)^2 \sum_{F \in \mathcal{F}_T} \eta_{TF}^{-1} (S_{TF}(\hat{v}_T), S_{TF}(\hat{w}_T))_F, \quad (5.3)$$

with  $\bar{\mu}_T(\hat{y}_T)$  an approximation of the local speed of sound in  $T$  evaluated using  $\hat{y}_T$ , the scaling factor  $\eta_{TF}$  equal to  $h_F^{-1}$  or  $h_T^{-1}$ , and  $\gamma > 0$  the scaling parameter. This function is nonlinear with respect to  $\hat{y}_T$  and is denoted  $\sigma$  instead of  $s$  to differentiate it from the linear stabilization.

The local stiffness form  $b_T$  also depends on the nonlinear speed of sound and is defined, for all  $\hat{y}_T, \hat{v}_T, \hat{w}_T \in \widehat{\mathcal{U}}_T^{l,k}$ , by

$$b_T(\hat{y}_T; \hat{v}_T, \hat{w}_T) := (\mu(y_T, \mathbf{G}_T^k(\hat{y}_T))^2 \mathbf{G}_T^k(\hat{v}_T), \mathbf{G}_T^k(\hat{w}_T))_T. \quad (5.4)$$

Both forms  $b_T$  and  $\sigma_T$  are nonlinear w.r.t.  $\hat{y}_T$  and linear w.r.t.  $\hat{v}_T$  and  $\hat{w}_T$ . The global forms  $b_h$  and  $\sigma_h$  are defined, for all  $\hat{y}_h, \hat{v}_h, \hat{w}_h \in \widehat{\mathcal{U}}_{h,0}^{l,k}$ , as

$$b_h(\hat{y}_h; \hat{v}_h, \hat{w}_h) := \sum_{T \in \mathcal{T}_h} b_T(\hat{y}_T; \hat{v}_T, \hat{w}_T), \quad \sigma_h(\hat{y}_h; \hat{v}_h, \hat{w}_h) := \sum_{T \in \mathcal{T}_h} \sigma_T(\hat{y}_T; \hat{v}_T, \hat{w}_T). \quad (5.5)$$

The space semi-discrete scheme for the nonlinear wave equation (3.1) consists of finding  $\hat{u}_h := (u_{\mathcal{T}}, u_{\mathcal{F}}) \in C^2(\bar{J}; \widehat{\mathcal{U}}_{h,0}^{l,k})$  such that, for all  $t \in \bar{J}$  and all  $\hat{w}_h := (w_{\mathcal{T}}, w_{\mathcal{F}}) \in \widehat{\mathcal{U}}_{h,0}^{l,k}$ ,

$$(\partial_t^2 u_{\mathcal{T}}, w_{\mathcal{T}})_{\Omega} + b_h(\hat{u}_h(t); \hat{u}_h(t), \hat{w}_h) + \sigma_h(\hat{u}_h(t); \hat{u}_h(t), \hat{w}_h) = (f(t), w_{\mathcal{T}})_{\Omega}. \quad (5.6)$$

**Time scheme.** The fully discrete solution is obtained, using the leapfrog scheme, by solving, for all  $n \geq 1$  and all  $\hat{w}_h \in \widehat{\mathcal{U}}_{h,0}^{l,k}$ ,

$$\frac{1}{\Delta t^2} (u_{\mathcal{T}}^{n+1} - 2u_{\mathcal{T}}^n + u_{\mathcal{T}}^{n-1}, w_{\mathcal{T}})_{\Omega} + b_h(\hat{u}_h^n; \hat{u}_h^n, \hat{w}_h) + \sigma_h(\hat{u}_h^{n-1}; \hat{u}_h^n, \hat{w}_h) = (f(t^n), w_{\mathcal{T}})_{\Omega}, \quad (5.7)$$

with  $u_{\mathcal{T}}^n, u_{\mathcal{T}}^{n-1}$  known from prior time steps or given by the initial conditions as follows:

$$u_{\mathcal{T}}^0 = \Pi_{\mathcal{T}}^l(u_0) \quad (5.8a)$$

$$b_h(\hat{u}_h^0; \hat{u}_h^0, (0, w_{\mathcal{F}})) + \sigma_h(u_{\mathcal{T}}^0; \hat{u}_h^0, (0, w_{\mathcal{F}})) = 0, \quad (5.8b)$$

$$(u_{\mathcal{T}}^1, w_{\mathcal{T}})_{\Omega} = (u_{\mathcal{T}}^0 + \Delta t \Pi_{\mathcal{T}}^l(v_0), w_{\mathcal{T}})_{\Omega} + \frac{\Delta t^2}{2} [(f(0), w_{\mathcal{T}})_{\Omega} - b_h(\hat{u}_h^0; \hat{u}_h^0, (w_{\mathcal{T}}, 0)) - \sigma_h(\hat{u}_h^0; \hat{u}_h^0, (w_{\mathcal{T}}, 0))]. \quad (5.8c)$$

Notice that, in equation (5.7), the stabilization term is linear w.r.t  $\hat{u}_h^n$  and that in the second equation (5.8b), the stabilization is linear with respect to  $\hat{u}_h^0$ . This linearization does not influence the accuracy because the speed of sound in the stabilization only aims at equilibrating the magnitude of the stabilization and stiffness terms. More precisely, the approximation of the local speed of sound in a cell  $T \in \mathcal{T}_h$  is evaluated as

$$\bar{\mu}_T(u_T^{n-1}) := \mu(u_T^{n-1}(\mathbf{x}_T), \mathbf{G}_T^k(\hat{u}_T^{n-1})(\mathbf{x}_T)), \quad (5.9)$$

where  $\mathbf{x}_T$  is the barycenter of  $T$ . If the speed of sound does not vary too much across the mesh and during the simulation, a constant value  $\bar{\mu}$  can be taken in all the cells and at all the time steps. For the stiffness operator  $b_h$ , the above linearization is not considered in order to preserve the second-order convergence rate in time of the discrete scheme.

As in the linear case, the scheme (5.7) with the initial conditions (5.8) can be rewritten in terms of the speed,  $v_{\mathcal{T}}^n$ , and the acceleration,  $a_{\mathcal{T}}^n$ , defined in (3.31):

$$(a_{\mathcal{T}}^n, w_{\mathcal{T}})_{\Omega} + b_h(\hat{u}_h^n; \hat{u}_h^n, \hat{w}_h) + \sigma_h(\hat{u}_h^{n-1}; \hat{u}_h^n, \hat{w}_h) = (f(t^n), w_{\mathcal{T}})_{\Omega}, \quad (5.10)$$

The fully discrete problem (5.7) can be written in algebraic form by considering the vector-valued non-linear stiffness operator  $\mathbf{B}(\mathbf{U}_{\mathcal{T}}^n, \mathbf{U}_{\mathcal{F}}^n)$  associated with the linear form  $b_h(\hat{u}_h^n; \hat{u}_h^n, \cdot)$ , with  $\mathbf{B}_{\mathcal{T}}, \mathbf{B}_{\mathcal{F}}$  respectively collecting the cell and face components. The linearized stabilization bilinear form  $\sigma_h(\hat{u}_h^{n-1}; \cdot, \cdot)$  leads to a symmetric matrix depending on  $\hat{u}_h^{n-1}$  and denoted by  $\mathcal{S}^{n-1}$ . Altogether, the algebraic formulation reads

$$\frac{1}{\Delta t^2} \begin{bmatrix} \mathcal{M} & 0 \\ 0 & 0 \end{bmatrix} \begin{pmatrix} \mathbf{U}^{n+1} - 2\mathbf{U}^n + \mathbf{U}^{n-1} \\ \cdot \end{pmatrix} + \begin{pmatrix} \mathbf{B}_{\mathcal{T}}(\mathbf{U}_{\mathcal{T}}^n, \mathbf{U}_{\mathcal{F}}^n) \\ \mathbf{B}_{\mathcal{F}}(\mathbf{U}_{\mathcal{T}}^n, \mathbf{U}_{\mathcal{F}}^n) \end{pmatrix} + \begin{bmatrix} \mathcal{S}_{\mathcal{T}\mathcal{T}}^{n-1} & \mathcal{S}_{\mathcal{T}\mathcal{F}}^{n-1} \\ \mathcal{S}_{\mathcal{F}\mathcal{T}}^{n-1} & \mathcal{S}_{\mathcal{F}\mathcal{F}}^{n-1} \end{bmatrix} \begin{pmatrix} \mathbf{U}_{\mathcal{T}}^n \\ \mathbf{U}_{\mathcal{F}}^n \end{pmatrix} = \begin{bmatrix} \mathbf{F}_{\mathcal{T}}^n \\ 0 \end{bmatrix}, \quad (5.11)$$

with  $\mathcal{M}$  the cell mass matrix. As a consequence of the structure of the global mass matrix, the face component is replaced by a “.” in the acceleration term. The submatrix  $\mathcal{S}_{\mathcal{T}\mathcal{T}}^{n-1}$  is block-diagonal, since  $\sigma_h(\hat{u}_h^{n-1}; \cdot, \cdot)$  does not couple cell degrees of freedom from different cells. It is useful to notice that in the mixed-order setting,  $\mathcal{S}_{\mathcal{F}\mathcal{F}}^{n-1}$  is still block-diagonal, since it still reduces on each face to the mass matrix associated with a local polynomial basis. Finally, as in the linear case, the cell mass matrix  $\mathcal{M}$  is also block-diagonal.

The resolution of the wave equation discretized with HHO and the leapfrog scheme proceeds in two steps for all  $n \geq 1$ :

1. Compute  $\mathbf{U}_{\mathcal{F}}^n$ , i.e., solve  $\mathbf{B}_{\mathcal{F}}(\mathbf{U}_{\mathcal{T}}^n, \mathbf{U}_{\mathcal{F}}^n) + \mathcal{S}_{\mathcal{F}\mathcal{F}}^{n-1} \mathbf{U}_{\mathcal{F}}^n = -\mathcal{S}_{\mathcal{F}\mathcal{T}}^{n-1} \mathbf{U}_{\mathcal{T}}^n$  using the second row in (5.11), where  $\mathbf{U}_{\mathcal{T}}^n$  is the data (known from the previous time step or the initial condition) and  $\mathbf{U}_{\mathcal{F}}^n$  the unknown;
2. Compute  $\mathbf{U}_{\mathcal{T}}^{n+1}$ , i.e., solve  $\frac{1}{\Delta t^2} \mathcal{M} \mathbf{U}_{\mathcal{T}}^{n+1} = \mathbf{F}_{\mathcal{T}}^n - \frac{1}{\Delta t^2} \mathcal{M} (\mathbf{U}_{\mathcal{T}}^{n-1} - 2\mathbf{U}_{\mathcal{T}}^n) - \mathbf{B}_{\mathcal{T}}(\mathbf{U}_{\mathcal{T}}^n, \mathbf{U}_{\mathcal{F}}^n) - \mathcal{S}_{\mathcal{T}\mathcal{T}}^{n-1} \mathbf{U}_{\mathcal{T}}^n - \mathcal{S}_{\mathcal{T}\mathcal{F}}^{n-1} \mathbf{U}_{\mathcal{F}}^n$  using the first row in (5.11).

### 5.1.2 Explicit time discretization with the leapfrog scheme

**The operator splitting in the mixed order setting.** As in the linear case, the mixed-order setting offers a convenient operator splitting, which consists, for  $m \geq 0$ , in finding  $u_{\mathcal{F}}^{n,m+1}$  such that

$$\sigma_h(\hat{u}_h^{n-1}; (0, u_{\mathcal{F}}^{n,m+1}), (0, w_{\mathcal{F}})) = -b_h((u_{\mathcal{T}}^n, u_{\mathcal{F}}^{n,m}); (u_{\mathcal{T}}^n, u_{\mathcal{F}}^{n,m}), (0, w_{\mathcal{F}})) - \sigma_h(\hat{u}_h^{n-1}; (u_{\mathcal{T}}^n, 0), (0, w_{\mathcal{F}})), \quad (5.12)$$

for all  $w_{\mathcal{F}} \in \mathcal{U}_{\mathcal{F},0}^k$ . The algebraic form is as follows: Setting  $\mathbf{U}_{\mathcal{F}}^{n,0} := \mathbf{U}_{\mathcal{F}}^{n-1}$ , we seek  $\mathbf{U}_{\mathcal{F}}^{n,m+1} \in \mathbb{R}^{N_{\mathcal{F}}}$  such that

$$\mathcal{S}_{\mathcal{F}\mathcal{F}}^{n-1} \mathbf{U}_{\mathcal{F}}^{n,m+1} = -\mathbf{B}_{\mathcal{F}}(\mathbf{U}_{\mathcal{T}}^n, \mathbf{U}_{\mathcal{F}}^{n,m}) - \mathcal{S}_{\mathcal{F}\mathcal{T}}^{n-1} \mathbf{U}_{\mathcal{T}}^n. \quad (5.13)$$

This splitting procedure is computationally effective since, as mentioned above, the face-face stabilization submatrix  $\mathcal{S}_{\mathcal{F}\mathcal{F}}^{n-1}$  is block-diagonal. The procedure is a fixed-point algorithm, so that its convergence is ensured for  $\delta < 1$ , where  $\delta$  is the Lipschitz constant of the vector-valued function  $(\mathcal{S}_{\mathcal{F}\mathcal{F}}^{n-1})^{-1} \mathbf{B}_{\mathcal{F}}(\mathbf{U}_{\mathcal{T}}^n, \cdot)$ , i.e.,

$$\|(\mathcal{S}_{\mathcal{F}\mathcal{F}}^{n-1})^{-1} (\mathbf{B}_{\mathcal{F}}(\mathbf{U}_{\mathcal{T}}^n, X) - \mathbf{B}_{\mathcal{F}}(\mathbf{U}_{\mathcal{T}}^n, Y))\| \leq \delta \|X - Y\| \quad \forall X, Y \in \mathbb{R}^{N_{\mathcal{F}}}. \quad (5.14)$$

**The operator splitting in the equal order setting.** This case requires some additional definitions, as in the linear case. Specifically, let  $\zeta_T$  be the local form such that for all  $T \in \mathcal{T}_h$ , and all  $\hat{y}_T \in \hat{\mathcal{U}}_T^{l,k}, v_{\partial T}, w_{\partial T} \in \mathcal{U}_{\partial T}^k$ ,

$$\begin{aligned} \zeta_T(\hat{y}_T; (0, v_{\partial T}), (0, w_{\partial T})) &= \gamma \mu_T (\hat{y}_T)^2 \sum_{F \in \mathcal{F}_T} \eta_{TF}^{-1} \left\{ ((I - \Pi_T^k) R_T^{k+1}(0, v_{\partial T})|_F, v_F)_F \right. \\ &\quad + (v_F, (I - \Pi_T^k) R_T^{k+1}(0, w_{\partial T})|_F)_F \\ &\quad \left. + (\Pi_F^k (I - \Pi_T^k) R_T^{k+1}(0, v_{\partial T}), \Pi_F^k (I - \Pi_T^k) R_T^{k+1}(0, w_{\partial T})|_F)_F \right\}, \end{aligned} \quad (5.15)$$

and let  $\sigma_T^*$  be the local form defined by

$$\sigma_T^*(\hat{y}_T; (0, v_{\partial T}), (0, w_{\partial T})) := \gamma \mu_T (\hat{y}_T)^2 \sum_{F \in \mathcal{F}_T} \eta_{TF}^{-1} (v_F, w_F)_F. \quad (5.16)$$

Then the equal-order stabilization form writes

$$\sigma_T(\hat{y}_T; (0, v_{\partial T}), (0, w_{\partial T})) = \sigma_T^*(\hat{y}_T; (0, v_{\partial T}), (0, w_{\partial T})) + \zeta_T(\hat{y}_T; (0, v_{\partial T}), (0, w_{\partial T})). \quad (5.17)$$

Let us introduce the global forms  $\sigma_h^*(\hat{y}_h; (0, v_{\mathcal{F}}), (0, w_{\mathcal{F}})) := \sum_{T \in \mathcal{T}_h} \sigma_T^*(\hat{y}_T; (0, v_{\partial T}), (0, w_{\partial T}))$  and  $\zeta_h(\hat{y}_h; (0, v_{\mathcal{F}}), (0, w_{\mathcal{F}})) := \sum_{T \in \mathcal{T}_h} \zeta_T(\hat{y}_T; (0, v_{\partial T}), (0, w_{\partial T}))$ , so that  $\sigma_h = \sigma_h^* + \zeta_h$ . This leads to the following iterative procedure, with the same initial condition as for the mixed-order setting: For all  $m \geq 0$ , find  $u_{\mathcal{F}}^{n,m+1} \in \mathcal{U}_{\mathcal{F},0}^k$  such that

$$\begin{aligned} \sigma_h^*(\hat{u}_h^{n-1}; (0, u_{\mathcal{F}}^{n,m+1}), (0, w_{\mathcal{F}})) &= -b_h((u_{\mathcal{T}}^n, u_{\mathcal{F}}^{n,m}); (u_{\mathcal{T}}^n, u_{\mathcal{F}}^{n,m}), (0, w_{\mathcal{F}})) \\ &\quad - \zeta_h(\hat{u}_h^{n-1}; (0, u_{\mathcal{F}}^{n,m}), (0, w_{\mathcal{F}})) - \sigma_h(\hat{u}_h^{n-1}; (u_{\mathcal{T}}^n, 0), (0, w_{\mathcal{F}})), \end{aligned} \quad (5.18)$$

for all  $w_{\mathcal{F}} \in \mathcal{U}_{\mathcal{F},0}^k$ . At the algebraic level, we define two matrices  $\mathcal{S}_{\mathcal{F}\mathcal{F}}^{*,n-1}$  and  $\mathcal{Z}_{\mathcal{F}\mathcal{F}}^{n-1}$  such that  $\mathcal{S}_{\mathcal{F}\mathcal{F}}^{n-1} = \mathcal{S}_{\mathcal{F}\mathcal{F}}^{*,n-1} + \mathcal{Z}_{\mathcal{F}\mathcal{F}}^{n-1}$ ,  $\mathcal{Z}_{\mathcal{F}\mathcal{F}}^{n-1}$  corresponds to the bilinear form  $\zeta_h(\hat{u}_h^{n-1}; \cdot, \cdot)$  and  $\mathcal{S}_{\mathcal{F}\mathcal{F}}^{*,n-1}$  to  $\sigma_h^*(\hat{u}_h^{n-1}; \cdot, \cdot)$ . Then the splitting procedure translates into the following iterative algorithm: For all  $m \geq 0$ , find  $U_{\mathcal{F}}^{n,m+1} \in \mathbb{R}^{N_{\mathcal{F}}}$  such that

$$\mathcal{S}_{\mathcal{F}\mathcal{F}}^{*,n-1} U_{\mathcal{F}}^{n,m+1} = -\mathbf{B}_F(U_{\mathcal{T}}^n, U_{\mathcal{F}}^{n,m}) - \mathcal{Z}_{\mathcal{F}\mathcal{F}}^{n-1} U_{\mathcal{F}}^{n,m} - \mathcal{S}_{\mathcal{F}\mathcal{T}}^{n-1} U_{\mathcal{T}}^n. \quad (5.19)$$

As for the mixed-order setting, the convergence condition is that  $\delta < 1$ , where  $\delta$  is the Lipschitz constant of the vector-valued function  $(\mathcal{S}_{\mathcal{F}\mathcal{F}}^{*,n-1})^{-1}(\mathcal{Z}_{\mathcal{F}\mathcal{F}}^{n-1}(\cdot) + \mathbf{B}_F(U_{\mathcal{T}}^n, \cdot))$ , i.e.,

$$\|(\mathcal{S}_{\mathcal{F}\mathcal{F}}^{*,n-1})^{-1}(\mathcal{Z}_{\mathcal{F}\mathcal{F}}^{n-1}(X - Y) + \mathbf{B}_F(U_{\mathcal{T}}^n, X) - \mathbf{B}_F(U_{\mathcal{T}}^n, Y))\| \leq \delta \|X - Y\|, \quad \forall X, Y \in \mathbb{R}^{N_{\mathcal{F}}}. \quad (5.20)$$

**Simplifying the splitting.** This splitting algorithm can be made more efficient, by neglecting the dependance of the speed of sound to the time step in the stabilization. Indeed, as mentioned earlier, if the value of  $\mu$  has the same order of magnitude in the entire domain, and does not vary strongly during the simulation, a unique value  $\bar{\mu}$  can be taken for all cells  $T \in \mathcal{T}_h$  and all time steps  $t^n$ . This simplifies the definition of the local stabilization operator  $\sigma_T$  from equation (5.3) into

$$s_T(\hat{v}_T, \hat{w}_T) = \gamma \bar{\mu}^2 \sum_{F \in \mathcal{F}_T} \eta_{TF}^{-1} (S_{TF}(\hat{v}_T), S_{TF}(\hat{w}_T))_F, \quad (5.21)$$

for all  $\hat{v}_T, \hat{w}_T \in \hat{\mathcal{U}}_T^{l,k}$ . The letter  $s$  is used again in this case, since this stabilization is the same linear operator as in the previous chapter. The functional expression of the splitting becomes, in the mixed order setting, for  $m \geq 0$ ,

$$s_h((0, u_{\mathcal{F}}^{n,m+1}), (0, w_{\mathcal{F}})) = -b_h((u_{\mathcal{T}}^n, u_{\mathcal{F}}^{n,m}); (u_{\mathcal{T}}^n, u_{\mathcal{F}}^{n,m}), (0, w_{\mathcal{F}})) - s_h((u_{\mathcal{T}}^n, 0), (0, w_{\mathcal{F}})), \quad (5.22)$$

and in the equal-order setting (with  $s_h^*$  and  $\zeta_h$  defined as in the linear case),

$$s_h^*((0, u_{\mathcal{F}}^{n,m+1}), (0, w_{\mathcal{F}})) = -b_h((u_{\mathcal{T}}^n, u_{\mathcal{F}}^{n,m}); (u_{\mathcal{T}}^n, u_{\mathcal{F}}^{n,m}), (0, w_{\mathcal{F}})) - \zeta_h((0, u_{\mathcal{F}}^{n,m}), (0, w_{\mathcal{F}})) - s_h((u_{\mathcal{T}}^n, 0), (0, w_{\mathcal{F}})), \quad (5.23)$$

for all  $w_{\mathcal{F}} \in \mathcal{U}_{\mathcal{F},0}^k$ . The algebraic expressions (5.13) in the mixed-order setting and (5.19) in the equal-order setting become, respectively, for  $m \geq 0$

$$\mathcal{S}_{\mathcal{F}\mathcal{F}} U_{\mathcal{F}}^{n,m+1} = -\mathbf{B}_F(U_{\mathcal{T}}^n, U_{\mathcal{F}}^{n,m}) - \mathcal{S}_{\mathcal{F}\mathcal{T}} U_{\mathcal{T}}^n \quad (5.24)$$

and

$$\mathcal{S}_{\mathcal{F}\mathcal{F}}^* U_{\mathcal{F}}^{n,m+1} = -\mathbf{B}_F(U_{\mathcal{T}}^n, U_{\mathcal{F}}^{n,m}) - \mathcal{Z}_{\mathcal{F}\mathcal{F}} U_{\mathcal{F}}^{n,m} - \mathcal{S}_{\mathcal{F}\mathcal{T}} U_{\mathcal{T}}^n. \quad (5.25)$$

This simplification allows to precompute the stabilization matrix and preinvert the block-diagonal iteration matrix  $\mathcal{S}_{\mathcal{F}\mathcal{F}}$  in the mixed order setting and  $\mathcal{S}_{\mathcal{F}\mathcal{F}}^*$  in the equal-order setting. This replaces an iterative solver by only local matrix-vector operations.



## 5.2 Nonlinear wave equation with a p-structure potential

In this section, we study the splitting algorithm on a nonlinear wave equation where the diffusion operator in space has a so-called p-structure:

$$\partial_t^2 u + \nabla \cdot ((\mu_0^2 + |\nabla u|^2)^{\frac{p-2}{2}} \nabla u) = f, \quad \text{in } \Omega \subset \mathbb{R}^2, \forall t \in J, \quad (5.26)$$

with  $\mu_0 \in \mathbb{R}_+$  and  $p \in (1, +\infty)$ . For the value  $p = 2$ , (5.26) corresponds to the linear wave equation (3.1). The associated Hamiltonian writes  $\mathcal{H}^p(g) := \frac{1}{p}(\mu_0^2 + |g|^2)^{\frac{p}{2}}$  for all  $g \in \mathbb{R}^2$ , and the nonlinear wave equation (5.26) can be rewritten as  $\partial_t^2 u + \nabla \cdot (\nabla_g \mathcal{H}^p(\nabla u)) = f$ . The local nonlinear stiffness form  $b_T$  on a cell  $T$  now writes

$$b_T(\hat{y}_T; \hat{v}_T, \hat{w}_T) := ((\mu_0^2 + |\mathbf{G}_T^k(\hat{y}_T)|^2)^{\frac{p-2}{2}} \mathbf{G}_T^k(\hat{v}_T), \mathbf{G}_T^k(\hat{w}_T)). \quad (5.27)$$

### 5.2.1 Setting

We set  $\Omega := (0, 1)^2$ ,  $\Theta := 0.8$ ,  $f := 0$  and we enforce homogeneous Dirichlet boundary conditions, null initial condition for  $u$ , and the initial velocity  $v_0(x, y) := 5 \sin(\pi x) \sin(\pi y)$ . In the linear case ( $p = 2$ ), the exact solution is  $u(x, y, t) = \frac{5}{\sqrt{2\pi}} \sin(\sqrt{2}\pi t) \sin(\pi x) \sin(\pi y)$ . When  $p \neq 2$ , the exact solution is not known and the error is computed using a reference solution (obtained with high polynomial orders (4, 3) and a fine mesh with  $h \approx 0.0003$ ). If  $\mu_0 = 0$ , the diffusion operator essentially behaves as the p-Laplace operator. On the contrary, if  $\mu_0$  is much larger than  $|\nabla u|$ , the problem is nearly linear and the behavior does not differ too much from a linear diffusion operator. In order to avoid such cases, the value of  $\mu_0$  is chosen so that it is not significantly smaller than  $\|\nabla u\|_{L^\infty(\Omega \times J)}$ . Since  $\|\nabla u\|_{L^\infty(\Omega \times J)}$  depends on the value of  $\mu_0$ , this can be asserted by trial and error. Here, we take  $\mu_0^2 := 0.5$ . Moreover, to scale the stabilization, we take the same value  $\bar{\mu}_T^2 := 5$  in each cell and at each time step. Based on our simulations (see Figure 5.1) which cover the range  $p \in [1.25, 5]$ , this appears to be a reasonable choice for scaling the stabilization.

In our numerical experiments, we consider an unstructured triangular mesh sequence, obtained by repeatedly refining an initial coarse mesh, thus preserving the same regularity in the entire mesh sequence. To compute the errors, we consider ten sensors placed in the triangle  $\{x \in [0, 0.5], 0 \leq y \leq x\}$ , and compute the value of the solution,  $u$ , the velocity,  $\partial_t u$ , the acceleration,  $\partial_t^2 u$ , and of both gradient components,  $\partial_x u, \partial_y u$ , over the simulation at these ten sensors. The sensors are placed only in the above triangle by symmetry arguments. In Figure 5.1, the solution at three sensors is shown for seven values of  $p$ . The values for which  $p \neq 2$  can be regrouped into pairs of conjugated values  $(p, p')$  such that  $\frac{1}{p} + \frac{1}{p'} = 1$ . Specifically, we choose  $p \in \{2.5, 3, 5\}$  and the conjugate values  $p' \in \{1.67, 1.5, 1.25\}$ . The larger  $|p - 2|$ , the harder the simulation, and the more refined the mesh needs to be in order to capture the nonlinear variations.

In this entire subsection, for the splitting and the semi-implicit scheme, we use the convergence criterion  $\epsilon = 10^{-10}$ .

### 5.2.2 Splitting efficiency for $p = 3$

We first make computations with various values of the stability parameter  $\gamma$  and  $p = 3$ . The minimal value  $\gamma^*$  leading to a converging splitting depends on the polynomial order and on the value of  $p$ , and cannot be deduced from the linear study. Table 5.1 reports the smallest value of  $\gamma$  observed numerically, denoted by  $\gamma^{\min}$  in what follows. Since this value is an estimation, it is rounded to the smallest larger integer.

**Lowest order: splitting vs semi-implicit.** The objective of the splitting being a reduction of computational costs, we first compare execution times between the splitting procedure and the semi-implicit

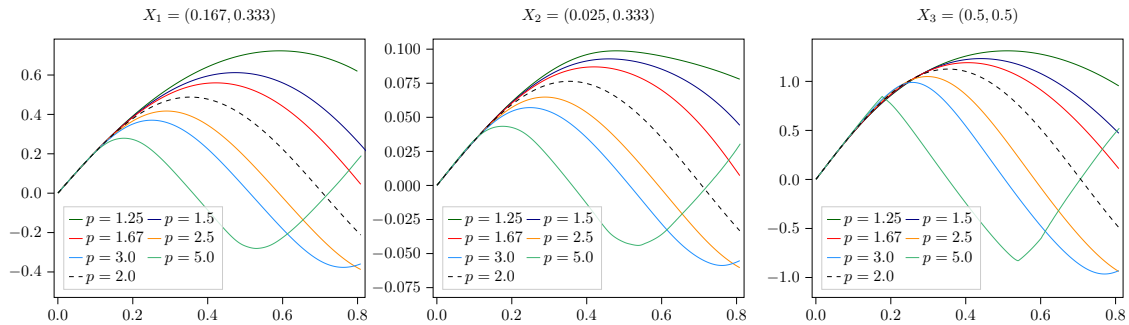


Figure 5.1: Nonlinear ( $p$ -structure) wave problem with  $p \in \{1.25, 1.5, 1.67, 2, 2.5, 3, 5\}$  - Solution at sensors placed at  $(0.167, 0.333)$ ,  $(0.025, 0.333)$  and  $(0.5, 0.5)$  on the time interval  $[0, 0.8]$ .

	(0,0)	(1,1)	(2,2)	(3,3)	(1,0)	(2,1)	(3,2)	(4,3)
$p = 3$	2	10	20	25	5	20	20	25

Table 5.1: Nonlinear ( $p$ -structure) wave problem with  $p = 3$  - Minimal value  $\gamma^{\min}$  for the splitting procedure to converge,  $k \in \{0, 1, 2, 3\}$ , equal- and mixed-order settings.

scheme (based on a Newton solver). Computations are run sequentially and for the lowest polynomial order only. The Newton linear system is solved with a direct solver since experiments with iterative solvers yield similar execution times. The time step is chosen so that the time discretization error and the space discretization error are equilibrated. For all values of  $\gamma$ , this leads to an equilibrated time step that is smaller than the critical time step. Thus, the execution time is not affected by the reduction of the critical time step for larger  $\gamma$ . On the one hand, Figure 5.2 reports the maximum of the discrete energy error, evaluated using the reference solution as the maximum value over the ten sensor points and all the time steps. On the other hand, Figure 5.2 also reports the mean number of splitting and Newton iterations per time step for each value of  $\gamma$  and each mesh. The results in Figure 5.2 deals with the lowest-order cases  $(0, 0)$  and  $(1, 0)$ . The first salient observation is that, as for the linear case, the behavior differs between the equal- and mixed-order settings. In the equal-order setting, the number of iterations is higher when  $\gamma$  increases, whereas there is an optimal value in the mixed-order setting, which is not  $\gamma^*$ . In the equal-order setting, the splitting procedure takes more time than the semi-implicit scheme on the coarsest meshes, but is (much) faster on the finest meshes. The gain in execution time (i.e. the ratio of the semi-implicit time over the splitting time) is reported in Table 5.2. In the mixed-order setting, the splitting procedure is always (much) faster than the semi-implicit scheme, and the gain in execution time increases as the mesh is refined. The number of iterations for the splitting and the semi-implicit scheme decreases when the mesh is refined. This is due to the time step being also refined to satisfy the requirement on error balancing. Hence, the static problem on the face unknowns becomes easier to solve. On the four coarsest meshes, in the equal-order setting, the number of splitting iterations is very large (more than 100 on the coarsest mesh). This explains the not so good performance of the splitting procedure on the coarsest meshes (see Figure 5.2a). In the mixed-order setting, the number of iterations is always quite low (less than 10), which explains the better performances observed in Figure 5.2b. Since the number of iterations for all the values of  $\gamma$  becomes small on the finest meshes, the value of  $\gamma$  does not impact the execution time in this case.

**Higher-orders: splitting versus semi-implicit.** The above experiments in the linear case and in the lowest-order nonlinear case show that the error does not depend on  $\gamma$  as long as this parameter remains in the range  $[1, 150]$ . Thus, for higher polynomial degrees, we only focus on the execution time. Figure 5.3 displays the average execution time of a single time step for each method and each polynomial order on the same family of unstructured triangular meshes as above. Here, we compare the sequential execution

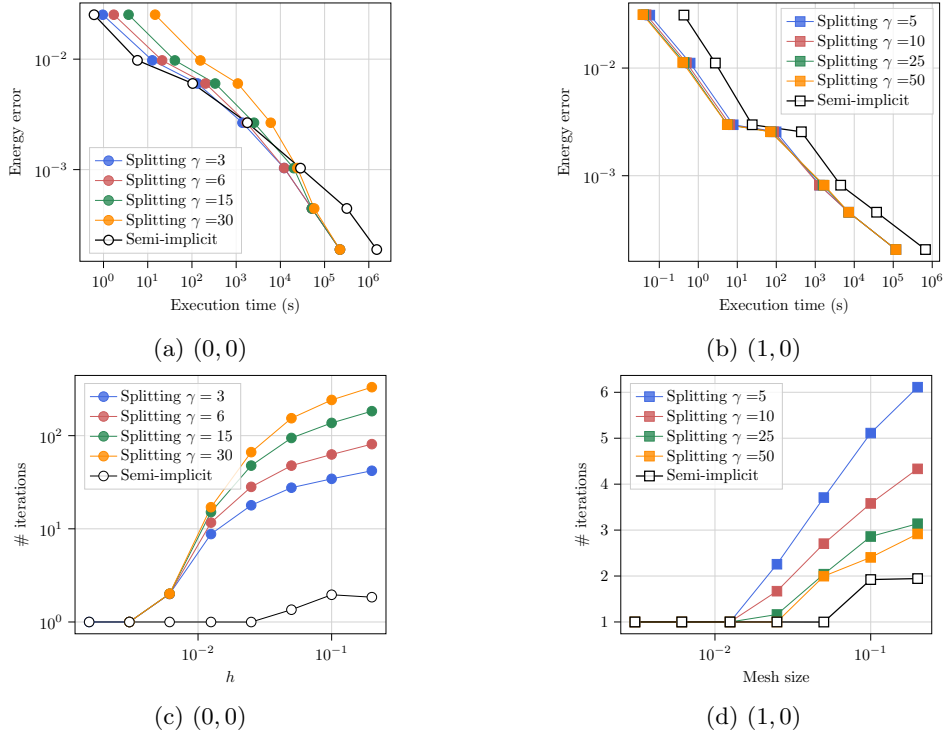


Figure 5.2: Nonlinear (p-structure) wave problem with  $p = 3$  - Energy error as a function of execution time (top) and mean number of iterations per time step as a function of mesh size (bottom),  $\gamma \in \{\gamma^{\min}, 2\gamma^{\min}, 5\gamma^{\min}, 10\gamma^{\min}\}$ , polynomial orders (0,0) and (1,0).

Mesh size $h$	0.0063	0.0031	0.0016
(0,0)	2.34	6.19	6.80
(1,0)	3.44	5.19	5.74

Table 5.2: Nonlinear (p-structure) wave problem with  $p = 3$  - Gain in execution time of the splitting procedure over the semi-implicit scheme on the three finest meshes,  $\gamma = \gamma^{\min}$ , polynomial orders (0,0) and (1,0).

times for the splitting procedure with  $\gamma \in \{\gamma^{\min}, 2\gamma^{\min}, 5\gamma^{\min}\}$  and the semi-implicit scheme with either a direct solver or a GMRES iterative solver with a Jacobi preconditioner. The gain in execution time is reported in Table 5.3 only for  $\gamma = \gamma^{\min}$  since the value of  $\gamma$  does not impact strongly the execution time on the considered meshes. The gain is given on the four finest meshes, exception made for the polynomial order  $k = 3$ , for which the semi-implicit scheme on the finest mesh is too expensive with our current implementation. As expected, the equal- and mixed-order settings behave differently on the coarse meshes. In the equal-order setting, the splitting procedure takes more time on the coarse meshes, but this execution time is reduced as the number of iterations diminishes, so that the splitting procedure becomes computationally efficient for mesh sizes smaller than  $h \approx 0.006$ . In the mixed-order setting, the execution time of the splitting procedure is always smaller or equivalent to the execution time of the semi-implicit scheme, but the gain (ratio of execution times) also increases as the mesh is refined. The gain on the four finest meshes is larger than 2 and reaches 6 for polynomial orders (1,1) and (2,2) when a direct solver is employed in the semi-implicit scheme.

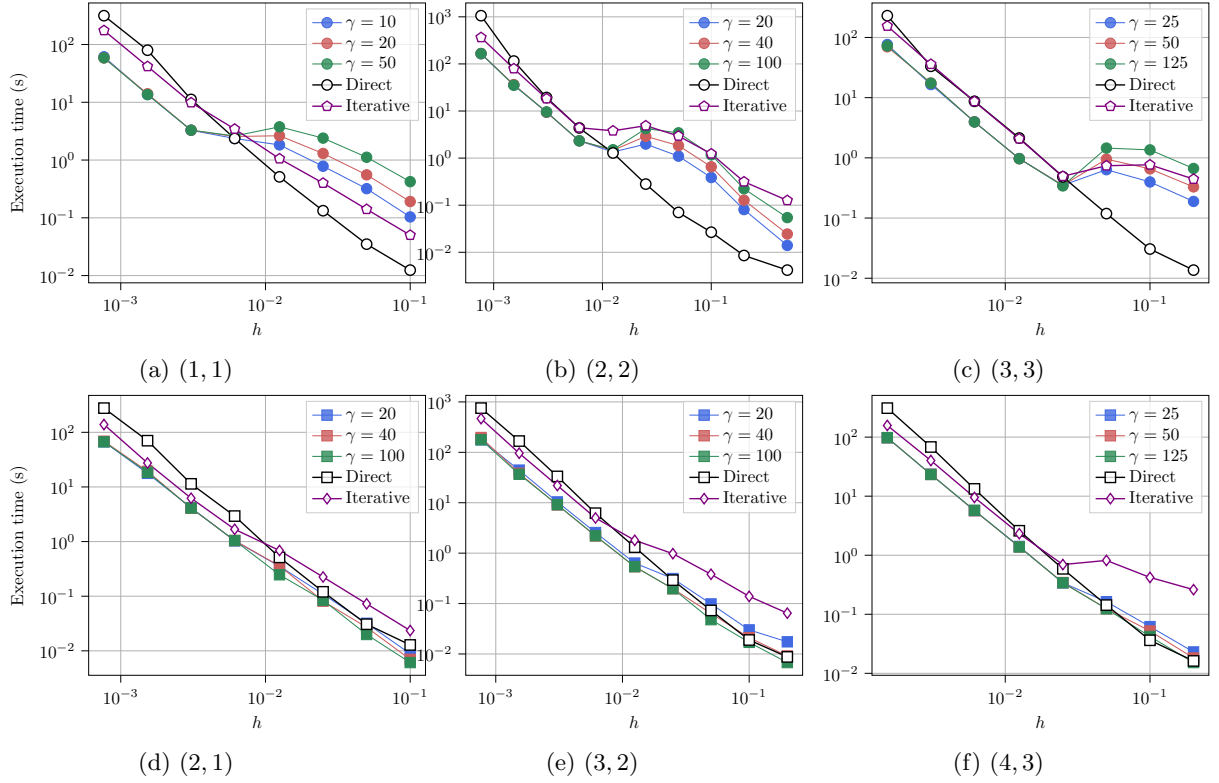


Figure 5.3: Nonlinear (p-structure) wave problem with  $p = 3$  - Comparison of execution times of a single time step with splitting procedure and semi-implicit scheme (direct or iterative linear solvers),  $\gamma \in \{\gamma^{\min}, 2\gamma^{\min}, 5\gamma^{\min}\}$ ,  $k \in \{1, 2, 3\}$  equal- and mixed-order setting.

	(1,1)		(2,1)		(2,2)		(3,2)		(3,3)		(4,3)	
Mesh size	direct	gmres	direct	gmres	direct	gmres	direct	gmres	direct	gmres	direct	gmres
0.00613	1.01	1.46	2.84	1.62	1.87	1.89	2.42	1.94	2.19	2.13	2.33	1.60
0.00306	3.41	3.01	2.73	1.49	2.04	1.94	3.21	2.11	2.02	2.13	2.90	1.71
0.00156	5.84	3.06	3.97	1.54	3.24	2.21	3.77	2.17	3.01	2.17	3.16	1.67
0.00078	5.15	2.85	4.13	2.07	6.43	2.23	4.25	2.63	NA	NA	NA	NA

Table 5.3: Nonlinear (p-structure) wave problem with  $p = 3$  - Gain in execution time of the splitting procedure over the semi-implicit scheme on the four finest meshes,  $\gamma = \gamma^{\min}$ ,  $k \in \{1, 2, 3\}$ , equal- and mixed-order settings.

**Comparison between orders and methods.** We now compare the discrete energy error as a function of the execution time for various polynomial orders using either the HHO method or conforming finite elements for space discretization. We consider  $k \in \{0, 1, 2\}$  for HHO in both mixed- and equal-order settings, and  $k \in \{1, 2\}$  for finite elements. Based on the above results, we only consider the splitting procedure for HHO simulations, and we choose the best value of  $\gamma$  in terms of execution time, i.e.  $\gamma = \gamma^{\min}$  in the equal-order setting and  $\gamma = 5\gamma^{\min}$  in the mixed-order setting. Figure 5.4 first highlights the same behavior as in the linear case when comparing the different HHO orders: higher orders and coarse meshes are more computationally efficient than lower orders on finer meshes. Another salient point is the comparison between mixed- and equal-order settings. Mixed-order settings turns out to be (much) more efficient. For instance, for the same error, the splitting procedure with polynomial orders (3,2) is ten times faster than the splitting procedure with polynomial orders (2,2). Furthermore, P1 finite elements

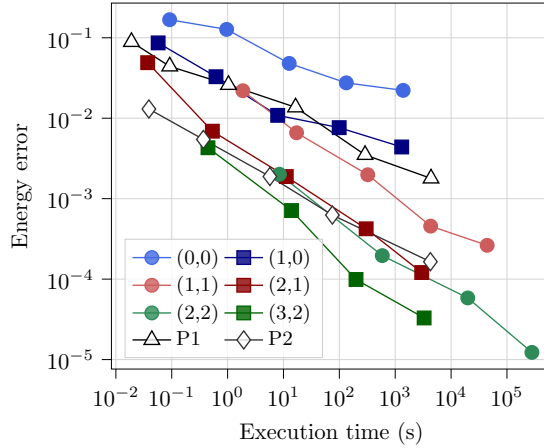


Figure 5.4: Nonlinear (p-structure) wave problem with  $p = 3$  - Discrete energy error as a function of execution time, HHO with  $(k, \gamma) \in \{(0, 2), (1, 10), (2, 20)\}$  in the equal-order setting and  $(k, \gamma) \in \{(0, 25), (1, 100), (2, 100)\}$  in the mixed-order setting, P1 and P2 finite elements.

are more efficient than the lowest equal-order HHO method, broadly equivalent to the lowest mixed-order HHO method, and less efficient than the higher-order HHO method. The efficiency of P2 finite elements is between the first- and second-order HHO settings and quite close to HHO with polynomial orders (2,1). Thus, for a given target error or time budget, high-order HHO with splitting remains the most effective option.

### 5.2.3 Impact of the nonlinearity on the efficiency of the splitting.

The above experiments have so far focused on the case  $p = 3$  corresponding to a mild nonlinearities. We now investigate other values of  $p$ . The same computations as in Figure 5.3 are performed, but only for  $k = 1$ . Other polynomial orders yield similar results and are not shown for brevity. The value of  $\gamma^{\min}$  for each value of  $p$  and both polynomial orders (1,1) and (2,1) are given by table 5.4. The considered values of  $\gamma$  are  $\{\gamma^{\min}, 2\gamma^{\min}, 5\gamma^{\min}\}$ , unless they are larger than 150. This limit is set to avoid numerical errors induced by larger values of  $\gamma$ . The computational setting is otherwise unchanged. Results are reported in Figure 5.5 and the gain in execution time for the three finest meshes in table 5.5 (the values for  $p = 3$  are already given by table 5.3 and are not repeated for brevity). The conclusions are similar to the mildly

	1.25	1.5	1.67	2.5	3	5
(1,1)	3	3	3	3	10	40
(2,1)	10	10	10	10	20	70

Table 5.4: Nonlinear (p-structure) wave problem with  $p \in \{1.25, 1.5, 1.67, 2.5, 3, 5\}$  - Minimal value  $\gamma^{\min}$  for the splitting procedure to converge,  $k = 1$ , equal- and mixed-order settings.

nonlinear case  $p = 3$  indicating that the variation in nonlinearity as quantified by  $p$  does not impact the performance of the splitting procedure in the equal-order setting. In the mixed-order setting,  $p < 2$  seems to be more favorable to the splitting procedure, whereas with  $p > 2$ , the gain in execution time is less pronounced.

In conclusion, this p-structure test case highlights the benefits of the splitting procedure on fine meshes compared to the semi-implicit scheme. The results also indicate that, when working on fine meshes, the choice of  $\gamma$  is not sensitive. Indeed, as long as  $\gamma > \gamma^{\min}$  and in a reasonable range, here  $[\gamma^{\min}, 150]$ , the error and the execution time are not affected by the choice of  $\gamma$ . Since  $\gamma^{\min}$  can be computed on coarse

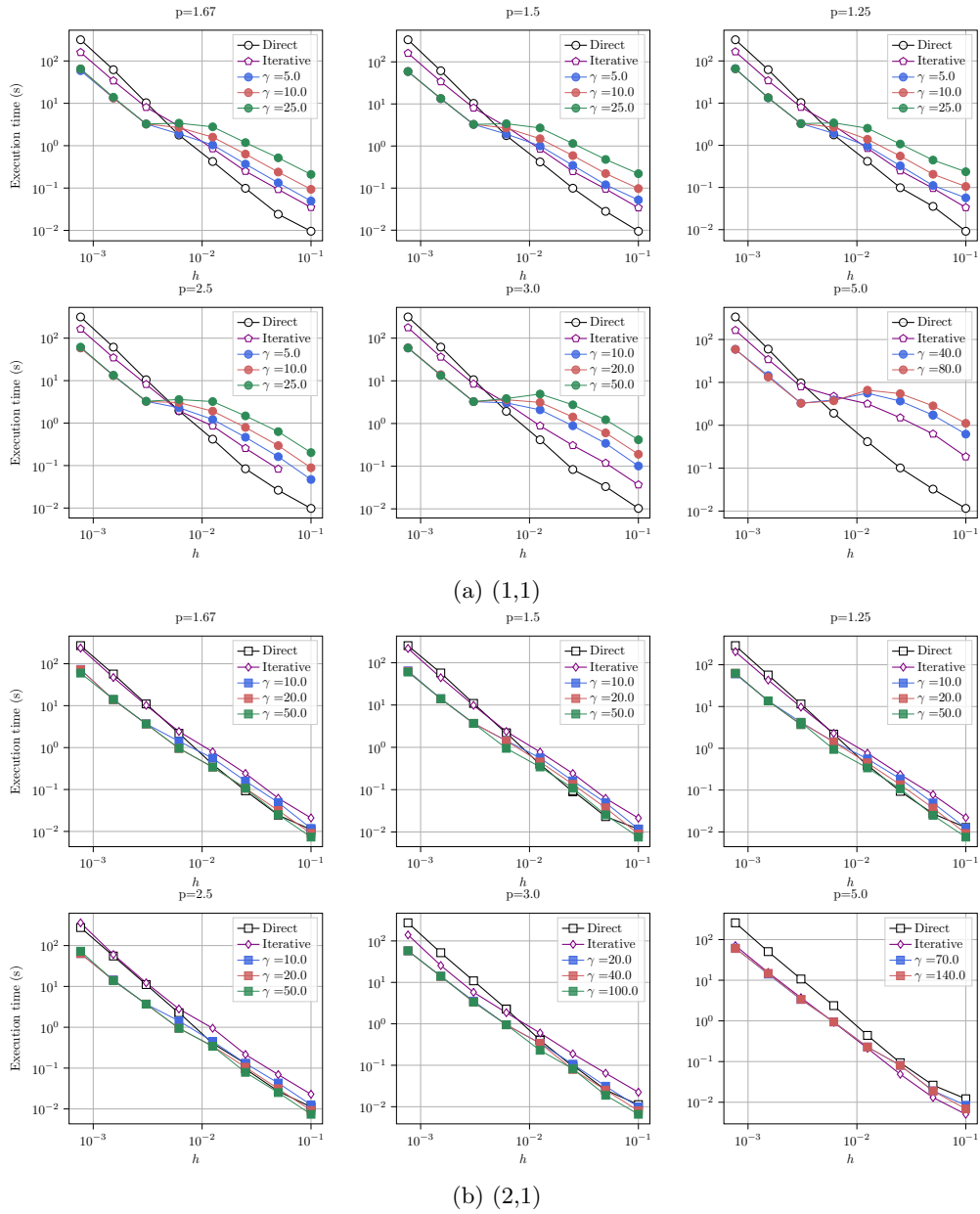


Figure 5.5: Nonlinear ( $p$ -structure) wave problem with  $p \in \{1.25, 1.5, 1.67, 2.5, 3, 5\}$  - Comparison of execution times for a single time step of the splitting procedure and of the semi-implicit scheme (with either a direct or an iterative solver), polynomial orders (1, 1) and (2, 1).

meshes and extrapolated on finer meshes, the value of  $\gamma$  is easy to set in practice.

### 5.3 Vibrating membrane problem

We now consider a model for a 2D vibrating membrane inspired from [53], where a one-dimensional string is considered with two unknowns, the transversal and the longitudinal displacements. Here, we neglect the longitudinal displacement, which leaves us only with the transversal displacement,  $u$ . This leads to

$p$	1.25		1.5		1.67		2.5		5	
Mesh size	direct	gmres	direct	gmres	direct	gmres	direct	gmres	direct	gmres
(1,1) - 0.00306	3.15	2.45	3.16	2.47	3.18	2.48	3.22	2.52	2.98	2.44
(1,1) - 0.00156	4.64	2.55	4.61	2.55	4.64	2.54	4.60	2.59	4.14	2.37
(1,1) - 0.00078	4.83	2.43	5.68	2.72	5.37	2.71	5.35	2.78	5.68	2.77
(2,1) - 0.00306	3.15	2.45	2.97	2.68	2.99	2.77	3.01	3.33	2.98	2.44
(2,1) - 0.00156	4.64	2.55	4.04	3.13	4.08	3.37	3.84	4.17	4.14	2.37
(2,1) - 0.00078	4.83	2.52	3.95	3.39	3.69	3.22	4.32	5.64	5.68	2.78

Table 5.5: Nonlinear (p-structure) wave problem with  $p \in \{1.25, 1.5, 1.67, 2.5, 5\}$  - Gain in execution time of the splitting procedure over the semi-implicit scheme on the four finest meshes,  $\gamma = \gamma^{\min}$ ,  $k = 1$ , equal- and mixed-order settings.

the following nonlinear wave equation (with zero source term):

$$\partial_t^2 u - \nabla \cdot (\mu(\nabla u) \nabla u) = 0, \quad \text{in } \Omega, \forall t \in J, \quad (5.28)$$

with  $\mu : \mathbb{R}^2 \rightarrow \mathbb{R}$  such that

$$\mu(g) = 1 - \alpha \frac{1}{\sqrt{|g|^2 + 1}}, \quad \forall g \in \mathbb{R}^2, \forall \alpha \in [0, 1). \quad (5.29)$$

The nonlinear function  $\mu$  is also considered in the literature in the context of mean-curvature flows. The associated Hamiltonian writes  $\mathcal{H}^\alpha(g) := \frac{1}{2}|g|^2 - \alpha \left[ \sqrt{|g|^2 + 1} - 1 \right]$  for all  $g \in \mathbb{R}^2$ , and equation (5.28) rewrites  $\partial_t^2 u - \nabla \cdot (\nabla_g \mathcal{H}^\alpha(\nabla u)) = 0$ . If  $\alpha = 0$ , (5.28) is equivalent to the linear acoustic wave equation. Moreover, if  $\alpha \neq 0$  and if the gradient  $\nabla u$  becomes very large, the nonlinear part of  $\mu$  becomes negligible with respect to the linear part. Thus, the most nonlinear behavior is expected for small deformations and  $\alpha$  close to one. Typically,  $\alpha = 0.8$  is considered to lead to a mildly nonlinear behavior, and  $\alpha = 0.99$  to a strongly nonlinear behavior. We consider  $\Omega := (0, 1)^2$ , and homogeneous Dirichlet boundary conditions on the displacement  $u$  are enforced. We consider a zero initial condition for  $u$  and an initial velocity

$$v_0(x, y) := e^{-\pi^2 f_p^2 r(x, y)^2}, \quad \text{with } r(x, y) := (0.5 - x)^2 + (0.5 - y)^2, \quad f_p := 3.33, \quad (5.30)$$

which simulates an impact at the center of the domain. The solutions at three sensors located at the points  $(0.167, 0.5)$ ,  $(0.333, 0.333)$  and  $(0.5, 0.5)$  are displayed in Figure 5.6 over the time interval  $[0, 1]$  for  $\alpha \in \{0, 0.8, 0.99\}$ . As  $\alpha$  increases, the wave propagation is slower and the wave front is sharper.

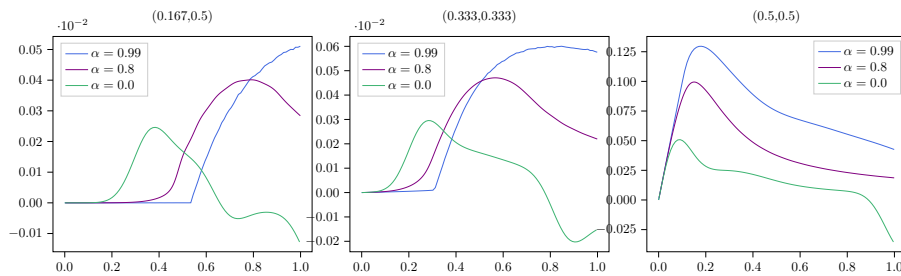


Figure 5.6: Vibrating membrane with  $\alpha \in \{0, 0.8, 0.99\}$  - Reference solution at sensors placed at  $(0.167, 0.5)$ ,  $(0.333, 0.333)$  and  $(0.5, 0.5)$  over the time interval  $[0, 1]$ .

### 5.3.1 Error as a function of the execution time

Figure 5.7 reports the discrete energy error on the displacement as a function of the execution time for P1 and P2 finite elements as well as the HHO method with polynomial orders  $k \in \{0, 1, 2\}$  in the mixed- and

equal-order settings. This error is computed as above, but using now 120 sensors positioned in the triangle  $\{x \in [0, 0.5], y \in [x, 0.5]\}$  using the points with barycentric coordinates  $(\frac{i_1}{6}, \frac{i_2}{6}, \frac{i_3}{6})$  with  $i_1 + i_2 + i_3 = 6$ . More points are considered in this test case in order to precisely capture the wave propagation. The same unstructured mesh sequence as for the p-structure test case is considered. The optimal value of  $\gamma$  is used for each polynomial order, namely  $\gamma = \gamma^{\min}$  in the equal-order setting and  $\gamma = 5\gamma^{\min}$  in the mixed-order setting ( $\gamma^{\min}$  being as before the smallest possible integer value of  $\gamma$  observed experimentally that leads to a converging splitting procedure). Table 5.6 reports the value of  $\gamma^{\min}$  and the optimal value of  $\gamma$  used in Figure 5.7. In the mildly nonlinear case,  $\alpha = 0.8$ , the curve associated with P1 finite elements is between

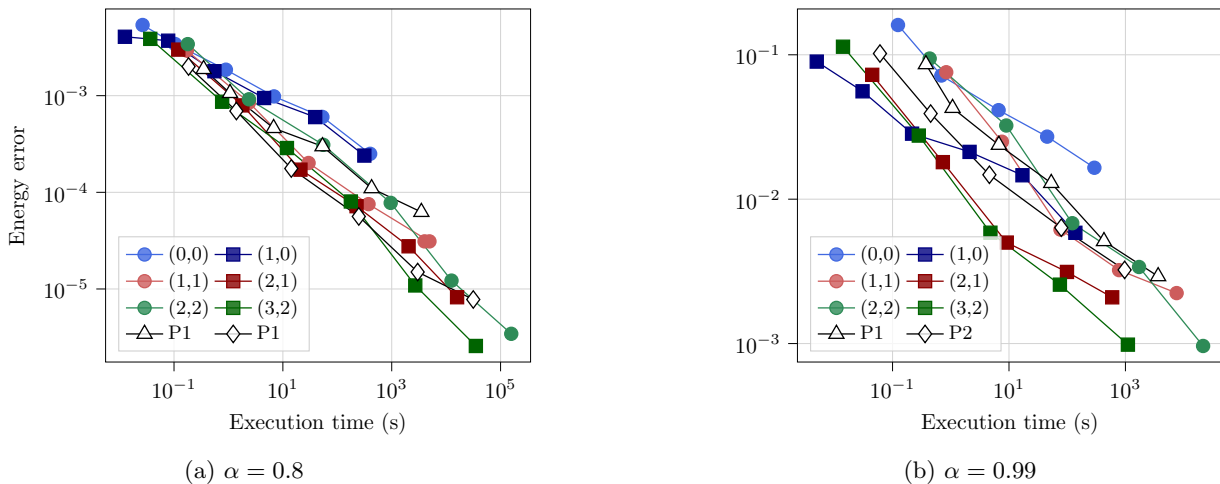


Figure 5.7: Vibrating membrane with  $\alpha \in \{0.8, 0.99\}$  - Discrete energy error as a function of the execution time, P1 finite elements compared to HHO with  $k \in \{0, 1, 2\}$ , mixed- and equal-order settings with splitting (only the best value of  $\gamma$  is displayed).

	(0, 0)	(1, 1)	(2, 2)	(1, 0)	(2, 1)	(3, 2)
$\alpha = 0.8, (\gamma^{\min}, \text{optimal value of } \gamma)$	(2,2)	(4,4)	(6,6)	(3,20)	(5,25)	(8,40)
$\alpha = 0.99, (\gamma^{\min}, \text{optimal value of } \gamma)$	(2,2)	(4,4)	(6,6)	(3,20)	(5,25)	(8,40)

Table 5.6: Vibrating membrane with  $\alpha \in \{0.8, 0.99\}$  -  $\gamma^{\min}$  defined as the smallest integer allowing the splitting procedure to converge, and optimal value (in execution time) of  $\gamma$ .

those associated with the lowest- and the first-order HHO method. Moreover, the first- and second-order HHO method and P2 finite elements are almost equivalent in terms of error as a function of the execution time. In the strongly nonlinear case,  $\alpha = 0.99$ , P1 finite elements are faster than the HHO method with polynomial order (0,0) and deliver a similar performance to that of the other equal-order HHO method. Moreover, mixed-order settings are (much) faster than equal-order settings and than P2 finite elements. The worse efficiency of the splitting in the equal-order setting can be explained by the results of Figure 5.8, which displays the mean number of splitting iterations per time step for each polynomial order and  $\alpha \in \{0.8, 0.99\}$  as a function of the mesh size  $h$ . In the equal-order setting, the number of iterations for  $\alpha = 0.99$  is much larger than that for  $\alpha = 0.8$ , by up to a factor of 10. Instead, in the mixed-order setting, the number of iterations at each time step remains smaller than 10 for both values of  $\alpha$ . To sum up, this experiment shows that the HHO method with splitting and high polynomial orders is more efficient than P1 finite elements and better or equivalent to P2 finite elements in the mildly nonlinear case. In the strongly nonlinear case, the HHO method with polynomial orders (2,1) and (3,2) is much faster than P2 finite elements. Hence, in this experiment as well, the mixed- and high-order HHO method is the most efficient choice.



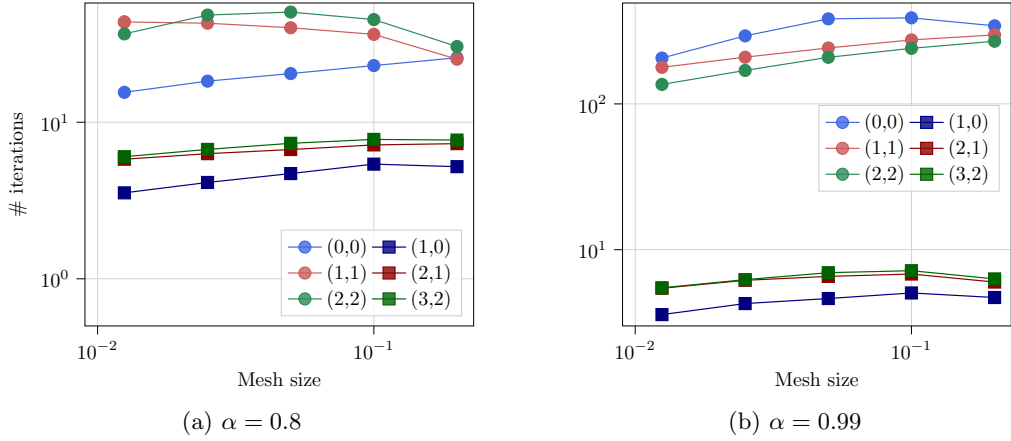


Figure 5.8: Vibrating membrane with  $\alpha \in \{0.8, 0.99\}$  - Mean number of iterations per time step for polynomial orders  $k \in \{0, 1, 2\}$ , mixed- and equal-order settings, with splitting (only the best value of  $\gamma$  is displayed).

### 5.3.2 Aitken acceleration

Considering the large number of iterations needed for the splitting procedure to converge in the equal-order setting, it can be interesting to use acceleration techniques to reduce the computational cost. We can consider for instance the  $\Delta^2$  Aitken acceleration introduced in [8], see, e.g., [134] for the implementation with vector-valued unknowns. Recalling that  $u_{\mathcal{F}}^{n,m}$  denotes the solution at iteration  $m$  of the splitting procedure and at time iteration  $n$ , we define  $\Delta u_{\mathcal{F}}^{n,m} := u_{\mathcal{F}}^{n,m+1} - u_{\mathcal{F}}^{n,m}$ ,  $\Delta^2 u_{\mathcal{F}}^{n,m} := u_{\mathcal{F}}^{n,m+1} - 2u_{\mathcal{F}}^{n,m} + u_{\mathcal{F}}^{n,m-1}$ , and

$$\tilde{u}_{\mathcal{F}}^m := u_{\mathcal{F}}^{n,m} - \frac{(\Delta u_{\mathcal{F}}^{n,m})^2}{\Delta^2 u_{\mathcal{F}}^{n,m}}. \quad (5.31)$$

The sequence of  $(\tilde{u}_{\mathcal{F}}^m)_{m \geq 0}$  converges faster to the same limit  $u_{\mathcal{F}}^{n+1}$  as  $(u_{\mathcal{F}}^{n,m})_{m \geq 0}$  as long as

$$\lim_{m \rightarrow \infty} \frac{u_{\mathcal{F}}^{n,m+1} - u_{\mathcal{F}}^{n+1}}{u_{\mathcal{F}}^{n,m} - u_{\mathcal{F}}^{n+1}} = \lambda \neq 1. \quad (5.32)$$

This definition cannot be directly translated into an algebraic equation, due to the division. One possible adaptation writes as follows in algebraic form, with  $U_{\mathcal{F}}^{n,m}$  the vector containing the degrees of freedom of  $u_{\mathcal{F}}^{n,m}$

$$\tilde{U}_{\mathcal{F}}^m := U_{\mathcal{F}}^{n,m} - \frac{\Delta U_{\mathcal{F}}^{n,m} \cdot \Delta^2 U_{\mathcal{F}}^{n,m}}{\|\Delta^2 U_{\mathcal{F}}^{n,m}\|^2} \Delta U_{\mathcal{F}}^{n,m}.$$

Considering the same membrane setting as above, Figure 5.9 shows the execution time (subfigures 5.9c and 5.9d) and the number of iterations for the splitting procedure with and without Aitken's acceleration (subfigures 5.9a and 5.9b). The optimal value of  $\gamma$  is the one from Table 5.6. As expected, Aitken's acceleration is more effective when the number of splitting iterations is large, i.e. in the equal-order setting and the strongly nonlinear case ( $\alpha = 0.99$ ). Instead, the gain is quite moderate in the mixed-order setting where the number of iterations is very small (less than 4). The equal-order setting with (1,1) benefits from Aitken's acceleration only in the strongly nonlinear case for the same reason. In this case, the splitting procedure with Aitken's acceleration is 30% faster than the splitting procedure without acceleration. In the equal-order setting with polynomial orders (2,2), Aitken's acceleration turns out to be very efficient since the number of iterations is large and does not decrease with mesh refinement. In this case, the gain in execution time is a factor of 5 for the mildly nonlinear case ( $\alpha = 0.8$ ) and more than 10 in the strongly nonlinear case ( $\alpha = 0.99$ ) and on the finest mesh. In conclusion, Aitken's acceleration

is an effective tool that helps improve the performance of the equal-order setting in the strongly nonlinear case.

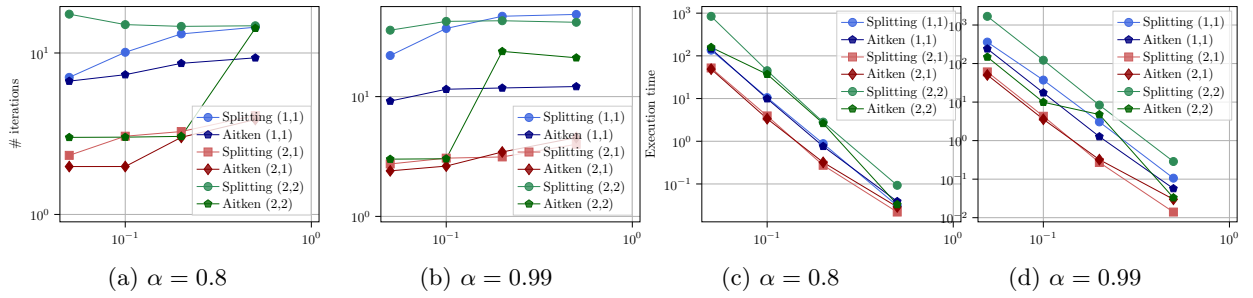


Figure 5.9: Vibrating membrane with  $\alpha \in \{0.8, 0.99\}$  - Mean number of iterations per time step, polynomial orders (1,1), (2,1), (2,2), splitting procedure with and without Aitken's acceleration, optimal value of  $\gamma$ .

## 5.4 Conclusions on the efficiency of the splitting procedure

The numerical experiments illustrated the effectivity of the splitting procedure. There are four points deserving to be put forward.

1. The behavior of the equal- and mixed-order is different. The mixed-order setting is more efficient than the equal-order setting due to the reduced number of splitting iterations to achieve convergence (Figure 5.4 for the p-structure case and Figure 5.7 for the vibrating membrane).
2. The tuning of the parameter  $\gamma$  does not bring any significant difficulty. Indeed, both linear and nonlinear experiments show that, for a value of  $\gamma$  smaller than 100, the quality of the solution is not impacted. All experiments show minimal values of  $\gamma$  ( $\gamma^*$  in the linear case where it can be computed and its approximation  $\gamma^{\min}$  in the nonlinear case) smaller than 100. Moreover, the results of Figures 5.2, 5.3 and 5.5 show that, on fine meshes, the value of  $\gamma$  does not impact the computation time. On coarse meshes, both linear and nonlinear experiments give the same optimal choice:  $\gamma = \gamma^{\min}$  in the equal-order setting and  $\gamma \in [3\gamma^{\min}; 5\gamma^{\min}]$  in the mixed-order setting.
3. In order to increase the accuracy of the solution, it is more efficient to increase the polynomial order than to refine the mesh (see Figure 4.8 in the linear case and Figures 5.4 and 5.7 in the nonlinear case).
4. For nonlinear problems, the splitting procedure combined with a mixed-order setting is always faster than the semi-implicit scheme. The gain in computation time depends on the nonlinearity and increases with mesh refinement. In the equal-order setting, the splitting procedure becomes more efficient than the semi-implicit scheme on refined meshes. Moreover, the HHO method with mixed-order setting, splitting procedure and polynomial order  $k \geq 1$  is always more efficient than the P1 finite element method. In most cases, equal-order computations are also faster. The splitting procedure with mixed-order and large  $k$  is therefore the most effective approach tested in the present work.

## Chapter 6

# Explicit HHO method for structural dynamics

The goal of this chapter is to approximate structural dynamics problems using the HHO method for the space semi-discretization and the leapfrog scheme for the time discretization, together with the splitting procedure presented in the two previous chapters. We first consider a linear elastic behavior and then finite-strain plasticity. Firstly, we study the space semi-discretization of the linear elastic wave equation with the HHO method and recall basic properties of this discretization (Section 6.1). Secondly, we present the time discretization of the elastic wave equation with the leapfrog scheme. In the same manner as for the linear acoustic wave equation in Chapter 3, the HHO-leapfrog scheme applied to the linear elastic wave equation leads to a semi-implicit time marching. Thus, a first goal is to extend the splitting procedure introduced in Chapter 4 to the linear elastic wave equation in order to retrieve a fully explicit time-stepping scheme. This is done in Section 6.2 together with a numerical study of the splitting parameters. Thirdly, we conduct numerical experiments comparing the HHO-leapfrog method with splitting to the standard finite element method on a contrasted material test case (Section 6.3). Fourthly, we present the full discretization of the structural dynamics equations with finite-strain plasticity and extend the splitting procedure to this case (Section 6.4). Finally, numerical examples illustrate the efficiency of the splitting procedure for this nonlinear equation (Section 6.5). These experiments also illustrate the robustness to locking of the HHO space semi-discretization.

### 6.1 The HHO method for the elastic wave equation

We are interested in the evolution of an elastic body whose reference configuration is the open bounded domain  $\Omega \subset \mathbb{R}^d$ ,  $d \in \{2, 3\}$ , during the time interval  $J := [0, \mathfrak{T}]$ ,  $\mathfrak{T} > 0$ . This body is subject to the action of a body force  $\mathbf{f} : \Omega \times J \rightarrow \mathbb{R}^d$  and a prescribed displacement  $\mathbf{u}_D : \partial\Omega_D \times J \rightarrow \mathbb{R}^d$  on part of its boundary  $\partial\Omega_D \subset \partial\Omega$ . We assume that  $\partial\Omega_D$  has positive measure in order to prevent rigid-body motions. The rest of the boundary is supposed to be free of traction. Under these external actions, the body is deformed. A point  $\mathbf{x} \in \Omega$  in the initial configuration is then mapped to a point  $\mathbf{x}' = \mathbf{x} + \mathbf{u}(\mathbf{x})$  where  $\mathbf{u} : \Omega \rightarrow \mathbb{R}^d$  is the displacement field.

#### 6.1.1 The linear elastic wave equation

In the context of linear elasticity, we consider only infinitesimal deformations, for which a relevant measure is the linearized strain tensor

$$\boldsymbol{\epsilon}(\mathbf{u}) := \boldsymbol{\nabla}_s \mathbf{u} := \frac{1}{2}(\boldsymbol{\nabla} \mathbf{u} + \boldsymbol{\nabla} \mathbf{u}^\top) \in \mathbb{R}_{\text{sym}}^{d \times d}, \quad (6.1)$$

where  $\mathbb{R}_{\text{sym}}^{d \times d}$  is the subset of symmetric matrices in  $\mathbb{R}^{d \times d}$ . The strain-stress relation is given by

$$\boldsymbol{\sigma}(\boldsymbol{\epsilon}) = \mathbb{A}\boldsymbol{\epsilon} := 2\mu\boldsymbol{\epsilon} + \lambda \text{tr}(\boldsymbol{\epsilon})\mathbf{I}_d \in \mathbb{R}_{\text{sym}}^{d \times d}, \quad (6.2)$$

where  $\mathbb{A}$  is the fourth-order stiffness tensor,  $\lambda \geq 0$  and  $\mu > 0$  are the Lamé coefficients and  $\mathbf{I}_d$  is the identity matrix in  $\mathbb{R}^{d \times d}$ . The Lamé coefficients can be computed from the Young modulus,  $E$ , and the Poisson ratio,  $\nu$ , using the following relations:

$$\mu = \frac{E}{2(1+\nu)}, \quad \lambda = \frac{E\nu}{(1+\nu)(1-2\nu)}. \quad (6.3)$$

In this setting, the linear elastic wave equation consists in finding  $\mathbf{u} : \Omega \times J \rightarrow \mathbb{R}^d$  such that (6.2) holds in  $\Omega \times J$  together with

$$\rho \partial_t^2 \mathbf{u} - \nabla \cdot (\boldsymbol{\sigma}(\boldsymbol{\epsilon}(\mathbf{u}))) = \mathbf{f}, \quad \text{in } \Omega \times J, \quad (6.4a)$$

$$\mathbf{u} = \mathbf{u}_D, \quad \text{on } \partial\Omega_D \times J, \quad (6.4b)$$

$$\mathbf{u}|_{t=0} = \mathbf{u}_0, \quad \text{in } \Omega, \quad (6.4c)$$

$$\partial_t \mathbf{u}|_{t=0} = \mathbf{v}_0, \quad \text{in } \Omega, \quad (6.4d)$$

where  $\rho > 0$  is the material density. For the sake of simplicity, we assume that  $\mathbf{u}_D = \mathbf{0}$  in this section and in Section 6.2. We also assume that  $\rho$ ,  $\lambda$  and  $\mu$  are piecewise constant on a polyhedral partition of  $\Omega$ .

The wave equation (6.4) describes the propagation of several types of elastic waves, such as P-waves (pressure waves) of speed  $c_P := \sqrt{\frac{\lambda+2\mu}{\rho}}$  and S-waves (shear waves) of speed  $c_S := \sqrt{\frac{\mu}{\rho}}$ . In the incompressible limit  $\lambda \rightarrow \infty$ , we have  $c_P \rightarrow \infty$ , whereas  $c_S$  remains finite. Table 6.1 displays the values of the density,  $\rho$ , the Young modulus,  $E$ , the Poisson ratio,  $\nu$ , and the Lamé coefficients (computed using (6.3)) for a selection of widely used materials. The presented values are good reference points to compare materials and study the impact of the material properties on the splitting procedure. Rubber is usually considered as nearly incompressible with a value of  $\nu \approx \frac{1}{2}$ . We are using here the softened value  $\nu = 0.499$  to avoid the singular case  $\nu = 0.5$  (so that the ratio  $\frac{\lambda}{\mu}$  for rubber is only one order of magnitude larger than that of gold). We also consider a reference material with unit Lamé coefficients and density. This material will be used in our numerical experiments with analytical solutions to simplify some computations.

Material	Steel	Copper	Diamond	Gold	Rubber	Reference
Density $\rho$ (kg/m <sup>3</sup> )	7800	8930	3520	18900	920	1
Young modulus $E$ (GPa)	200	120	1220	80	2	2.5
Poisson ratio $\nu$	0.3	0.35	0.2	0.42	0.499	0.25
Lamé coefficients $(\mu, \lambda)$ (GPa)	(77, 115)	(43, 101)	(508, 339)	(27, 144)	(0.74, 366)	(1, 1)

Table 6.1: Material properties for a selection of materials: density, Young modulus, Poisson ratio and Lamé coefficients.

Assuming  $\mathbf{f} \in L^2(J; \mathbf{L}^2(\Omega))$ , we seek the solution of (6.4) in  $U := H^2(J; \mathbf{L}^2(\Omega)) \cap L^2(J; \mathbf{H}_0^1(\Omega))$  such that

$$(\rho \partial_t^2 \mathbf{u}, \mathbf{w})_\Omega + (\boldsymbol{\sigma}(\boldsymbol{\epsilon}(\mathbf{u})), \boldsymbol{\epsilon}(\mathbf{w}))_\Omega = (\mathbf{f}, \mathbf{w})_\Omega, \quad \text{a.e. } t \in J, \quad \forall \mathbf{w} \in \mathbf{H}_0^1(\Omega), \quad (6.5)$$

with the initial conditions

$$\mathbf{u}(0) = \mathbf{u}_0, \quad \partial_t \mathbf{u}(0) = \mathbf{v}_0. \quad (6.6)$$

**Remark 6.1.1** (Rewriting using Lamé coefficients). Notice that (6.5) can be rewritten

$$(\rho \partial_t^2 \mathbf{u}, \mathbf{w})_\Omega + 2\mu(\boldsymbol{\epsilon}(\mathbf{u}), \boldsymbol{\epsilon}(\mathbf{w}))_\Omega + \lambda(\nabla \cdot \mathbf{u}, \nabla \cdot \mathbf{w})_\Omega = (\mathbf{f}, \mathbf{w})_\Omega, \quad \text{a.e. } t \in J, \quad \forall \mathbf{w} \in H_0^1(\Omega). \quad (6.7)$$

This is the classical weak formulation for the linear elastic wave equation. The divergence is obtained by taking the trace of the symmetric gradient  $\epsilon(\mathbf{u})$ . In view of extension to plasticity, we keep the formulation (6.5) in what follows.  $\square$

## 6.1.2 HHO operators

### 6.1.2.1 HHO unknowns

We consider a sequence of meshes  $(\mathcal{T}_h)_{h>0}$  following the same principles as in Section 2.2.1. Each mesh  $\mathcal{T}_h$  is assumed to be compatible with the partition on which  $\rho, \mu, \lambda$  are piecewise constant, and with the partition of  $\partial\Omega$  induced by the subset  $\partial\Omega_D$ . Let the integer  $k \geq 1$  be the polynomial order of the face unknowns and let  $l \in \{k, k+1\}$  be the order of the cell unknowns. Remark that, so as to prevent rigid-body motions, the lowest polynomial order  $k = 0$  is not available. Let  $\mathbb{P}_d^l(T; \mathbb{R}^d)$  (resp.  $\mathbb{P}_{d-1}^k(F; \mathbb{R}^d)$ ) denote the set of vector-valued  $d$ -variate (resp.  $(d-1)$ -variate) polynomials of degree at most  $l$  (resp.  $k$ ) restricted to the cell  $T \in \mathcal{T}_h$  (resp. to the face  $F \in \mathcal{F}_h$ ). The linear space composed of all the cell degrees of freedom is denoted  $\mathbf{U}_{\mathcal{T}}^l$ , and the linear space composed of all the face degrees of freedom is denoted  $\mathbf{U}_{\mathcal{F}}^k$ . These spaces are defined as Cartesian products in the form

$$\mathbf{U}_{\mathcal{T}}^l := \prod_{T \in \mathcal{T}_h} \mathbb{P}_d^l(T; \mathbb{R}^d), \quad \mathbf{U}_{\mathcal{F}}^k := \prod_{F \in \mathcal{F}_h} \mathbb{P}_{d-1}^k(F; \mathbb{R}^d), \quad (6.8)$$

and we slightly abuse the notation by viewing an element  $\mathbf{w}_{\mathcal{T}} = (\mathbf{w}_T)_{T \in \mathcal{T}_h} \in \mathbf{U}_{\mathcal{T}}^l$  as a function defined a.e. over  $\Omega$  such that  $\mathbf{w}_{\mathcal{T}}|_T := \mathbf{w}_T$  for all  $T \in \mathcal{T}_h$ . The collection of all the cell and face degrees of freedom is the hybrid space

$$\widehat{\mathbf{U}}_h^{l,k} := \mathbf{U}_{\mathcal{T}}^l \times \mathbf{U}_{\mathcal{F}}^k. \quad (6.9)$$

A generic element of  $\widehat{\mathbf{U}}_h^{l,k}$  is denoted  $\widehat{\mathbf{w}}_h := (\mathbf{w}_{\mathcal{T}}, \mathbf{w}_{\mathcal{F}}) \in \mathbf{U}_{\mathcal{T}}^l \times \mathbf{U}_{\mathcal{F}}^k$  and, as before, variables with hats refer to hybrid variables. For a given cell  $T \in \mathcal{T}_h$ , we also define a local hybrid space of degrees of freedom

$$\widehat{\mathbf{U}}_T^{l,k} := \mathbb{P}_d^l(T; \mathbb{R}^d) \times \mathbf{U}_{\partial T}^k, \quad \mathbf{U}_{\partial T}^k := \prod_{F \in \mathcal{F}_T} \mathbb{P}_{d-1}^k(F; \mathbb{R}^d). \quad (6.10)$$

Then  $\widehat{\mathbf{w}}_T := (\mathbf{w}_T, \mathbf{w}_{\partial T} = (\mathbf{w}_F)_{F \in \mathcal{F}_T}) \in \widehat{\mathbf{U}}_T^{l,k}$  denotes a generic local hybrid unknown in  $T$ , composed of one cell unknown and the collection of the face unknowns for all the faces in  $\mathcal{F}_T$ . As above, we slightly abuse the notation by viewing an element  $\mathbf{w}_{\partial T} = (\mathbf{w}_F)_{F \in \mathcal{F}_T} \in \mathbf{U}_{\partial T}^k$  as a function defined a.e. over  $\partial T$  such that  $\mathbf{w}_{\partial T}|_F := \mathbf{w}_F$  for all  $F \in \mathcal{F}_T$ . Let

$$\mathbf{U}_{\mathcal{F},0}^k := \{\mathbf{v}_{\mathcal{F}} \in \mathbf{U}_{\mathcal{F}}^k, \text{ s.t } \mathbf{v}_F = \mathbf{0}, \forall F \subset \partial\Omega_D\}. \quad (6.11)$$

The space of hybrid unknowns respecting the Dirichlet conditions is denoted

$$\widehat{\mathbf{U}}_{h,0}^{l,k} := \mathbf{U}_{\mathcal{T}}^l \times \mathbf{U}_{\mathcal{F},0}^k. \quad (6.12)$$

We keep the same notation for the  $L^2$ -orthogonal projection onto polynomial spaces as for the acoustic wave equation: For all  $T \in \mathcal{T}_h$ ,  $\Pi_T^l$  is the projection onto  $\mathbb{P}_d^l(T; \mathbb{R}^d)$  and for all  $F \in \mathcal{F}_h$ ,  $\Pi_F^k$  is the projection onto  $\mathbb{P}_{d-1}^k(F; \mathbb{R}^d)$ . We define the local projector  $\hat{I}_T^{k,l} : \mathbf{H}^1(T) \rightarrow \widehat{\mathbf{U}}_T^{l,k}$  for all  $T \in \mathcal{T}_h$  such that, for all  $\mathbf{w} \in \mathbf{H}^1(T)$ ,

$$\hat{I}_T^{k,l}(\mathbf{w}) := (\Pi_T^l(\mathbf{w}), (\Pi_F^k(\mathbf{w}))_{F \in \mathcal{F}_T}) \in \widehat{\mathbf{U}}_T^{l,k}. \quad (6.13)$$

Then the global projection operator can be defined as  $\hat{I}_h^{k,l} : \mathbf{H}^1(\Omega) \rightarrow \widehat{\mathbf{U}}_h^{l,k}$  such that, for all  $\mathbf{w} \in \mathbf{H}^1(\Omega)$ ,

$$\hat{I}_h^{k,l}(\mathbf{w}) := ((\Pi_T^l(\mathbf{w}))_{T \in \mathcal{T}_h}, (\Pi_F^k(\mathbf{w}))_{F \in \mathcal{F}_h}). \quad (6.14)$$

### 6.1.2.2 Symmetric gradient reconstruction operator.

In the same manner as for the acoustic wave equation, we define a gradient reconstruction operator locally on each cell  $T \in \mathcal{T}_h$ . But, since the deformation measure is the symmetric gradient, we directly define a symmetric gradient reconstruction operator

$$\mathbf{E}_T^k : \widehat{\mathbf{U}}_T^{l,k} \rightarrow \mathbf{Z}(T; \mathbb{R}_{\text{sym}}^{d \times d}), \quad (6.15)$$

where  $\mathbf{Z}(T; \mathbb{R}_{\text{sym}}^{d \times d}) := \mathbb{P}_d^k(T; \mathbb{R}_{\text{sym}}^{d \times d})$  is such that, for all  $\mathbf{q} \in \mathbb{P}_d^k(T; \mathbb{R}_{\text{sym}}^{d \times d})$ ,

$$(\mathbf{E}_T^k(\hat{\mathbf{v}}_T), \mathbf{q})_T = (\nabla_s \mathbf{v}_T, \mathbf{q})_T + (\mathbf{v}_{\partial T} - \mathbf{v}_T, \mathbf{q} \cdot \mathbf{n}_T)_{\partial T}. \quad (6.16)$$

**Lemma 6.1.2** (Approximation properties for  $\mathbf{E}_T^k \circ \hat{I}_T^{k,l}$ ). *The following holds true for all  $h > 0$  and all  $T \in \mathcal{T}_h$ :*

$$\mathbf{E}_T^k(\hat{I}_T^{k,l}(\mathbf{v})) = \Pi_T^k(\nabla_s \mathbf{v}). \quad (6.17)$$

Moreover, there exists a real number  $C > 0$  such that, for all  $h > 0$ , all  $T \in \mathcal{T}_h$  and all  $\mathbf{v} \in \mathbf{H}^{k+2}(T)$ , we have

$$\|\nabla_s \mathbf{v} - \mathbf{E}_T^k(\hat{I}_T^{k,l}(\mathbf{v}))\|_T + h_T^{\frac{1}{2}} \|\nabla_s \mathbf{v} - \mathbf{E}_T^k(\hat{I}_T^{k,l}(\mathbf{v}))\|_{\partial T} \leq Ch_T^{k+1} |\mathbf{v}|_{\mathbf{H}^{k+2}(T)}. \quad (6.18)$$

### 6.1.2.3 Stabilization operator.

Let  $T \in \mathcal{T}_h$ . For all  $\hat{\mathbf{w}}_T \in \widehat{\mathbf{U}}_T^{l,k}$ , we set  $\delta_T(\hat{\mathbf{w}}_T) := \mathbf{w}_{\partial T} - \mathbf{w}_T|_{\partial T}$  on  $\partial T$  and  $\delta_{TF}(\hat{\mathbf{w}}_T) := \delta_T(\hat{\mathbf{w}}_T)|_F$  for all  $F \in \mathcal{F}_T$ . In the mixed-order setting, we define the local stabilization operator  $\mathbf{S}_{TF}^{\text{MO}}$  as

$$\mathbf{S}_{TF}^{\text{MO}}(\hat{\mathbf{w}}_T) := \Pi_F^k(\delta_{TF}(\hat{\mathbf{w}}_T)), \quad \forall \hat{\mathbf{w}}_T \in \widehat{\mathbf{U}}_T^{l,k}. \quad (6.19)$$

In the equal-order setting, the definition of the stabilization requires the computation of a displacement reconstruction operator. There are two main choices for this displacement operator. The first one is the traditional HHO reconstruction operator introduced in [82]. It is denoted  $\mathbf{R}_T^{\text{SG},k+1} : \widehat{\mathbf{U}}_T^{l,k} \rightarrow \mathbb{P}_d^{k+1}(T; \mathbb{R}^d)$  and verifies, for all  $\hat{\mathbf{v}}_T \in \widehat{\mathbf{U}}_T^{l,k}$  and for all  $\mathbf{w} \in \mathbb{P}_d^{k+1}(T; \mathbb{R}^d)$ ,

$$\left( \nabla_s \mathbf{R}_T^{\text{SG},k+1}(\hat{\mathbf{v}}_T), \nabla_s \mathbf{w} \right)_T = (\nabla_s \mathbf{v}_T, \nabla_s \mathbf{w})_T + (\mathbf{v}_{\partial T} - \mathbf{v}_T, \nabla_s \mathbf{w} \cdot \mathbf{n}_T)_{\partial T}, \quad (6.20a)$$

$$\int_T \mathbf{R}_T^{\text{SG},k+1}(\hat{\mathbf{v}}_T) = \int_T \mathbf{v}_T, \quad \int_T \nabla_{\text{ss}} \left( \mathbf{R}_T^{\text{SG},k+1}(\hat{\mathbf{v}}_T) \right) = \frac{1}{2} \int_{\partial T} \mathbf{n}_T \otimes \mathbf{v}_{\partial T} - \mathbf{v}_{\partial T} \otimes \mathbf{n}_T, \quad (6.20b)$$

where the superscript SG stands for symmetric gradient,  $\nabla_{\text{ss}}$  denotes the skew-symmetric part of the gradient and  $\otimes$  the tensor product operation. Whatever choice is made for the right-hand side in (6.20b) is not important in what follows.

The second choice is to consider the full gradient reconstruction operator  $\mathbf{R}_T^{\text{FG},k+1}$  as the vector-valued version of the reconstruction operator used for the acoustic wave equation, the superscript FG standing for full gradient. The reconstruction operator  $\mathbf{R}_T^{\text{FG},k+1} : \widehat{\mathbf{U}}_T^{l,k} \rightarrow \mathbb{P}_d^{k+1}(T; \mathbb{R}^d)$  verifies, for all  $\hat{\mathbf{v}}_T \in \widehat{\mathbf{U}}_T^{l,k}$  and for all  $\mathbf{w} \in \mathbb{P}_d^{k+1}(T; \mathbb{R}^d)$ ,

$$\left( \nabla \mathbf{R}_T^{\text{FG},k+1}(\hat{\mathbf{v}}_T), \nabla \mathbf{w} \right)_T = (\nabla \mathbf{v}_T, \nabla \mathbf{w})_T + (\mathbf{v}_{\partial T} - \mathbf{v}_T, \nabla \mathbf{w} \cdot \mathbf{n}_T)_{\partial T}. \quad (6.21a)$$

$$\int_T \mathbf{R}_T^{\text{FG},k+1}(\hat{\mathbf{v}}_T) = \int_T \mathbf{v}_T. \quad (6.21b)$$

In what follows, we denote by  $\mathbf{R}_T^{x,k+1}$  either  $\mathbf{R}_T^{\text{SG},k+1}$  or  $\mathbf{R}_T^{\text{FG},k+1}$ . We will see in Section 6.2 that only  $\mathbf{R}_T^{\text{FG},k+1}$  leads to a converging splitting procedure. In both cases, the local stabilization operator writes

$$\mathbf{S}_{TF}^x(\hat{\mathbf{w}}_T) := \Pi_F^k(\delta_{TF}(\hat{\mathbf{w}}_T)) + (I - \Pi_T^k) \mathbf{R}_T^{x,k+1}(0, \delta_T(\hat{\mathbf{w}}_T))|_F, \quad \forall \hat{\mathbf{w}}_T \in \widehat{\mathbf{U}}_T^{l,k} := \widehat{\mathbf{U}}_T^k, \quad (6.22)$$

with  $x$  being either SG or FG.

In both equal- and mixed-order settings, we define, for all  $\hat{\mathbf{v}}_T \in \widehat{\mathbf{U}}_T^{l,k}$ , the cell stabilization operator  $\mathbf{S}_{\partial T}^x := (\mathbf{S}_{TF}^x)_{F \in \mathcal{F}_T}$  collecting the operators  $\mathbf{S}_{TF}^x$  from all the faces  $F \in \mathcal{F}_T$ . Then, for all  $\hat{\mathbf{v}}_T, \hat{\mathbf{w}}_T \in \widehat{\mathbf{U}}_T^{l,k}$ , the local stabilization form is defined as

$$\mathbf{s}_T^x(\hat{\mathbf{v}}_T, \hat{\mathbf{w}}_T) = \gamma \mu_T (\eta_T^{-1} \mathbf{S}_{\partial T}^x(\hat{\mathbf{v}}_T), \mathbf{S}_{\partial T}^x(\hat{\mathbf{w}}_T))_{\partial T}, \quad (6.23)$$

with, on each face, the scaling factor  $\eta_T|_F := \eta_{TF}$  equal to  $h_F$  or  $h_T$ ,  $\gamma > 0$  a scaling parameter,  $\mu_T := \mu|_T$  the Lamé coefficient in  $T$ ,  $x = \text{MO}$  in the mixed-order setting and  $x \in \{\text{SG}, \text{FG}\}$  in the equal-order setting. We consider the choice  $\eta_{TF} = h_T^{-1}$  in what follows. In the original HHO method for elasticity [82], the value  $\gamma = 2$  is considered.

#### 6.1.2.4 Stability and boundedness.

A direct verification shows that the map  $\|\cdot\|_{\text{HHO}} : \widehat{\mathbf{U}}_h^{l,k} \rightarrow \mathbb{R}_+$  such that

$$\|\hat{\mathbf{v}}_h\|_{\text{HHO}}^2 := \sum_{T \in \mathcal{T}_h} \mu_T \{ \|\nabla_s \mathbf{v}_T\|_T^2 + h_T^{-1} \|\mathbf{v}_{\partial T} - \mathbf{v}_T\|_{\partial T}^2 \}, \quad \forall \hat{\mathbf{v}}_h \in \widehat{\mathbf{U}}_h^{l,k}, \quad (6.24)$$

defines a norm on  $\widehat{\mathbf{U}}_{h,0}^{l,k}$  (and a seminorm on  $\widehat{\mathbf{U}}_h^{l,k}$ ), and we have the following important result [1, Lem. 1].

**Lemma 6.1.3** (Stability and boundedness). *There are  $0 < \alpha \leq \varpi < \infty$ , independent of  $\mu$ , such that, for all  $h > 0$  and all  $\hat{\mathbf{v}}_h \in \widehat{\mathbf{U}}_{h,0}^{l,k}$ ,*

$$\alpha \|\hat{\mathbf{v}}_h\|_{\text{HHO}}^2 \leq \frac{1}{2} \|\mu^{\frac{1}{2}} \mathbf{E}_T^k(\hat{\mathbf{v}}_h)\|_{\Omega}^2 + \|\mu^{\frac{1}{2}} h^{-\frac{1}{2}} \mathbf{S}_{\partial \mathcal{T}_h}^x(\hat{\mathbf{v}}_h)\|_{\partial \mathcal{T}_h}^2 \leq \varpi \|\hat{\mathbf{v}}_h\|_{\text{HHO}}^2, \quad (6.25)$$

with  $\mathbf{S}_{\partial \mathcal{T}_h}^x := (\mathbf{S}_{\partial T}^x)_{T \in \mathcal{T}_h}$ ,

$$\|\mu^{\frac{1}{2}} h^{-\frac{1}{2}} \mathbf{S}_{\partial \mathcal{T}_h}^x(\hat{\mathbf{v}}_h)\|_{\partial \mathcal{T}_h}^2 := \sum_{T \in \mathcal{T}_h} \|\mu_T^{\frac{1}{2}} h_T^{-\frac{1}{2}} \mathbf{S}_{\partial T}^x(\hat{\mathbf{v}}_T)\|_{\partial T}^2 := \sum_{T \in \mathcal{T}_h} \mu_T h_T^{-1} (\mathbf{S}_{\partial T}^x(\hat{\mathbf{v}}_T), \mathbf{S}_{\partial T}^x(\hat{\mathbf{v}}_T))_{\partial T},$$

where we considered  $x = \text{SG}$  in the equal-order setting and  $x = \text{MO}$  in the mixed-order setting, and the global symmetric gradient reconstruction such that  $\mathbf{E}_T^k(\hat{\mathbf{v}}_h)|_T := \mathbf{E}_T^k(\hat{\mathbf{v}}_T)$  for all  $T \in \mathcal{T}_h$  with  $\mathbf{E}_T^k$  defined in (6.16).

Regarding the stability in the equal-order setting with  $x = \text{FG}$ , no proof is available in the literature. However, we did verify numerically that  $\|\mathbf{E}_T^k(\cdot)\|_T^2 + \|h_T^{-\frac{1}{2}} \mathbf{S}_{\partial T}^{\text{SG}}(\cdot)\|_{\partial T}^2$  and  $\|\mathbf{E}_T^k(\cdot)\|_T^2 + \|h_T^{-\frac{1}{2}} \mathbf{S}_{\partial T}^{\text{FG}}(\cdot)\|_{\partial T}^2$  yield equivalent norms by computing the generalized eigenvalues of the equivalent algebraic problem on a single mesh cell. Let us denote  $\mathcal{A}_T^{\text{SG}}$  the matrix associated with the norm  $\|\mathbf{E}_T^k(\cdot)\|_T^2 + \|h_T^{-\frac{1}{2}} \mathbf{S}_{\partial T}^{\text{SG}}(\cdot)\|_{\partial T}^2$  and  $\mathcal{A}_T^{\text{FG}}$  the matrix associated with the norm  $\|\mathbf{E}_T^k(\cdot)\|_T^2 + \|h_T^{-\frac{1}{2}} \mathbf{S}_{\partial T}^{\text{FG}}(\cdot)\|_{\partial T}^2$ . Figure 6.1 displays the sorted eigenvalues of the generalized eigenproblem  $\mathcal{A}_T^{\text{FG}} X = \lambda \mathcal{A}_T^{\text{SG}} X$  on a square and a right-triangular cell, in the equal-order setting,  $k \in \{1, 2, 3\}$ . The mixed-order setting is not considered, since  $\mathbf{S}_{\partial T}^{\text{FG}}$  is only used in the equal-order setting. We observe that all the eigenvalues are finite and larger than zero. Therefore, using Lemma 6.1.3, we conjecture the following result.

**Conjecture 6.1.4** (Full Gradient stability and boundedness). *There are  $0 < \alpha \leq \varpi < \infty$ , independent of  $\mu$ , such that, for all  $h > 0$  and all  $\hat{\mathbf{v}}_h \in \widehat{\mathbf{U}}_{h,0}^{l,k}$ ,*

$$\alpha' \|\hat{\mathbf{v}}_h\|_{\text{HHO}}^2 \leq \frac{1}{2} \|\mu^{\frac{1}{2}} \mathbf{E}_T^k(\hat{\mathbf{v}}_h)\|_{\Omega}^2 + \|\mu^{\frac{1}{2}} \eta_T^{-\frac{1}{2}} \mathbf{S}_{\partial \mathcal{T}_h}^{\text{FG}}(\hat{\mathbf{v}}_h)\|_{\partial \mathcal{T}_h}^2 \leq \varpi' \|\hat{\mathbf{v}}_h\|_{\text{HHO}}^2. \quad (6.26)$$

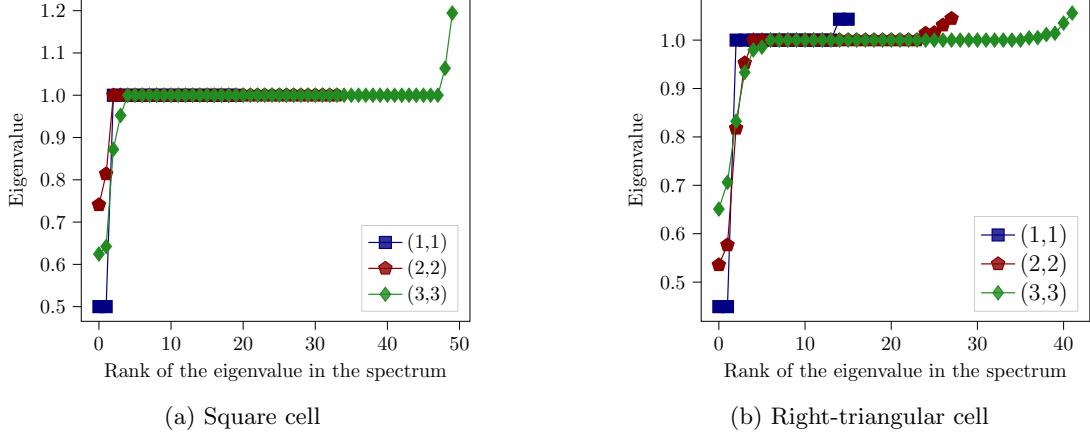


Figure 6.1: Linear elasticity - Eigenvalues of the generalized eigenproblem  $\mathcal{A}_T^{\text{FG}} X = \lambda \mathcal{A}_T^{\text{SG}} X$ , equal-order setting with polynomial orders  $k \in \{1, 2, 3\}$ , square and right-triangular cells.

**Remark 6.1.5** (Variants). Considering the above definitions, more HHO variants can be defined. Firstly, it is possible to use the symmetric gradient reconstruction operator  $\mathbf{R}_T^{\text{SG},k+1}$  to define the following symmetric gradient reconstruction operator:

$$\mathbf{E}_T^{\text{SG},k} := \nabla_s \mathbf{R}_T^{\text{SG},k+1}. \quad (6.27)$$

This is actually what is done in the traditional HHO formulation introduced in [82]. In this case, one is led to define the broken elliptic projection  $\mathcal{E}_T^{\text{SG},k+1} : \mathbf{H}^1(\Omega) \rightarrow \mathcal{U}_T^{k+1}$ , so that, for all  $\mathbf{w} \in \mathbf{H}^1(\Omega)$  and all  $T \in \mathcal{T}_h$ ,  $\mathcal{E}_T^{\text{SG},k+1}(\mathbf{w})|_T$  is uniquely defined by the relations

$$\left( \nabla_s(\mathcal{E}_T^{\text{SG},k+1}(\mathbf{w}) - \mathbf{w}), \nabla_s \mathbf{q} \right)_T = 0, \forall \mathbf{q} \in \mathbb{P}_d^{k+1, \text{RM}}(T; \mathbb{R}^d), \quad (6.28a)$$

$$\int_T \{\mathcal{E}_T^{\text{SG},k+1}(\mathbf{w}) - \mathbf{w}\} = \mathbf{0}, \quad \int_T \nabla_{\text{ss}}(\mathcal{E}_T^{\text{SG},k+1}(\mathbf{w}) - \mathbf{w}) = \mathbf{0}, \quad (6.28b)$$

where  $\mathbb{P}_d^{k+1, \text{RM}}(T; \mathbb{R}^d) := \{\mathbf{q} \in \mathbb{P}_d^{k+1}(T; \mathbb{R}^d) \mid \int_T \mathbf{q} = \mathbf{0}, \int_T \nabla_{\text{ss}} \mathbf{q} = \mathbf{0}\}$  denotes the subset of polynomials in  $\mathbb{P}_d^{k+1}(T; \mathbb{R}^d)$  with zero rigid-body motions. Then, the symmetric gradient reconstruction operator verifies the following polynomial consistency property [82, Lem. 2]:

$$\mathbf{E}_T^{\text{SG},k}(\hat{I}_T^{k,l}(\mathbf{v})) = \nabla_s \mathcal{E}_T^{\text{SG},k+1}(\mathbf{v}), \quad \forall \mathbf{v} \in \mathbf{H}^1(T), \quad (6.29)$$

leading to the approximation property (6.18) with  $\mathbf{E}_T^k = \mathbf{E}_T^{\text{SG},k}$ . This gradient reconstruction also verifies a stability property in the form

$$\alpha \|\hat{\mathbf{v}}_h\|_{\text{HHO}}^2 \leq \frac{1}{2} \|\mu^{\frac{1}{2}} \mathbf{E}_T^{\text{SG},k}(\hat{\mathbf{v}}_h)\|_{\Omega}^2 + |\mu^{\frac{1}{2}} \eta_T^{-\frac{1}{2}} \mathbf{S}_{\partial T}^{\text{SG}}(\hat{\mathbf{v}}_h)|_{\Omega}^2 \leq \varpi \|\hat{\mathbf{v}}_h\|_{\text{HHO}}^2, \quad \forall \hat{\mathbf{v}}_h \in \hat{\mathcal{U}}_{h,0}^{l,k}, \quad (6.30)$$

see [82, Lem. 4]. This is not the main reconstruction considered in this thesis, since, as shown in Section 6.2, this reconstruction does not lead to a converging splitting procedure in the equal-order setting.

Secondly, one could also want to use the full gradient reconstruction to define yet another symmetric gradient reconstruction operator as follows:

$$\mathbf{E}_T^{\text{FG},k} := \nabla_s \mathbf{R}_T^{\text{FG},k+1}. \quad (6.31)$$



In this case, one is led to define the broken elliptic projection  $\mathcal{E}_{\mathcal{T}}^{\text{FG},k+1} : \mathbf{H}^1(\Omega) \rightarrow \mathbf{U}_{\mathcal{T}}^{k+1}$  such that, for all  $\mathbf{w} \in \mathbf{H}^1(\Omega)$  and all  $T \in \mathcal{T}_h$ ,  $\mathcal{E}_{\mathcal{T}}^{\text{FG},k+1}(\mathbf{w})|_T$  is uniquely defined by the relations

$$\left( \nabla(\mathcal{E}_T^{\text{FG},k+1}(\mathbf{w}) - \mathbf{w}), \nabla \mathbf{q} \right)_T = 0, \forall \mathbf{q} \in \mathbb{P}_d^{k+1,0}(T; \mathbb{R}^d), \quad (6.32a)$$

$$\int_T \{\mathcal{E}_{\mathcal{T}}^{\text{FG},k+1}(\mathbf{w}) - \mathbf{w}\} = \mathbf{0}, \quad (6.32b)$$

where  $\mathbb{P}_d^{k+1,0}(T; \mathbb{R}^d) := \{\mathbf{q} \in \mathbb{P}_d^{k+1}(T; \mathbb{R}^d) \mid \int_T \mathbf{q} = \mathbf{0}\}$  denotes the subset of polynomials in  $\mathbb{P}_d^{k+1}(T; \mathbb{R}^d)$  with zero translations. There is no available proof concerning stability. However, we did verify numerically the norm equivalence between  $\|\mathbf{E}_T^k(\cdot)\|_T^2 + \|h_T^{-\frac{1}{2}} \mathbf{S}_{\partial T}^{\text{SG}}(\cdot)\|_{\partial T}^2$  and  $\|\mathbf{E}_T^{\text{FG},k}(\cdot)\|_T^2 + \|h_T^{-\frac{1}{2}} \mathbf{S}_{\partial T}^{\text{FG}}(\cdot)\|_{\partial T}^2$  in the equal-order setting and between  $\|\mathbf{E}_T^k(\cdot)\|_T^2 + \|h_T^{-\frac{1}{2}} \mathbf{S}_{\partial T}^{\text{MO}}(\cdot)\|_{\partial T}^2$  and  $\|\mathbf{E}_T^{\text{FG},k}(\cdot)\|_T^2 + \|h_T^{-\frac{1}{2}} \mathbf{S}_{\partial T}^{\text{MO}}(\cdot)\|_{\partial T}^2$  in the mixed-order setting, in the same manner as for Conjecture 6.1.4. Specifically, we denote  $\mathcal{A}_T^{\text{FG},x}$  the matrix associated with the norm  $\|\mathbf{E}_T^{\text{FG},k}(\cdot)\|_T^2 + \|h_T^{-\frac{1}{2}} \mathbf{S}_{\partial T}^x(\cdot)\|_{\partial T}^2$ , where the double superscript refers to the full gradient reconstruction being used in  $\mathbf{E}_T^{\text{FG},k}$  and either  $\mathbf{S}_{TF}^{\text{FG},k}$  or  $\mathbf{S}_{TF}^{\text{MO},k}$  being used in the stabilization. Figure 6.2 displays the sorted eigenvalues of the generalized eigenproblem  $\mathcal{A}_T^{\text{FG},x} X = \lambda \mathcal{A}_T^{\text{SG}} X$  for polynomial orders  $k \in \{1, 2, 3\}$ , equal- and mixed-order settings, on a square and a right-triangular cell. In this case also, all the eigenvalues are finite and nonzero, which leads us to conjecture a stability property in the form

$$\alpha \|\hat{\mathbf{v}}_h\|_{\text{HHO}}^2 \leq \frac{1}{2} \|\mu^{\frac{1}{2}} \mathbf{E}_{\mathcal{T}}^{\text{FG},k}(\hat{\mathbf{v}}_h)\|_{\Omega}^2 + \|\mu^{\frac{1}{2}} h^{-\frac{1}{2}} \mathbf{S}_{\partial \mathcal{T}_h}^x(\hat{\mathbf{v}}_h)\|_{\partial \mathcal{T}_h}^2 \leq \varpi \|\hat{\mathbf{v}}_h\|_{\text{HHO}}^2, \quad \forall \hat{\mathbf{v}}_h \in \hat{\mathbf{U}}_{h,0}^{l,k}. \quad (6.33)$$

Unfortunately, as discussed in Remark 6.1.6, consistency is not ensured for this HHO discretization when

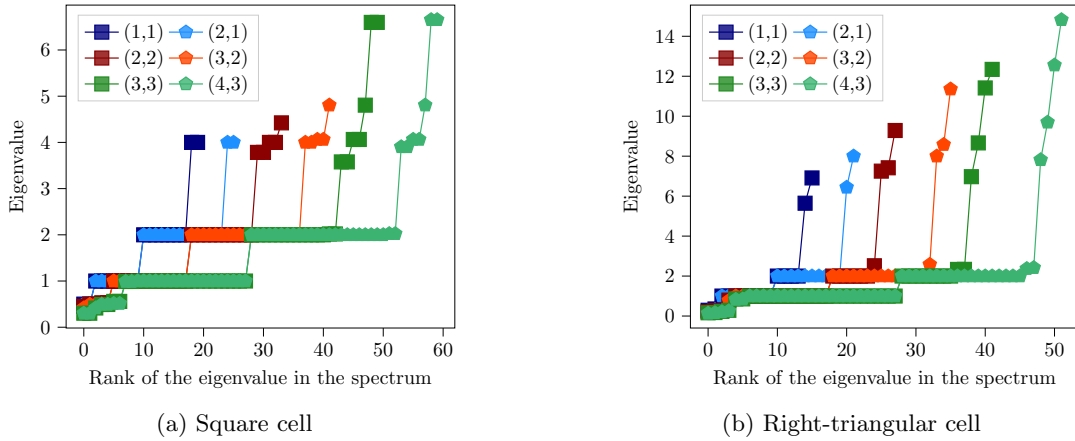


Figure 6.2: Linear elasticity - Eigenvalues of the generalized eigenproblem  $\mathcal{A}_T^{\text{FG},x} X = \lambda \mathcal{A}_T^{\text{SG}} X$ , equal- and mixed-order settings with polynomial orders  $k \in \{1, 2, 3\}$ , square and right-triangular cells.

considering polynomial orders  $k \geq 2$ .  $\square$

### 6.1.3 HHO global discretization

#### 6.1.3.1 Functional formulation.

Using the symmetric gradient reconstruction operator defined in (6.16) in each cell  $T \in \mathcal{T}_h$ , we define the local stiffness bilinear form  $b_T$ , such that, for all  $\hat{\mathbf{v}}_T, \hat{\mathbf{w}}_T \in \hat{\mathbf{U}}_T^{l,k}$ ,

$$b_T(\hat{\mathbf{v}}_T, \hat{\mathbf{w}}_T) := \left( \boldsymbol{\sigma}(\mathbf{E}_T^k(\hat{\mathbf{v}}_T)), \mathbf{E}_T^k(\hat{\mathbf{w}}_T) \right)_T. \quad (6.34)$$

The global bilinear forms  $b_h$  and  $s_h^x$  are defined, for all  $\hat{\mathbf{v}}_h, \hat{\mathbf{w}}_h \in \widehat{\mathbf{U}}_h^{l,k}$ , as

$$b_h(\hat{\mathbf{v}}_h, \hat{\mathbf{w}}_h) := \sum_{T \in \mathcal{T}_h} b_T(\hat{\mathbf{v}}_T, \hat{\mathbf{w}}_T), \quad s_h^x(\hat{\mathbf{v}}_h, \hat{\mathbf{w}}_h) := \sum_{T \in \mathcal{T}_h} s_T^x(\hat{\mathbf{v}}_T, \hat{\mathbf{w}}_T), \quad (6.35)$$

and the global stiffness bilinear form  $a_h^x$  is defined as

$$a_h^x(\hat{\mathbf{v}}_h, \hat{\mathbf{w}}_h) = b_h(\hat{\mathbf{v}}_h, \hat{\mathbf{w}}_h) + s_h^x(\hat{\mathbf{v}}_h, \hat{\mathbf{w}}_h), \quad (6.36)$$

where  $x = \text{MO}$  in the mixed-order setting and  $x \in \{\text{SG}, \text{FG}\}$  in the equal-order setting. Assuming  $\mathbf{f} \in C^0(\bar{J}, \mathbf{L}^2(\Omega))$ , the space semi-discrete scheme for (6.5) consists of finding  $\hat{\mathbf{u}}_h := (\mathbf{u}_{\mathcal{T}}, \mathbf{u}_{\mathcal{F}}) \in C^2(\bar{J}; \widehat{\mathbf{U}}_{h,0}^{l,k})$  such that, for all  $t \in \bar{J}$  and all  $\hat{\mathbf{w}}_h := (\mathbf{w}_{\mathcal{T}}, \mathbf{w}_{\mathcal{F}}) \in \widehat{\mathbf{U}}_{h,0}^{l,k}$ ,

$$(\rho \partial_t^2 \mathbf{u}_{\mathcal{T}}(t), \mathbf{w}_{\mathcal{T}})_{\Omega} + a_h^x(\hat{\mathbf{u}}_h(t), \hat{\mathbf{w}}_h) = (\mathbf{f}(t), \mathbf{w}_{\mathcal{T}})_{\Omega}. \quad (6.37)$$

Notice that the homogeneous Dirichlet boundary condition is enforced by the condition  $\hat{\mathbf{u}}_h(t) \in \widehat{\mathbf{U}}_{h,0}^{l,k}$  at all times  $t \in \bar{J}$ . The initial conditions are enforced on the cell degrees of freedom as follows:

$$\mathbf{u}_{\mathcal{T}}(0) := \Pi_{\mathcal{T}}^l(\mathbf{u}_0), \quad \partial_t \mathbf{u}_{\mathcal{T}}(0) := \Pi_{\mathcal{T}}^l(\mathbf{v}_0). \quad (6.38)$$

**Remark 6.1.6** (Consistency). Consistency is proved in the mixed-order setting [1, Lem. 6] and in the equal-order setting with  $x = \text{SG}$  [82, Thm. 8]. In the equal-order setting with  $x = \text{FG}$ , no proof is available. However, we verified numerically that the expected orders of convergence are indeed obtained. We consider the static equation  $a_h^{\text{FG}}(\hat{\mathbf{u}}_h, \hat{\mathbf{w}}_h) = (\mathbf{f}, \mathbf{w}_{\mathcal{T}})_{\Omega}$  with analytical solution

$$\mathbf{u}(x, y) := \phi(x, y), \quad \phi(x, y) := [-\sin(\pi x) \cos(\pi y), \cos(\pi x) \sin(\pi y)]^{\top}, \quad (6.39)$$

with source function  $\mathbf{f} := 2\pi^2 \mu \phi$ , initial conditions  $\mathbf{u}_0 := \mathbf{0}$ ,  $\mathbf{v}_0 := \mathbf{0}$  and reference material parameters. Figure 6.3 displays the  $L^2$ - and  $H^1$ -errors as a function of the mesh size for a series of refined square meshes, polynomial orders  $k \in \{1, 2, 3\}$  and equal-order setting. Convergence orders in  $h^{k+1}$  for the  $H^1$ -error and in  $h^{k+2}$  for the  $L^2$ -error are observed.  $\square$

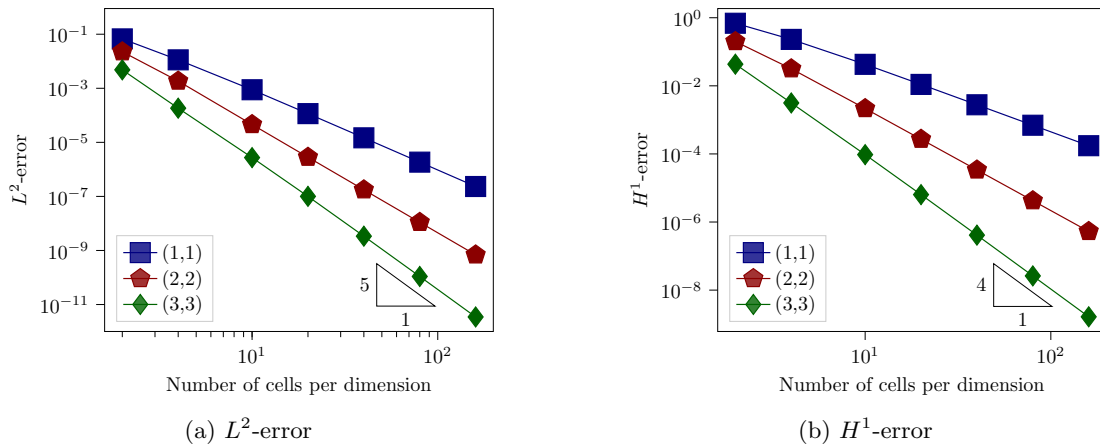


Figure 6.3: Linear elasticity (analytical solution for the static equation) -  $L^2$  and  $H^1$ -errors, full tensor stiffness and full gradient stabilization, polynomial orders  $k \in \{1, 2, 3\}$ , equal-order setting.

**Remark 6.1.7** (Variants). Following Remark 6.1.5, one could also consider the local bilinear form

$$b_T^{\text{FG}}(\hat{\mathbf{v}}_T, \hat{\mathbf{w}}_T) := \left( \boldsymbol{\sigma}(\mathbf{E}_T^{\text{FG},k}(\hat{\mathbf{v}}_T)), \mathbf{E}_T^{\text{FG},k}(\hat{\mathbf{w}}_T) \right)_T, \quad (6.40)$$

and define the global bilinear form  $a_h^{\text{FG},x} := \sum_{T \in \mathcal{T}_h} b_T^{\text{FG}} + s_h^x$ , with  $x = \text{MO}$  in the mixed-order setting and  $x = \text{FG}$  in the equal-order setting. However, our numerical experiments show that this HHO formulation lacks consistency for polynomial orders  $k \geq 2$ . Indeed, using the same analytical test case as before, we solve  $a_h^{\text{FG},x}(\hat{\mathbf{u}}_h, \hat{\mathbf{w}}_h) = (\mathbf{f}, \mathbf{w}_\mathcal{T})_\Omega$  and compute the  $L^2$ - and  $H^1$ -errors, reported in Figure 6.4 for polynomial orders  $k \in \{1, 2, 3\}$ , equal- and mixed-order settings. All polynomial orders show the same convergence rates: 2 for the  $H^1$ -error and 3 for the  $L^2$ -error. This means that for polynomial orders  $k \geq 2$ , the convergence rates are not optimal, but they are optimal for  $k = 1$  (a proof of this fact is still open). Considering these results, this HHO discretization is not used in what follows.  $\square$

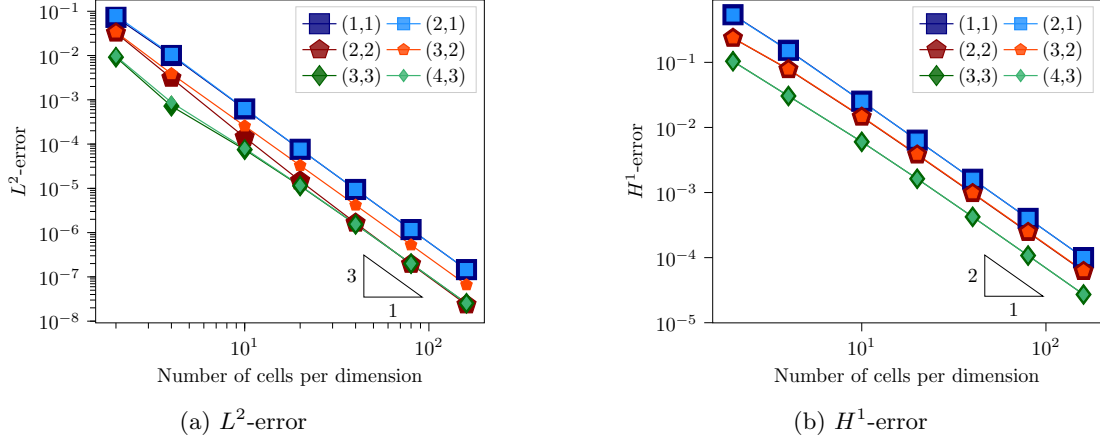


Figure 6.4: Linear elasticity (analytical solution for the static equation) -  $L^2$  and  $H^1$ -errors, full gradient stiffness, polynomial orders  $k \in \{1, 2, 3\}$ , equal- and mixed-order settings.

### 6.1.3.2 Algebraic formulation.

Let  $N_\mathcal{T} := \dim(\mathbf{U}_\mathcal{T}^l)$ ,  $N_\mathcal{F} := \dim(\mathbf{U}_{\mathcal{F},0}^k)$  and let  $\{\phi_i\}_{1 \leq i \leq N_\mathcal{T}}$ ,  $\{\psi_i\}_{1 \leq i \leq N_\mathcal{F}}$  be bases of  $\mathbf{U}_\mathcal{T}^l$  and  $\mathbf{U}_{\mathcal{F},0}^k$ , respectively. Let  $(\mathbf{U}_\mathcal{T}(t), \mathbf{U}_\mathcal{F}(t)) \in \mathbb{R}^{N_\mathcal{T}} \times \mathbb{R}^{N_\mathcal{F}}$  be the vector of the (time-dependent) degrees of freedom of the space semi-discrete solution  $\hat{\mathbf{u}}_h(t)$  on these bases, and  $\mathbf{F}_\mathcal{T}(t) \in \mathbb{R}^{N_\mathcal{T}}$  the vector having components  $((\mathbf{f}(t), \phi_i)_\Omega)_{1 \leq i \leq N_\mathcal{T}}$ . Let also  $\mathcal{M}$  denote the mass matrix of the vector-valued displacement on the cells (weighted by the density  $\rho$ ),  $\mathbf{u}_\mathcal{T}(t)$ , and let  $\mathcal{A}^x$  denote the stiffness matrix associated with the bilinear form  $a_h^x$ . We have  $\mathcal{A}^x = \mathcal{B} + \mathcal{S}^x$ , where  $\mathcal{B}$  is associated with the stiffness bilinear form  $b_h$  and  $\mathcal{S}^x$  with the stabilization bilinear form  $s_h^x$ . The algebraic formulation of (6.37) writes

$$\begin{bmatrix} \mathcal{M} & 0 \\ 0 & 0 \end{bmatrix} \begin{pmatrix} \partial_t^2 \mathbf{U}(t) \\ \cdot \end{pmatrix} + \begin{bmatrix} \mathcal{A}_{\mathcal{T}\mathcal{T}}^x & \mathcal{A}_{\mathcal{T}\mathcal{F}}^x \\ \mathcal{A}_{\mathcal{F}\mathcal{T}}^x & \mathcal{A}_{\mathcal{F}\mathcal{F}}^x \end{bmatrix} \begin{pmatrix} \mathbf{U}_\mathcal{T}(t) \\ \mathbf{U}_\mathcal{F}(t) \end{pmatrix} = \begin{bmatrix} \mathbf{F}_\mathcal{T}(t) \\ 0 \end{bmatrix}. \quad (6.41)$$

### 6.1.3.3 Time discretization with the leapfrog scheme.

We recall the following definitions: let  $N$  be the number of discrete time intervals such that  $(t^n)_{n \in \{0:N\}}$  are the discrete time nodes with  $t^0 = 0$  and  $t^N := \mathfrak{T}$ . For the sake of simplicity, we consider a fixed time step  $\Delta t := \frac{\mathfrak{T}}{N}$ . The time discrete unknown  $\hat{\mathbf{u}}_h^n = (\mathbf{u}_\mathcal{T}^n, \mathbf{u}_\mathcal{F}^n) \in \hat{\mathbf{U}}_{h,0}^{l,k}$  is meant to be an approximation of the space semi-discrete HHO solution  $\hat{\mathbf{u}}_h(t^n) \in \hat{\mathbf{U}}_{h,0}^{l,k}$ . We set  $\mathbf{f}^n := \mathbf{f}(t^n)$  for all  $n \in \{0:N\}$ . The leapfrog scheme consists of solving, for all  $n \in \{1:N-1\}$ ,

$$\frac{1}{\Delta t^2} (\rho(\mathbf{u}_\mathcal{T}^{n+1} - 2\mathbf{u}_\mathcal{T}^n + \mathbf{u}_\mathcal{T}^{n-1}), \mathbf{w}_\mathcal{T})_\Omega + a_h^x(\hat{\mathbf{u}}_h^n, \hat{\mathbf{w}}_h) = (\mathbf{f}^n, \mathbf{w}_\mathcal{T})_\Omega, \quad \forall \hat{\mathbf{w}}_h \in \hat{\mathbf{U}}_{h,0}^{l,k}, \quad (6.42)$$

where the unknowns are  $\mathbf{u}_{\mathcal{T}}^{n+1}$  and  $\mathbf{u}_{\mathcal{F}}^n$ , whereas  $\mathbf{u}_{\mathcal{T}}^n$  and  $\mathbf{u}_{\mathcal{F}}^{n-1}$  are known from prior time steps or given by the initial conditions as follows:

$$\mathbf{u}_{\mathcal{T}}^0 = \Pi_{\mathcal{T}}^l(\mathbf{u}_0), \quad (6.43a)$$

$$a_h^x(\hat{\mathbf{u}}_h^0, (\mathbf{0}, \mathbf{w}_{\mathcal{F}})) = 0, \quad \forall \mathbf{w}_{\mathcal{F}} \in \mathcal{U}_{\mathcal{F},0}^k, \quad (6.43b)$$

$$(\rho \mathbf{u}_{\mathcal{T}}^1, \mathbf{w}_{\mathcal{T}})_{\Omega} = (\rho(\mathbf{u}_{\mathcal{T}}^0 + \Delta t \Pi_{\mathcal{T}}^l(\mathbf{v}_0)), \mathbf{w}_{\mathcal{T}})_{\Omega} + \frac{\Delta t^2}{2} [(\mathbf{f}^0, \mathbf{w}_{\mathcal{T}})_{\Omega} - a_h^x(\hat{\mathbf{u}}_h^0, (\mathbf{w}_{\mathcal{T}}, \mathbf{0}))], \quad \forall \mathbf{w}_{\mathcal{T}} \in \mathcal{U}_{\mathcal{T}}^l. \quad (6.43c)$$

We notice from (6.42) that the leapfrog scheme yields a semi-implicit scheme with the following static coupling between the cell and the face unknowns:

$$a_h^x((\mathbf{0}, \mathbf{u}_{\mathcal{F}}^n), (\mathbf{0}, \mathbf{w}_{\mathcal{F}})) = -a_h^x((\mathbf{u}_{\mathcal{T}}^n, \mathbf{0}), (\mathbf{0}, \mathbf{w}_{\mathcal{F}})), \quad \forall \mathbf{w}_{\mathcal{F}} \in \mathcal{U}_{\mathcal{F},0}^k. \quad (6.44)$$

With  $\mathbf{U}_{\mathcal{T}}^n$  and  $\mathbf{U}_{\mathcal{F}}^n$  respectively collecting the degrees of freedom of  $\mathbf{u}_{\mathcal{T}}^n$  and  $\mathbf{u}_{\mathcal{F}}^n$  for all  $n \in \{1:N-1\}$ , (6.42) translates into the algebraic equation

$$\frac{1}{\Delta t^2} \begin{bmatrix} \mathcal{M} & 0 \\ 0 & 0 \end{bmatrix} \begin{pmatrix} \mathbf{U}_{\mathcal{T}}^{n+1} - 2\mathbf{U}_{\mathcal{T}}^n + \mathbf{U}_{\mathcal{T}}^{n-1} \\ \cdot \end{pmatrix} + \begin{bmatrix} \mathcal{A}_{\mathcal{T}\mathcal{T}}^x & \mathcal{A}_{\mathcal{T}\mathcal{F}}^x \\ \mathcal{A}_{\mathcal{F}\mathcal{T}}^x & \mathcal{A}_{\mathcal{F}\mathcal{F}}^x \end{bmatrix} \begin{pmatrix} \mathbf{U}_{\mathcal{T}}^n \\ \mathbf{U}_{\mathcal{F}}^n \end{pmatrix} = \begin{bmatrix} \mathbf{F}_{\mathcal{T}}^n \\ 0 \end{bmatrix}. \quad (6.45)$$

The static coupling (6.44) rewrites at each time step  $n$  as

$$\mathcal{A}_{\mathcal{F}\mathcal{F}}^x \mathbf{U}_{\mathcal{F}}^n = -\mathcal{A}_{\mathcal{F}\mathcal{T}}^x \mathbf{U}_{\mathcal{T}}^n. \quad (6.46)$$

## 6.2 Splitting for the elastic wave equation

In this section, we adapt the splitting introduced for the acoustic wave equation in Section 4.1 to the elastic case. The devising of the splitting procedure does not change. It still relies on the block-diagonal face-face stabilization matrix in the mixed-order setting and on a block-diagonal part of the face-face stabilization matrix in the equal-order setting. The difference lays in the value of the stiffness and stabilization operators, since they use a symmetric gradient reconstruction operator instead of a plain gradient one.

### 6.2.1 Design of the splitting

#### 6.2.1.1 Splitting in the mixed-order setting

In the same manner as in Chapter 4, the splitting proceeds as follows: for all  $n \in \{1:N-1\}$ , set  $\mathbf{u}_{\mathcal{F}}^{n,0} = \mathbf{u}_{\mathcal{F}}^{n-1}$  and iterate on  $m \geq 0$  by finding  $\mathbf{u}_{\mathcal{F}}^{n,m+1} \in \mathcal{U}_{\mathcal{F},0}^k$  such that

$$s_h^{\text{MO}}((\mathbf{0}, \mathbf{u}_{\mathcal{F}}^{n,m+1}), (\mathbf{0}, \mathbf{w}_{\mathcal{F}})) = -b_h((\mathbf{0}, \mathbf{u}_{\mathcal{F}}^{n,m}), (\mathbf{0}, \mathbf{w}_{\mathcal{F}})) - a_h^{\text{MO}}((\mathbf{u}_{\mathcal{T}}^n, \mathbf{0}), (\mathbf{0}, \mathbf{w}_{\mathcal{F}})), \quad \forall \mathbf{w}_{\mathcal{F}} \in \mathcal{U}_{\mathcal{F},0}^k. \quad (6.47)$$

The algebraic form is as follows: Setting  $\mathbf{U}_{\mathcal{F}}^{n,0} := \mathbf{U}_{\mathcal{F}}^{n-1}$ , we seek  $\mathbf{U}_{\mathcal{F}}^{n,m+1} \in \mathbb{R}^{N_{\mathcal{F}}}$  such that

$$\mathcal{S}_{\mathcal{F}\mathcal{F}}^{\text{MO}} \mathbf{U}_{\mathcal{F}}^{n,m+1} = -\mathcal{B}_{\mathcal{F}\mathcal{F}} \mathbf{U}_{\mathcal{F}}^{n,m} - \mathcal{A}_{\mathcal{F}\mathcal{T}}^{\text{MO}} \mathbf{U}_{\mathcal{T}}^n. \quad (6.48)$$

This splitting procedure is computationally effective since, as in the acoustic wave case, the face-face stabilization submatrix  $\mathcal{S}_{\mathcal{F}\mathcal{F}}^{\text{MO}}$  is block-diagonal. This splitting procedure is a fixed-point algorithm using the affine function  $\mathbf{g} : \mathbf{X} \rightarrow (\mathcal{S}_{\mathcal{F}\mathcal{F}}^{\text{MO}})^{-1}(\mathcal{B}_{\mathcal{F}\mathcal{F}} \mathbf{X} - \mathcal{A}_{\mathcal{F}\mathcal{T}}^{\text{MO}} \mathbf{U}_{\mathcal{T}}^n)$  for all  $\mathbf{X} \in \mathbb{R}^{N_{\mathcal{F}}}$ . Its convergence is thus ensured for  $\delta < 1$ , where  $\delta$  is the Lipschitz constant of  $\mathbf{g}$ . Since  $\mathbf{g}$  is an affine map, this Lipschitz constant is found by computing the spectral radius of  $(\mathcal{S}_{\mathcal{F}\mathcal{F}}^{\text{MO}})^{-1} \mathcal{B}_{\mathcal{F}\mathcal{F}}$  and the convergence criterion is

$$\rho((\mathcal{S}_{\mathcal{F}\mathcal{F}}^{\text{MO}})^{-1} \mathcal{B}_{\mathcal{F}\mathcal{F}}) < 1. \quad (6.49)$$

Using the same arguments as for the acoustic wave equation, this condition can be satisfied by scaling the stabilization  $s_h^{\text{MO}}$  by  $\gamma$  large enough. The values of  $\gamma$  leading to a converging splitting procedure are studied in Section 6.2.2.

### 6.2.1.2 Splitting in the equal-order setting

This case requires some additional definitions, as in the scalar-valued case. Specifically, let  $\zeta_T^x$  be the local bilinear form such that for all  $T \in \mathcal{T}_h$ , and all  $\mathbf{v}_{\partial T}, \mathbf{w}_{\partial T} \in \mathcal{U}_{\partial T}^k$ ,

$$\begin{aligned} \zeta_T^x((\mathbf{0}, \mathbf{v}_{\partial T}), (\mathbf{0}, \mathbf{w}_{\partial T})) &= \gamma \mu_T^2 \sum_{F \in \mathcal{F}_T} \eta_{TF}^{-1} \left\{ \left( (I - \Pi_T^k) \mathbf{R}_T^{x, k+1}(\mathbf{0}, \mathbf{v}_{\partial T})|_F, \mathbf{w}_F \right)_F \right. \\ &\quad + \left( \mathbf{v}_F, (I - \Pi_T^k) \mathbf{R}_T^{x, k+1}(\mathbf{0}, \mathbf{w}_{\partial T})|_F \right)_F \\ &\quad \left. + \left( \Pi_F^k (I - \Pi_T^k) \mathbf{R}_T^{x, k+1}(\mathbf{0}, \mathbf{v}_{\partial T})|_F, \Pi_F^k (I - \Pi_T^k) \mathbf{R}_T^{x, k+1}(\mathbf{0}, \mathbf{w}_{\partial T})|_F \right)_F \right\}, \end{aligned} \quad (6.50)$$

and let  $s_T^*$  be the local bilinear form defined by

$$s_T^*((\mathbf{0}, \mathbf{v}_{\partial T}), (\mathbf{0}, \mathbf{w}_{\partial T})) := \gamma \mu_T^2 \sum_{F \in \mathcal{F}_T} \eta_{TF}^{-1} (\mathbf{v}_F, \mathbf{w}_F)_F. \quad (6.51)$$

Then the equal-order stabilization bilinear form writes

$$s_T^x = s_T^* + \zeta_T^x. \quad (6.52)$$

Let us introduce the global bilinear forms  $s_h^*((\mathbf{0}, \mathbf{v}_F), (\mathbf{0}, \mathbf{w}_F)) := \sum_{T \in \mathcal{T}_h} s_T^*((\mathbf{0}, \mathbf{v}_{\partial T}), (\mathbf{0}, \mathbf{w}_{\partial T}))$  and  $\zeta_h^x((\mathbf{0}, \mathbf{v}_F), (\mathbf{0}, \mathbf{w}_F)) := \sum_{T \in \mathcal{T}_h} \zeta_T^x((\mathbf{0}, \mathbf{v}_{\partial T}), (\mathbf{0}, \mathbf{w}_{\partial T}))$ , so that  $s_h^x = s_h^* + \zeta_h^x$ . This leads to the following iterative procedure, with the same initial condition as for the mixed-order setting: For all  $m \geq 0$ , find  $\mathbf{u}_{\mathcal{F}}^{n, m+1} \in \mathcal{U}_{\mathcal{F}, 0}^k$  such that

$$\begin{aligned} s_h^*((\mathbf{0}, \mathbf{u}_{\mathcal{F}}^{n, m+1}), (\mathbf{0}, \mathbf{w}_{\mathcal{F}})) &= -b_h((\mathbf{u}_{\mathcal{T}}^n, \mathbf{u}_{\mathcal{F}}^{n, m}), (\mathbf{0}, \mathbf{w}_{\mathcal{F}})) \\ &\quad - \zeta_h^x((\mathbf{0}, \mathbf{u}_{\mathcal{F}}^{n, m}), (\mathbf{0}, \mathbf{w}_{\mathcal{F}})) - s_h^x((\mathbf{u}_{\mathcal{T}}^n, \mathbf{0}), (\mathbf{0}, \mathbf{w}_{\mathcal{F}})), \quad \forall \mathbf{w}_{\mathcal{F}} \in \mathcal{U}_{\mathcal{F}, 0}^k. \end{aligned} \quad (6.53)$$

At the algebraic level, we define two matrices  $\mathcal{S}_{\mathcal{F}\mathcal{F}}^*$  and  $\mathcal{Z}_{\mathcal{F}\mathcal{F}}^x$  such that  $\mathcal{S}_{\mathcal{F}\mathcal{F}}^x = \mathcal{S}_{\mathcal{F}\mathcal{F}}^* + \mathcal{Z}_{\mathcal{F}\mathcal{F}}^x$ ,  $\mathcal{Z}_{\mathcal{F}\mathcal{F}}^x$  corresponds to the bilinear form  $\zeta_h^x(\cdot, \cdot)$  and  $\mathcal{S}_{\mathcal{F}\mathcal{F}}^*$  to  $s_h^*(\cdot, \cdot)$ . Then the splitting procedure translates into the following iterative algorithm: For all  $m \geq 0$ , find  $\mathbf{U}_{\mathcal{F}}^{n, m+1} \in \mathbb{R}^{N_{\mathcal{F}}}$  such that

$$\mathcal{S}_{\mathcal{F}\mathcal{F}}^* \mathbf{U}_{\mathcal{F}}^{n, m+1} = -\mathcal{B}_{\mathcal{F}\mathcal{F}} \mathbf{U}_{\mathcal{F}}^{n, m} - \mathcal{Z}_{\mathcal{F}\mathcal{F}}^x \mathbf{U}_{\mathcal{F}}^{n, m} - \mathcal{A}_{\mathcal{F}\mathcal{T}}^x \mathbf{U}_{\mathcal{T}}^n. \quad (6.54)$$

The convergence condition is

$$\rho((\mathcal{S}_{\mathcal{F}\mathcal{F}}^*)^{-1} (\mathcal{B}_{\mathcal{F}\mathcal{F}} + \mathcal{Z}_{\mathcal{F}\mathcal{F}}^x)) < 1, \quad (6.55)$$

which can be achieved by using  $\gamma$  large enough under the additional condition

$$\alpha := \rho((\mathcal{S}_{\mathcal{F}\mathcal{F}}^*)^{-1} \mathcal{Z}_{\mathcal{F}\mathcal{F}}^x) < 1. \quad (6.56)$$

The admissible values of  $\gamma$  are studied in Section 6.2.2.

### 6.2.1.3 Verification of the condition $\alpha < 1$ (equal-order setting)

**Full gradient reconstruction.** For  $x = \text{FG}$ , the stabilization is the vector-valued version of the stabilization used for the acoustic wave equation. Since the eigenvalues of  $(\mathcal{S}_{\mathcal{F}\mathcal{F}}^*)^{-1} \mathcal{Z}_{\mathcal{F}\mathcal{F}}^{\text{FG}}$  remain unchanged, the condition  $\alpha < 1$  is ensured for the full gradient reconstruction for all the mesh shapes and polynomial orders considered in Chapter 4.

Order $(k, l)$	(1,1)	(2,2)	(3,3)
Squares	$\approx 1.0$	0.6	0.74
Right triangles	1.1	1.25	0.98
Unstructured triangles	1.1	0.93	0.77
Unstructured quadrangles	1.32	0.98	1.03

Table 6.2: Linear elasticity (reference material) -  $\alpha$ , the spectral radius of  $(\mathcal{S}_{\mathcal{F}\mathcal{F}}^*)^{-1}\mathcal{Z}_{\mathcal{F}\mathcal{F}}^{\text{SG}}$ ,  $k \in \{1, 2, 3\}$ , equal-order setting.

**Symmetric gradient reconstruction.** For the symmetric gradient reconstruction  $x = \text{SG}$ , on the other hand, a new study on the value of  $\alpha$  must be carried out. Table 6.2 displays the value of  $\alpha$  for  $k \in \{1, 2, 3\}$  on the domain  $\Omega := (0, 1)^2$  for the reference material ( $\mu = \lambda = 1$ ). We consider a mesh size  $h := 0.1$  and four different cell shapes: squares, right triangles (square mesh with all cells divided along the diagonal), general quadrangles and triangles, both belonging to shape-regular sequences of unstructured meshes. The triangular meshes are generated using `gmsh`, and the quadrangular meshes are also created using `gmsh` by merging pairs of adjacent triangles.

As seen in Table 6.2, the value of  $\alpha$  is not always smaller than one. This means that the splitting may not converge for these cases. It could, however, be possible to find a range of values for  $\gamma$  for which the spectral radius of  $(\mathcal{S}_{\mathcal{F}\mathcal{F}}^*)^{-1}(\mathcal{B}_{\mathcal{F}\mathcal{F}} + \mathcal{Z}_{\mathcal{F}\mathcal{F}}^{\text{SG}})$  is smaller than 1, provided the eigenvalues of  $\mathcal{B}_{\mathcal{F}\mathcal{F}}$  and  $\mathcal{Z}_{\mathcal{F}\mathcal{F}}^{\text{SG}}$  compensate. But as  $\gamma \rightarrow \infty$ , the spectral radius of  $(\mathcal{S}_{\mathcal{F}\mathcal{F}}^*)^{-1}(\mathcal{B}_{\mathcal{F}\mathcal{F}} + \mathcal{Z}_{\mathcal{F}\mathcal{F}}^{\text{SG}})$  converges to that of  $(\mathcal{S}_{\mathcal{F}\mathcal{F}}^*)^{-1}\mathcal{Z}_{\mathcal{F}\mathcal{F}}^{\text{SG}}$ ; hence, there is a value  $\gamma^{\text{MAX}}$  such that, for all  $\gamma > \gamma^{\text{MAX}}$ , the splitting does not converge. To illustrate, Figure 6.5 displays the spectral radius of  $(\mathcal{S}_{\mathcal{F}\mathcal{F}}^*)^{-1}(\mathcal{B}_{\mathcal{F}\mathcal{F}} + \mathcal{Z}_{\mathcal{F}\mathcal{F}}^{\text{SG}})$  as a function of  $\gamma$  for  $k \in \{1, 2, 3\}$  for the same meshes and polynomial orders as those considered in Table 6.2. In all the cases where  $\alpha > 1$ , there is no value of  $\gamma$  such that the splitting converges.

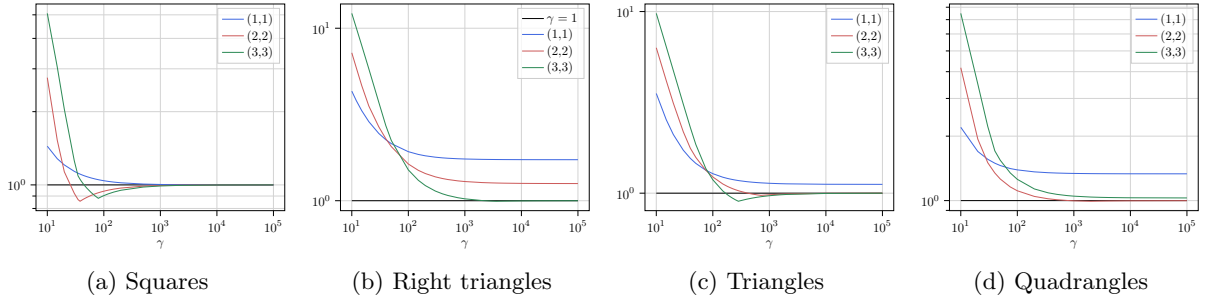


Figure 6.5: Linear elasticity (reference material) - Spectral radius of  $(\mathcal{S}_{\mathcal{F}\mathcal{F}}^*)^{-1}(\mathcal{B}_{\mathcal{F}\mathcal{F}} + \mathcal{Z}_{\mathcal{F}\mathcal{F}}^{\text{SG}})$ ,  $k \in \{1, 2, 3\}$ , equal-order setting,  $h = 0.1$ .

**Conclusion on the choice of the type of reconstruction.** The above results lead us to consider only the full gradient reconstruction in the equal-order setting in all the rest of this Chapter. The mixed-order setting is unaffected since the stabilization does not use the reconstructed potential. Since for each polynomial order, only one reconstruction is considered in the following, we drop the superscript  $x$ . We denote the stabilization  $s_h^{\text{MO}}$  in the mixed-order setting and the stabilization  $s_h^{\text{FG}}$  in the equal-order setting by the generic notation  $s_h$ , and proceed in the same manner for the associated matrices.

## 6.2.2 Numerical study on the stability parameter

As for the acoustic wave equation, we define  $\gamma^*$  as the smallest value of  $\gamma$  such that the splitting procedure converges. In this section, we consider the factors that may influence the value of  $\gamma^*$  and study them

numerically. The study on the linear elasticity equation also gives some insights on the expected behavior for more complicated equations, for instance when considering finite-strain plasticity.

The parameters influencing the value of  $\gamma^*$  and the CFL condition are the polynomial order and the mesh regularity. These parameters have already been investigated in the case of the acoustic wave equation and are expected to produce a similar impact for linear elasticity. We recall that, in the case of the acoustic wave equation, the value of  $\gamma^*$  increases with the polynomial order and is somewhat larger on unstructured meshes than on structured meshes. In addition, we observed that downgrading the mesh regularity tightens the CFL condition. These are the expected behaviors for elastodynamics as well.

### 6.2.2.1 Convergence of $\gamma^*$ with the mesh size

Figure 6.6 displays the value of  $\gamma^*$  on a mesh sequence of iteratively refined right triangles with polynomial orders  $k \in \{1, 2, 3\}$  in the equal- and mixed-order settings. We enforce homogeneous Dirichlet boundary conditions on the top and bottom edges of  $\Omega$ . This figure also displays via horizontal lines the value of  $\gamma^*$  on a single right-triangular mesh cell without boundary conditions. We first remark that, as for the acoustic wave equation,  $\gamma^*$  increases with the polynomial order. We also notice that the value of  $\gamma^*$  converges towards the value on a single cell without boundary conditions as the mesh is refined. This illustrates that the choice of the boundary conditions does not impact the value of  $\gamma^*$  on refined meshes. Comparing the equal- and mixed-order settings, the difference of the values for  $\gamma^*$  is more pronounced than for the acoustic wave equation (see Figure 4.1), but remains smaller than the differences caused by varying the polynomial orders.

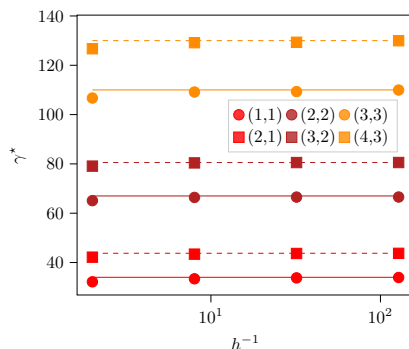


Figure 6.6: Linear elasticity (reference material) - Spectral radius as a function of the mesh size with right-triangular cells, and homogeneous Dirichlet conditions in 2D, compared to the reference value without boundary conditions,  $k \in \{1, 2, 3\}$ , equal- and mixed-order settings.

### 6.2.2.2 Impact of mesh regularity on $\gamma^*$ .

We observed in Chapter 4 that the mesh regularity and the shape of the cells impact the value of  $\gamma^*$ . Table 6.3 reports the value of  $\gamma^*$  for a square cell, a right-isocenes triangular cell as well as estimates for general quadrangular and triangular cells belonging to shape-regular sequences of unstructured meshes, still for reference material parameters, polynomial orders  $k \in \{1, 2, 3\}$  and equal- and mixed-order settings. For the two latter meshes, the reported value of  $\gamma^*$  is the largest observed value, rounded to the above integer on all the mesh cells and for all the meshes in the sequence. We notice again that the value of  $\gamma^*$  increases with the polynomial order. Moreover, for a given polynomial order, the smallest value of  $\gamma^*$  is obtained on square cells, and the largest value on right-triangular meshes. Once again, the value of  $\gamma^*$  for a given value of  $k$  is smaller in the equal-order setting than in the mixed-order setting, and the difference is larger than for the acoustic wave equation.

Order $(l, k)$	(1,1)	(2,2)	(3,3)	(2,1)	(3,2)	(4,3)
Squares	10	20	39	19	38	61
Right triangles	34	67	110	44	79	129
Unstructured triangles	27	50	84	36	64	103
Unstructured quadrangles	18	40	61	24	48	81

Table 6.3: Linear elasticity (reference material) -  $\gamma^*$  computed via the spectral radius of  $(\mathbf{S}_{\mathcal{FF}}^*)^{-1}(\mathbf{B}_{\mathcal{FF}} + \mathbf{Z}_{\mathcal{FF}})$  in the equal-order setting and of  $\mathbf{S}_{\mathcal{FF}}^{-1}\mathbf{B}_{\mathcal{FF}}$  in the mixed-order setting, on a single cell (first two lines) and on a shape-regular mesh sequence (last two lines),  $k \in \{1, 2, 3\}$ .

### 6.2.2.3 Impact of mesh regularity on the CFL.

We recall the result on the stability of the leapfrog scheme, using

$$\mathcal{D}(\gamma) := \mathcal{M}^{-1}(\mathcal{A}_{\mathcal{TT}}(\gamma) - \mathcal{A}_{\mathcal{TF}}(\gamma)\mathcal{A}_{\mathcal{FF}}(\gamma)^{-1}\mathcal{A}_{\mathcal{FT}}(\gamma)), \quad (6.57)$$

where  $\mathcal{A}(\gamma)$  is the stiffness matrix defined in (6.41) with the dependence on  $\gamma$  made explicit. The stability condition on  $\Delta t$  then reads

$$\Delta t(\gamma) \leq \Delta t^{\text{opt}}(\gamma) := \frac{2}{\sqrt{\rho(\mathcal{D}(\gamma))}}. \quad (6.58)$$

Table 6.4 illustrates the tightening of the CFL condition by displaying the ratio  $\Delta t^{\text{opt}}(\gamma^*)/\Delta t^{\text{opt}}(1)$  for  $k \in \{1, 2, 3\}$  on the meshes with  $h = 0.1$  from the mesh sequences considered in the previous experiments, still using the reference material. As for the acoustic wave equation, we remark that increasing  $\gamma^*$  affects negatively the CFL condition. The optimal time step is reduced by a factor of  $\simeq 2$  on squares meshes and by a factor between 2 and 4 on the other meshes. Increasing the polynomial order slightly alleviates the tightening of the stability condition. Since the value of  $\gamma^*$  is smaller in the equal-order setting than in the mixed-order setting for a given value of  $k$ , so is the impact on the stability condition. This difference is particularly acute on square and right-triangular cells.

Order $(k, l)$	(1,1)	(2,2)	(3,3)	(2,1)	(3,2)	(4,3)
Squares	0.61	0.74	0.72	0.38	0.44	0.47
Right triangles	0.37	0.47	0.55	0.17	0.23	0.32
Unstructured triangles	0.34	0.42	0.53	0.25	0.26	0.32
Unstructured quadrangles	0.29	0.35	0.44	0.28	0.30	0.36

Table 6.4: Linear elasticity (reference material) -  $\Delta t^{\text{opt}}(\gamma^*)/\Delta t^{\text{opt}}(1)$ ,  $h = 0.1$ ,  $k \in \{1, 2, 3\}$ , equal- and mixed-order setting.

### 6.2.2.4 Impact of material parameters

The ratio of the Lamé coefficients in the stiffness matrix plays an important role. The first observation is that, in the incompressible limit  $\nu \uparrow \frac{1}{2}$ , we have  $\frac{\lambda}{\mu} \rightarrow \infty$ . Since  $c_P := \sqrt{\frac{\lambda+2\mu}{\rho}}$  and  $c_S := \sqrt{\frac{\mu}{\rho}}$ , this means that, as  $\nu \uparrow \frac{1}{2}$ , the velocity of P-waves,  $c_P$ , becomes much larger than the velocity of S-waves,  $c_S$ . The main consequence is that the CFL condition imposes that  $\Delta t \rightarrow 0$ . This is the reason why materials with values of  $\nu$  too close to  $\frac{1}{2}$  are challenging, and we leave this difficulty, which is not specific to the HHO method, to future work.

Table 6.5 displays the value of  $\gamma^*$  on right triangular meshes for  $k \in \{1, 2\}$  in the equal- and mixed-order settings for all the materials considered in Table 6.1. We remark that the value of  $\gamma^*$  is indeed impacted by the value of the Poisson coefficient. Indeed,  $\gamma^*$  is much larger for gold and rubber, which have larger values of  $\nu$ . Instead, for steel and copper, which have a similar value for  $\nu$ , the values of  $\gamma^*$  are



similar. Diamond has the smallest value of  $\nu$  and also the smallest value of  $\gamma^*$  for all polynomial orders. The reference material does not behave differently, since its Poisson ratio is close to the Poisson ratio of steel. Altogether, this illustrates the fact that the Poisson coefficient is the main material parameter influencing the value of  $\gamma^*$ . This test also shows that, when considering usual metals such as steel and copper, the value of  $\gamma^*$  computed on steel is a good reference, but is not accurate enough to be used for all the materials since it could be underestimated. In the forthcoming experiments in this chapter, we will indicate the type of material used, and recompute  $\gamma^*$  if necessary.

Material	Steel	Copper	Diamond	Gold	Rubber	Reference
(1,1)	45	50	30	85	6083	34
(2,2)	89	98	60	167	11949	67
(2,1)	59	65	39	113	8374	43
(3,2)	108	119	72	207	15150	81

Table 6.5: Linear elasticity -  $\gamma^*$  computed via the spectral radius of  $(\mathcal{S}_{\mathcal{FF}}^*)^{-1}(\mathcal{B}_{\mathcal{FF}} + \mathcal{Z}_{\mathcal{FF}})$  in the equal-order setting and of  $\mathcal{S}_{\mathcal{FF}}^{-1}\mathcal{B}_{\mathcal{FF}}$  in the mixed-order setting, on a single right-triangular cell,  $k \in \{1, 2\}$ .

Figure 6.7 shows the value of  $\gamma^*$  as a function of  $\nu \in [0; 0.5)$  on a right-triangular cell for a Young modulus  $E = 2.5$ , polynomial orders  $k \in \{1, 2\}$  and equal- and mixed-order settings. The results confirm that  $\gamma^*$  grows unboundedly as  $\nu \uparrow \frac{1}{2}$ .

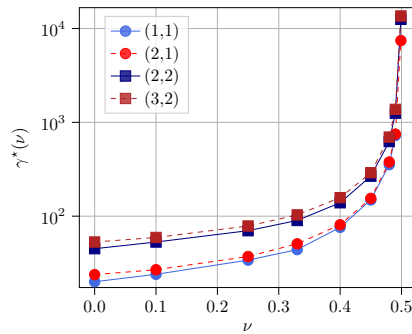


Figure 6.7: Linear elasticity -  $\gamma^*$ , computed via the spectral radius of  $(\mathcal{S}_{\mathcal{FF}}^*)^{-1}(\mathcal{B}_{\mathcal{FF}} + \mathcal{Z}_{\mathcal{FF}})$  in the equal-order setting and of  $\mathcal{S}_{\mathcal{FF}}^{-1}\mathcal{B}_{\mathcal{FF}}$  in the mixed-order setting, as a function of the Poisson ratio  $\nu$ ,  $k \in \{1, 2\}$ ,  $h = 0.1$ , right-triangular cell.

### 6.2.2.5 3D meshes

The previous experiments have been performed on 2D meshes. The same experiments can be done on 3D meshes. We only consider Cartesian hexahedra and unstructured tetrahedra with  $k = 1$ , in the equal- and mixed-order settings. Table 6.6 displays the value of  $\gamma^*$ , the ratio  $\Delta t^{\text{opt}}(\gamma^*)/\Delta t^{\text{opt}}(1)$  for the reference material on meshes with  $h = 0.2$ , and with the two reconstruction operators. The value of  $\gamma^*$  on hexahedra is comparable to the one on square meshes for the same polynomial order, whereas it is three times larger on tetrahedra than on triangles. Moreover, the tightening of the CFL condition is more pronounced on 3D meshes than on 2D meshes. Altogether, the results on 3D meshes seem to indicate that hexahedral meshes may lead to better performances than tetrahedral meshes when using the splitting procedure, especially in the equal-order setting.

$(k, l)$	(1,1)	(2,1)
$\gamma^*$ <b>hexahedra</b>	16	34
$\Delta t^{\text{opt}}(\gamma^*)/\Delta t^{\text{opt}}(1)$ <b>hexahedra</b>	0.41	0.21
$\gamma^*$ <b>tetrahedra</b>	101	123
$\Delta t^{\text{opt}}(\gamma^*)/\Delta t^{\text{opt}}(1)$ <b>tetrahedra</b>	0.21	0.16

Table 6.6: Linear elasticity (reference material) - Value of  $\gamma^*$  and ratio  $\Delta t^{\text{opt}}(\gamma^*)/\Delta t^{\text{opt}}(1)$  on 3D meshes (Cartesian hexahedra and unstructured tetrahedra),  $h = 0.2$ ,  $k = 1$  equal and mixed-order settings.

### 6.3 Numerical experiments on the elastic wave equation

This section presents numerical experiments using the splitting procedure for the elastic wave equation. We start by verifying the orders of convergence in space and in time. We then investigate the optimal choice of the splitting parameter  $\gamma$  in terms of error and computational time. Finally, we present a test case similar to the heterogeneous wave propagation problem in Chapter 5, in order to compare the performances of the explicit HHO-leapfrog scheme with standard finite elements.

In this section, we consider a nonhomogeneous Dirichlet boundary prescription  $\mathbf{u}_D$  on  $\partial\Omega_D$ . We briefly expose the impact of this condition on the algebraic formulation and show that there is no change in the splitting procedure. In our implementation, the nonhomogeneous Dirichlet boundary condition is strongly enforced and the corresponding face unknowns are then eliminated. This allows us to keep the same number of degrees of freedom in the global matrix as for homogeneous boundary conditions. We define

$$\mathbf{u}_{\mathcal{F},D}^k := \prod_{F \notin \Omega_D} \{\mathbf{0}\} \times \prod_{F \subset \Omega_D} \mathbb{P}_{d-1}^k(F; \mathbb{R}^d). \quad (6.59)$$

Let  $N_{\mathcal{F},D} := \dim(\mathbf{u}_{\mathcal{F},D}^k)$ , and let  $\mathcal{A}_{\bullet\mathcal{F},D} \in \mathbb{R}^{(N_{\mathcal{T}}+N_{\mathcal{F}}) \times N_{\mathcal{F},D}}$  be the block of the stiffness matrix associated with  $a_h((\mathbf{0}, \mathbf{v}_{\mathcal{F}}), \hat{\mathbf{w}}_h)$  for all  $\hat{\mathbf{w}}_h \in \hat{\mathbf{U}}_{h,0}^{l,k}, \mathbf{v}_{\mathcal{F}} \in \mathbf{U}_{\mathcal{F},D}^k$ . Let  $\mathbf{U}_D^n \in \mathbb{R}_{\mathcal{F},D}^N$  be the vector containing the degrees of freedom of  $\Pi_F^k(\mathbf{u}_D(t^n))$  on all faces  $F \subset \partial\Omega_D$ . We obtain the following fully discrete algebraic formulation:

$$\begin{bmatrix} \mathcal{M} & 0 \\ 0 & 0 \end{bmatrix} \begin{pmatrix} \partial_t^2 \mathbf{U}(t) \\ \cdot \end{pmatrix} + \begin{bmatrix} \mathcal{A}_{\mathcal{T}\mathcal{T}} & \mathcal{A}_{\mathcal{T}\mathcal{F}} \\ \mathcal{A}_{\mathcal{F}\mathcal{T}} & \mathcal{A}_{\mathcal{F}\mathcal{F}} \end{bmatrix} \begin{pmatrix} \mathbf{U}_{\mathcal{T}}(t) \\ \mathbf{U}_{\mathcal{F}}(t) \end{pmatrix} = \begin{bmatrix} \mathbf{F}_{\mathcal{T}}(t) - \mathcal{A}_{\mathcal{T}\mathcal{F},D} \mathbf{U}_D^n \\ -\mathcal{A}_{\mathcal{F}\mathcal{F},D} \mathbf{U}_D^n \end{bmatrix}. \quad (6.60)$$

Since the only change with respect to (6.41) appears on the right-hand side of (6.60), the convergence of the splitting procedure is unaffected. In the mixed-order setting, the algebraic formulation is as follows: Setting  $\mathbf{U}_{\mathcal{F}}^{n,0} := \mathbf{U}_{\mathcal{F}}^{n-1}$ , we seek  $\mathbf{U}_{\mathcal{F}}^{n,m+1} \in \mathbb{R}^{N_{\mathcal{F}}}$  such that, for  $m \geq 0$ ,

$$\mathcal{S}_{\mathcal{F}\mathcal{F}} \mathbf{U}_{\mathcal{F}}^{n,m+1} = -\mathcal{B}_{\mathcal{F}\mathcal{F}} \mathbf{U}_{\mathcal{F}}^{n,m} - \mathcal{A}_{\mathcal{F}\mathcal{T}} \mathbf{U}_{\mathcal{T}}^n - \mathcal{A}_{\mathcal{F}\mathcal{F},D} \mathbf{U}_D^n. \quad (6.61)$$

#### 6.3.1 Convergence order

**Convergence in space.** We consider the time-dependent equation (6.5) and focus on the convergence in space of the HHO space semi-discretization by considering the exact solution

$$\mathbf{u}(x, y, t) := t^2 \phi(x, y), \quad \phi(x, y) := [-\sin(\pi x) \cos(\pi y), \cos(\pi x) \sin(\pi y)]^\top, \quad (6.62)$$

over the domain  $\Omega := (0, 1)^2$ . This solution is obtained by imposing the Dirichlet boundary conditions  $\mathbf{u}_D := \mathbf{u}|_{\partial\Omega}$ , the initial conditions  $\mathbf{u}_0 := \mathbf{0}$ ,  $\mathbf{v}_0 := \mathbf{0}$ , and the source term  $\mathbf{f}(x, y, t) := 2(\mu t^2 \pi^2 + \rho) \phi(x, y)$ . There is no time discretization error since  $\mathbf{u}$  is quadratic in time. We consider the reference material and the final time  $\mathfrak{T} := 0.1$ . Figure 6.8 displays the  $L^2$ - and energy-errors on a sequence of refined triangular meshes,  $k \in \{1, 2, 3\}$ , mixed-order setting. Optimal convergence orders in space are obtained. We considered here only the mixed-order setting for the HHO method for brevity.

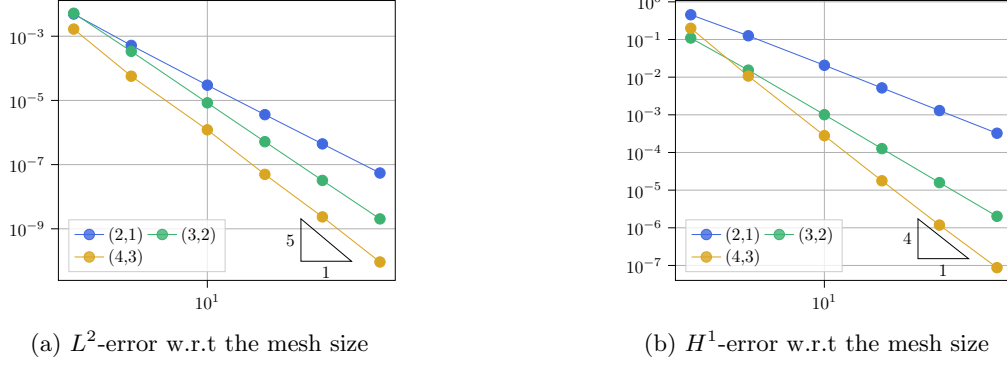


Figure 6.8: Linear elasticity (reference material), analytical test case with quadratic time-dependence - Convergence of the  $L^2$ - and  $H^1$ -errors at the final time  $\mathfrak{T} = 0.1$  as a function of the mesh size,  $k \in \{1, 2, 3\}$ , mixed-order setting.

**Convergence in space and in time.** We now study the convergence in space and time by considering the exact solution

$$\mathbf{u}(x, y, t) := \frac{1}{\sqrt{2\pi}} \sin \sqrt{2\pi t} \phi(x, y),$$

over the domain  $\Omega := (0, 1)^2$ . This solution is obtained by imposing the Dirichlet boundary conditions  $\mathbf{u}_D := \mathbf{u}|_{\partial\Omega}$ , the initial conditions  $\mathbf{u}_0 := \mathbf{0}$ ,  $\mathbf{v}_0 := \phi$ , and the source term  $\mathbf{f}(x, y, t) := \mathbf{0}$ . There are both space and time discretization errors. We consider the reference material and the final time  $\mathfrak{T} := 0.1$ . Firstly, we take a time step  $\Delta t$  proportional to the mesh size  $h$ . The errors are reported in Figure 6.9, which displays the  $L^2$ - and energy-errors on a sequence of refined triangular meshes, with polynomial orders  $k \in \{1, 2, 3\}$  in the mixed-order setting. The convergence rate is of second-order since  $h$  and  $\Delta t$  are proportional and since the leapfrog has second-order accuracy. In order to obtain better convergence rates in space, we must take a time step sufficiently small so that the time discretization error is of the same order as, or smaller than, the space discretization error. This is the case in the second numerical experiment, with errors reported in Figure 6.10 showing optimal convergence rates in space.

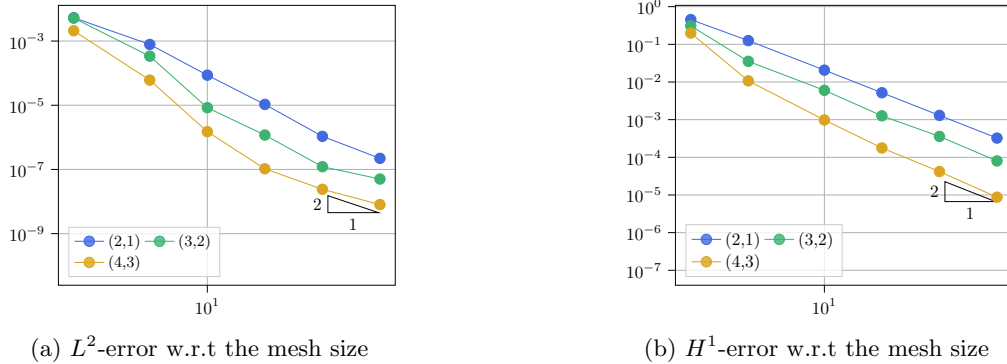


Figure 6.9: Linear elasticity (reference material), analytical test case with nonpolynomial time-dependence, time step proportional to the mesh size - Convergence of the  $L^2$ - and  $H^1$ -errors at the final time  $\mathfrak{T} = 0.1$  as a function of the mesh size,  $k \in \{1, 2, 3\}$ , mixed-order setting.

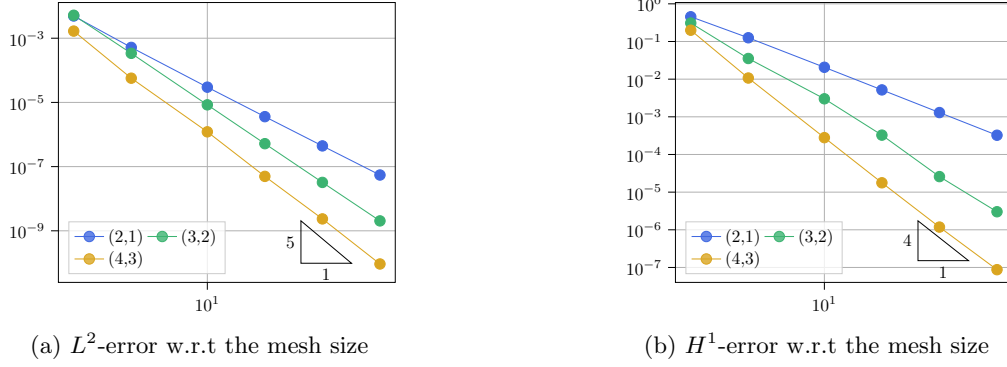


Figure 6.10: Linear elasticity (reference material), analytical test case with nonpolynomial time-dependence, small time step - Convergence of the  $L^2$ - and  $H^1$ -errors at the final time  $\mathfrak{T} = 0.1$  as a function of the mesh size,  $k \in \{1, 2, 3\}$ , mixed-order setting.

### 6.3.2 Choice of the splitting parameter

We want to select the best value of  $\gamma$  for the splitting procedure in terms of error and execution time. We proceed as in Section 4.2.2. We approximate the exact solution (6.62) over a series of refined meshes composed of square cells. Figure 6.11 displays the mean number of splitting iterations per time step (first row), and the total number of splitting iterations during the simulation (second row), for polynomial orders  $k \in \{1, 2, 3\}$  and  $\gamma \in \{1.5\gamma^*, 3\gamma^*, 5\gamma^*, 10\gamma^*\}$ . We recall that the value of  $\gamma^*$  for each polynomial order is given in Table 6.3. For each value of  $\gamma$  and for all the meshes, we set  $\Delta t := 0.8\Delta t^{\text{opt}}(\gamma)$ , with  $\Delta t^{\text{opt}}(\gamma)$  defined in (6.58).

As for the acoustic wave equation, in the mixed-order setting, we observe that the mean number of iterations per time step decreases as  $\gamma$  increases. Since the CFL condition imposes that the time step also decreases, there is a tradeoff between having less iterations per time step and performing more time steps. The value of  $\gamma$  yielding the less total splitting iterations is  $\gamma = 10\gamma^*$  for the polynomial order  $k = 1$  and  $\gamma = 3\gamma^*$  for the polynomial orders  $k \in \{2, 3\}$ . In the equal-order setting, the behavior also resembles that for the acoustic wave equation: the larger  $\gamma$ , the more splitting iterations needed for the splitting to converge. Interestingly, for polynomial orders  $k \in \{2, 3\}$ , the mean number of splitting iterations reaches a maximum on medium-sized meshes, and decreases as the mesh is further refined. This means that the impact of the value of  $\gamma$  on the computational efficiency decreases for finer meshes. But, since the CFL condition tightens with the value of  $\gamma$ , it is still more efficient to take  $\gamma$  as small as possible, for instance  $\gamma = 1.5\gamma^*$ .

Since the choice  $\gamma = 10\gamma^*$  in the mixed-order setting lies beyond the interval  $[3\gamma^*, 5\gamma^*]$  identified for the acoustic wave equation, we verify that such a large value of  $\gamma$  does not affect the error. Figure 6.12 displays, for the same test case and the same meshes as in Figure 6.11, the energy error as a function of the execution time for  $k = 1$  and  $\gamma \in \{1.5\gamma^*, 3\gamma^*, 5\gamma^*, 10\gamma^*\}$ . The curves for  $\gamma \in \{3\gamma^*, 5\gamma^*, 10\gamma^*\}$  almost overlap, and the curve for  $\gamma = 1.5\gamma^*$  is slightly above the three other curves for the finer meshes. The larger values of  $\gamma$  thus yield similar errors to the lower values. Hence, the value of  $\gamma$  can be taken in the interval  $[3\gamma^*, 10\gamma^*]$  with little impact on the precision of the discrete scheme. Since  $\gamma = 10\gamma^*$  yields the smallest total number of iterations for the polynomial order  $k = 1$ , this should be the choice. For the polynomial orders  $k \in \{2, 3\}$ , the choice should be  $\gamma = 3\gamma^*$ , which is also the optimal value in the acoustic case.

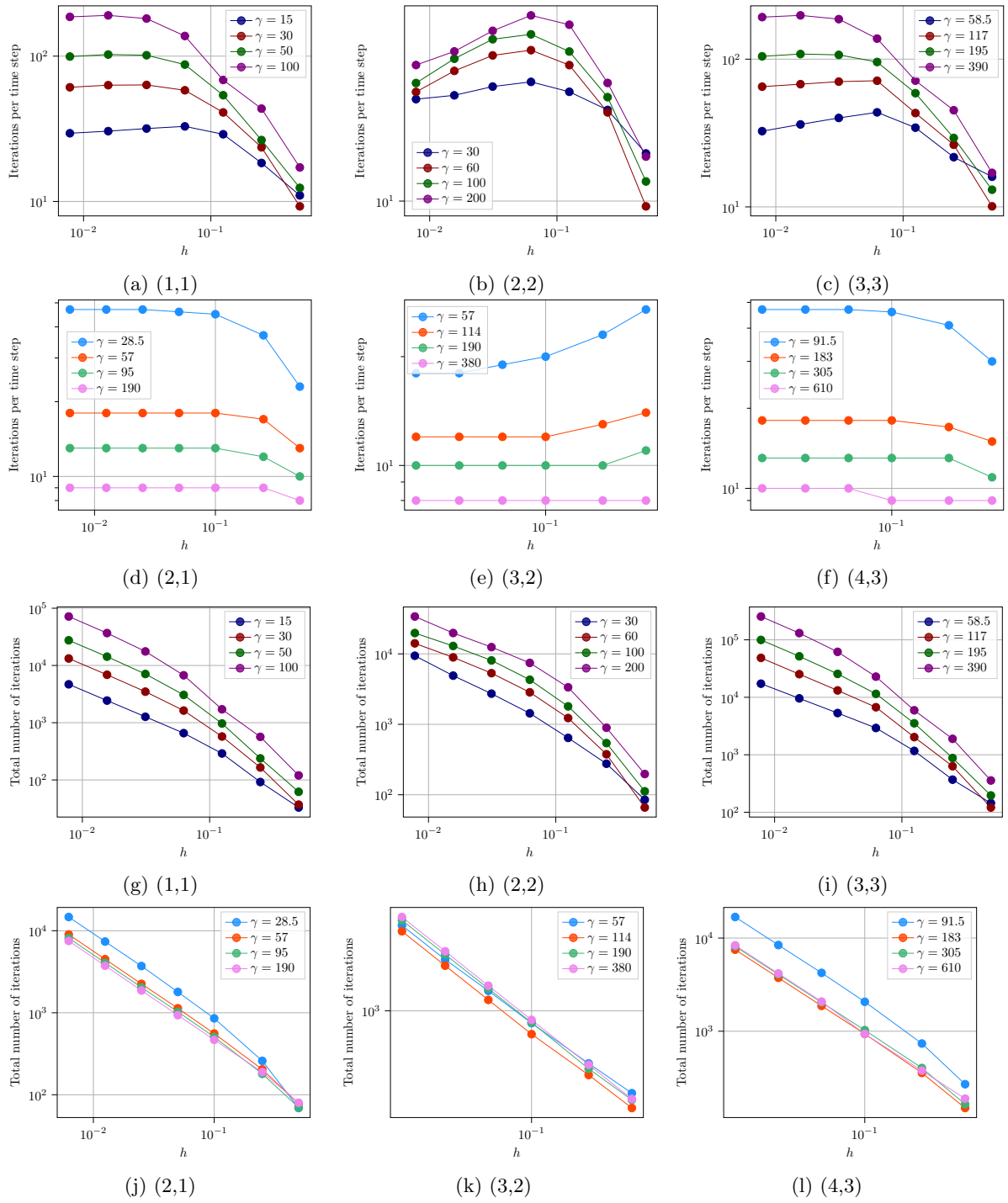


Figure 6.11: Linear elasticity (reference material), analytical test case with quadratic time-dependence - Mean number of splitting iterations in log scale per time step (top two rows) and total number of splitting iterations in log scale (bottom two rows), as a function of the mesh size,  $k \in \{1, 2, 3\}$ , equal- and mixed-order settings.

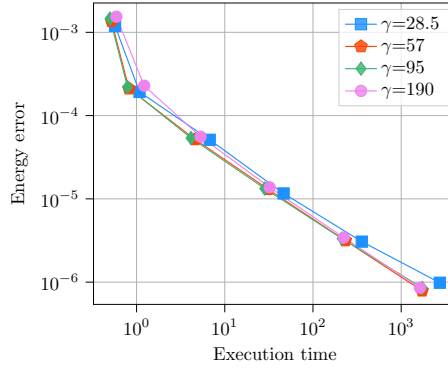


Figure 6.12: Linear elasticity (reference material), analytical test case with quadratic time-dependence - Energy error at  $\mathfrak{T} = 0.1$  as a function of the execution time,  $\gamma \in \{1.5\gamma^*, 3\gamma^*, 5\gamma^*, 10\gamma^*\}$ ,  $k = 1$ , mixed-order setting.

### 6.3.3 Acceleration of the splitting procedure

Results of Figure 6.11 show that, especially in the equal-order setting, the number of splitting iterations at each time step can be large, and reach more than 100. As for the nonlinear wave equation in Section 5.3.2, it is possible to use the  $\Delta^2$  Aitken acceleration procedure (see [8] for the introduction to the method and [134] for a possible implementation) to reduce the number of splitting iteration at each time step. We consider the same test case as in the previous section. Figure 6.13 displays the mean number of splitting iterations at each time step for  $\gamma \in \{1.5\gamma^*, 3\gamma^*, 5\gamma^*, 10\gamma^*\}$ , polynomial orders  $k \in \{1, 2, 3\}$  and mixed- and equal-order settings. The same time step as for Figure 6.11 is used, so that the total number of splitting iterations is reduced by the same factor as the number of iterations at each time step. Since the number of splitting iterations is much smaller than in Figure 6.11, a linear scale is used for the  $y$ -axis. Table 6.7 displays the gain between the plain splitting procedure and the accelerated splitting procedure with Aitken in terms of number of splitting iterations, on the finest mesh with mesh size  $h = 0.0078$ , for  $\gamma \in \{1.5\gamma^*, 3\gamma^*, 5\gamma^*, 10\gamma^*\}$ . These results show that the Aitken procedure accelerates the splitting procedure for all polynomial orders and all values of  $\gamma$ . In the equal-order setting, the acceleration factor is between 2 and 6 and is larger for larger values of  $\gamma$ . It is, however, not sufficient to make large values of  $\gamma$  faster than  $\gamma = 1.5\gamma^*$  for polynomial orders  $k \in \{1, 3\}$ , and  $\gamma = 1.5\gamma^*$  remains the choice for these polynomial orders. For  $k = 2$ , on the contrary,  $\gamma = 10\gamma^*$  with Aitken acceleration requires less splitting iterations than  $\gamma = 1.5\gamma^*$ . It is hence interesting to look at the total number of splitting iterations for this polynomial order. In the same manner, in the mixed-order setting, the gain is larger for  $\gamma = 1.5\gamma^*$ , up to a factor 8, and could make the choice  $\gamma = 1.5\gamma^*$  more attractive.

Polynomial order $(l, k)$	(1,1)	(2,2)	(3,3)	(2,1)	(3,2)	(4,3)
$\gamma = 1.5\gamma^*$	2.6	2.6	2.9	8.0	3.0	5.4
$\gamma = 3\gamma^*$	4.1	3.1	4.1	3.3	2.0	3.0
$\gamma = 5\gamma^*$	4.9	3.7	5.5	2.8	1.7	2.2
$\gamma = 10\gamma^*$	5.6	5.6	5.8	3.0	1.5	1.7

Table 6.7: Linear elasticity (reference material), analytical test case with quadratic time-dependence - Gain in number of splitting iterations for the plain splitting procedure versus the Aitken accelerated splitting procedure for  $\gamma \in \{1.5\gamma^*, 3\gamma^*, 5\gamma^*, 10\gamma^*\}$ , mesh size  $h = 0.0078$ , polynomial orders  $k \in \{1, 2, 3\}$ , equal- and mixed-order settings.

Figure 6.14 displays the total number of splitting iterations in the equal-order setting with  $k = 2$  and in the mixed-order setting for  $k \in \{1, 2, 3\}$ . In the equal-order setting with  $k = 2$ , the less number of splitting iterations is still obtained with  $\gamma = 1.5\gamma^*$ . In the mixed-order setting, for  $k \in \{1, 2\}$ ,  $\gamma = 1.5\gamma^*$

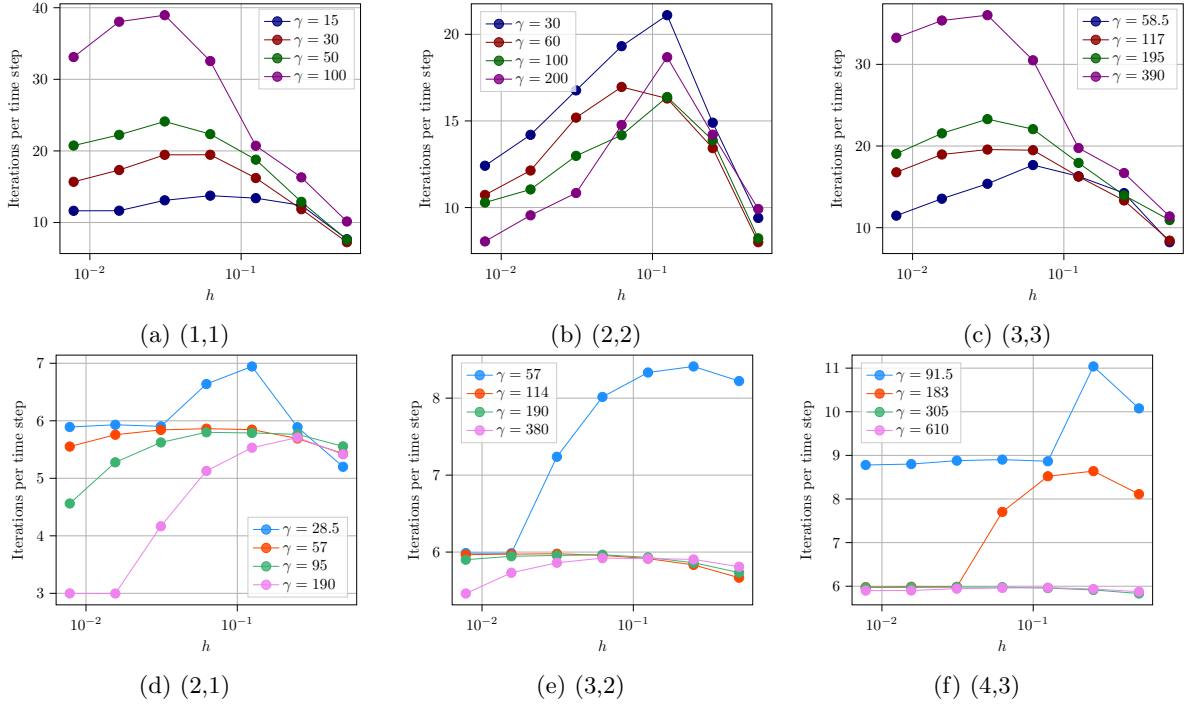


Figure 6.13: Linear elasticity (reference material), analytical test case with quadratic time-dependence - Mean number of splitting iterations with Aitken acceleration per time step, as a function of the mesh size,  $k \in \{1, 2, 3\}$ , equal- and mixed-order settings.

is also the choice yielding the smallest number of splitting iterations, whereas for  $k = 3$ ,  $\gamma = 1.5\gamma^*$  and  $\gamma = 3\gamma^*$  yield an equivalent number of splitting iterations.

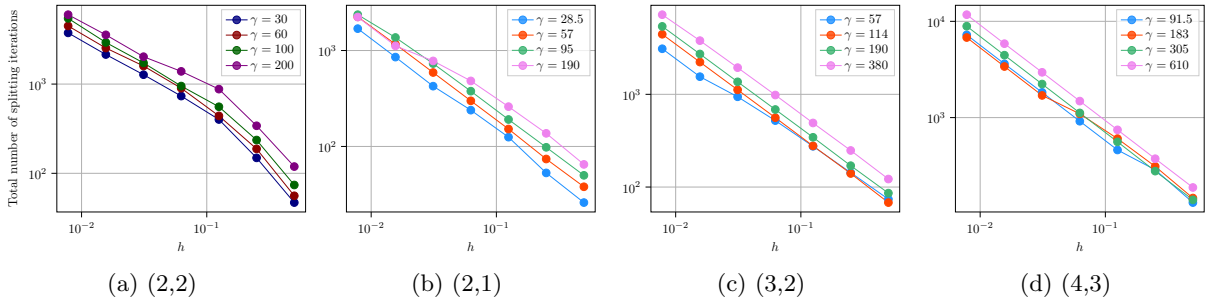


Figure 6.14: Linear elasticity (reference material), analytical test case with quadratic time-dependence - Total number of splitting iterations with Aitken acceleration per time step (in log scale), as a function of the mesh size,  $k \in \{1, 2, 3\}$  in the mixed-order setting and  $k = 2$  in the equal-order setting.

These results illustrate the efficiency of the Aitken acceleration to reduce the number of splitting iterations. For all polynomial orders, on the finest mesh, the mean number of splitting iterations for  $\gamma = 1.5\gamma^*$  is comprised between 3 and 13 with the Aitken acceleration. It also changes the choice for the optimal value of  $\gamma$  in the mixed-order setting. For all polynomial orders, equal- and mixed-order settings, the best value of  $\gamma$  with the Aitken acceleration is  $\gamma = 1.5\gamma^*$ . This is the choice made in the rest of this section.

### 6.3.4 Comparison with finite elements

We now focus on a wave propagation problem in a heterogeneous medium. This test case is analogous to the one presented in Section 4.3. The goal is to evaluate the performances of the splitting procedure compared to the standard finite element method. It is expected that the HHO method gives better performances when considering the propagation of elastic waves in a body composed of contrasted materials.

**Presentation of the test case.** We consider the following test case (see Figure 6.15):  $\Omega := (-1, 1)^2$  with two subdomains  $\Omega_1 := (-1, 1) \times (-1, 0)$  and  $\Omega_2 := (-1, 1) \times (0, 1)$  with respective material properties  $\lambda_1 := \mu_1 := 1$ ,  $\rho_1 := 1$  and  $\lambda_2 := \mu_2 := 4$ ,  $\rho_2 := 1$ , which give respective speeds of sound  $c_{P,1} = \sqrt{3}$ ,  $c_{S,1} = 1$  and  $c_{P,2} = 2\sqrt{3}$ ,  $c_{S,2} = 2$ . These values of the Lamé coefficients give  $E_1 = 2.5$ ,  $\nu_1 = 0.25$  and  $E_2 = 10$ ,  $\nu_2 = 0.25$ . A source is placed in the upper subdomain  $\Omega_2$  at the point of coordinates  $S := (0, 0.5)$ . The wave first propagates in  $\Omega_2$ , then is partially reflected at the interface between  $\Omega_1$  and  $\Omega_2$  and partially transmitted to the lower subdomain  $\Omega_1$ . We place 9 sensors in  $\Omega$  as illustrated in Figure 6.15. The simulations are run until  $\mathfrak{T} := 0.6$ . Reflecting boundary conditions are enforced on

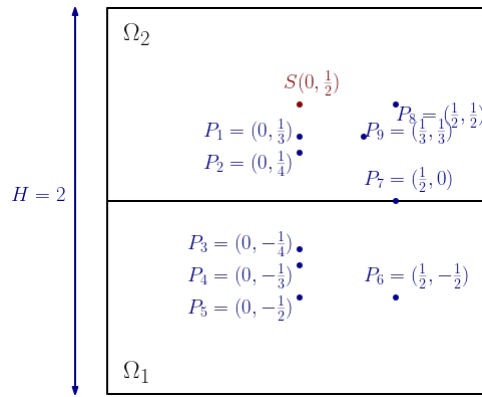


Figure 6.15: Ricker wavelet propagation in heterogeneous medium - Schematics of the physical domain: two-layer medium with a Dirac source at  $S$  and 9 sensors placed at points  $P_1$  to  $P_9$ .

$\partial\Omega$ . The initial source is a Ricker wavelet in time and a Dirac measure in space centered at the point  $S$ . The shape of a Ricker wavelet is illustrated in Chapter 4 in Figure 4.13b. We take a peak frequency  $f_p := 20\sqrt{3}$ . We use the open-source software `Gar6More2D` [89] to compute the semi-analytical solution at the 9 sensor points. The comparison is valid only until the waves reflecting on the boundary of  $\Omega$  reach the sensors. These maximal validity times can be computed by geometric arguments (see Table 6.8) using the speed of P-waves, since they are faster than S-waves. When comparing wave simulations with the HHO method or the finite element method, we need the initial data  $\mathbf{u}_0$  and  $\mathbf{v}_0$ . In order to simulate the same setting as in `Gar6more2D`, the initial condition is also a Ricker wavelet, but in space, expressed in terms of the variable  $r = \sqrt{(x - x_S)^2 + (y - y_S)^2}$ , the distance to the source. One can show that this initial condition is equivalent to the initial source used in `Gar6More2D`, with only a time shift. This time shift is a given data of the semi-analytical simulation with `Gar6More2D`. For a given time peak frequency  $f_p$  and a speed of sound for pressure waves  $c_p$  in the domain of the source (here  $c_{P,2}$ ), the corresponding wave number  $k_p$  is given by  $k_p = c_p^{-1} f_p = 10$ . The initial condition reads

$$\mathbf{u}_0(x, y) = e^{-\pi^2 r^2 k_p^2} [(x - x_S), (y - y_S)]^T, \quad \mathbf{v}_0(x, y) := \mathbf{0}. \quad (6.63)$$

We notice that the initial condition is on the displacement and not on the velocity as it was the case with the acoustic waves; this is due to the method used in `Gar6more2D` for elastic waves. Figure 6.16 displays the semi-analytical solution  $\mathbf{u}$  at the sensor  $P_9$  in the time interval  $[0, 0.6]$ . The semi-analytical solution is in agreement with the times reported in Table 6.8 for the P-wave to reach the sensor  $P_9$ . The third wave reaching the sensor at  $t = 0.32$  is the S-wave reflected at the interface.



Sensor	$P_1$	$P_2$	$P_7$	$P_8$	$P_9$
from $S$ to $P$	0.048	0.07	0.20	0.14	0.11
reflection of P-waves on the interface to $P$	0.24	0.21	0.20	0.32	0.26
reflection of P-waves on the top boundary to $P$	0.33	0.36	0.46	0.32	0.52

Table 6.8: Ricker wavelet propagation in heterogeneous medium - Times for different wave paths to reach the sensors  $P_1, P_2, P_7, P_8$  and  $P_9$ .

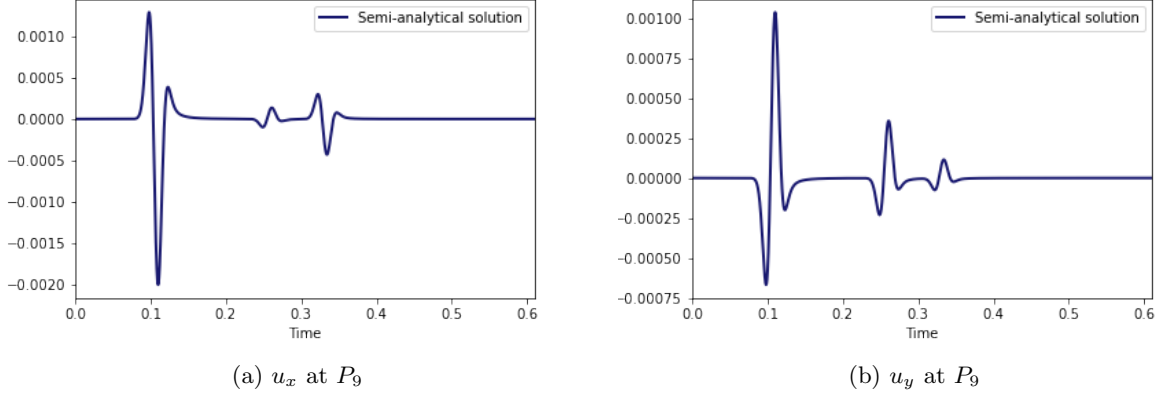


Figure 6.16: Ricker wavelet propagation in heterogeneous medium - Semi-analytical solution:  $x$ - and  $y$ -components of the displacement at sensor  $P_9$ .

Figure 6.17 displays the solution on the computational domain  $\Omega$  computed using  $P2$  finite elements. Reflections at the boundaries and at the interface are composed of  $P$ - and  $S$ -waves. The  $S$ -waves are particularly visible on Figure 6.17c (reflection at the interface) and on Figure 6.17g (reflection at the top boundary).

**Error on a fixed mesh.** We start by running the HHO and finite element simulations on the same unstructured triangular mesh with  $h = 0.02$ . Since the Poisson ratio is  $\nu = 0.25$  in both domains, the splitting parameter  $\gamma^*$  is constant over  $\Omega$  and corresponds to the value for the reference material. The HHO simulations are run with the splitting procedure, using  $\gamma = 1.5\gamma^*$  for all polynomial orders. For both HHO and finite element simulations, the time step is taken so that  $\Delta t = 0.8\Delta t^{\text{opt}}(\gamma)$ . The obtained solution is shown in Figure 6.18, where we consider  $u_x$  and  $u_y$  at  $P_9$ ,  $u_x$  at  $P_6$  and  $u_y$  at  $P_4$ . The solution for P1 finite elements is very inaccurate and is only shown on Figure 6.18a for readability. The solution for P2 finite elements shows dispersion at  $P_6$ , which is located in the lower subdomain  $\Omega_2$ . The HHO solution with polynomial order  $k = 1$  shows oscillations after the first wave front at  $P_4$  and  $P_6$ , whereas at  $P_9$  and for higher orders, the HHO solution matches well the semi-analytical solution. This illustrates (again) the fact that finite elements are more dispersive than HHO in heterogeneous media and that higher orders for space semi-discretization handle dispersion better than lower orders. We compute the maximum over the 9 sensors of the  $\ell^2$ -in-time energy-error using the simulation results and the semi-analytical solution and report this error in Table 6.9. As shown in Figure 6.18, P1 finite elements yield the largest error, and HHO with polynomial orders  $k \in \{2, 3\}$  yields a much smaller error than the lower-order settings.

**Error with a fixed number of degrees of freedom.** Since the number of degrees of freedom for finite elements and HHO are not the same, we now compare errors with (approximately) the same number of degrees of freedom. We consider two cases: a coarse case with 36K degrees of freedom and a fine case with 280K degrees of freedom. The corresponding mesh sizes are reported in Table 6.10. The maximum

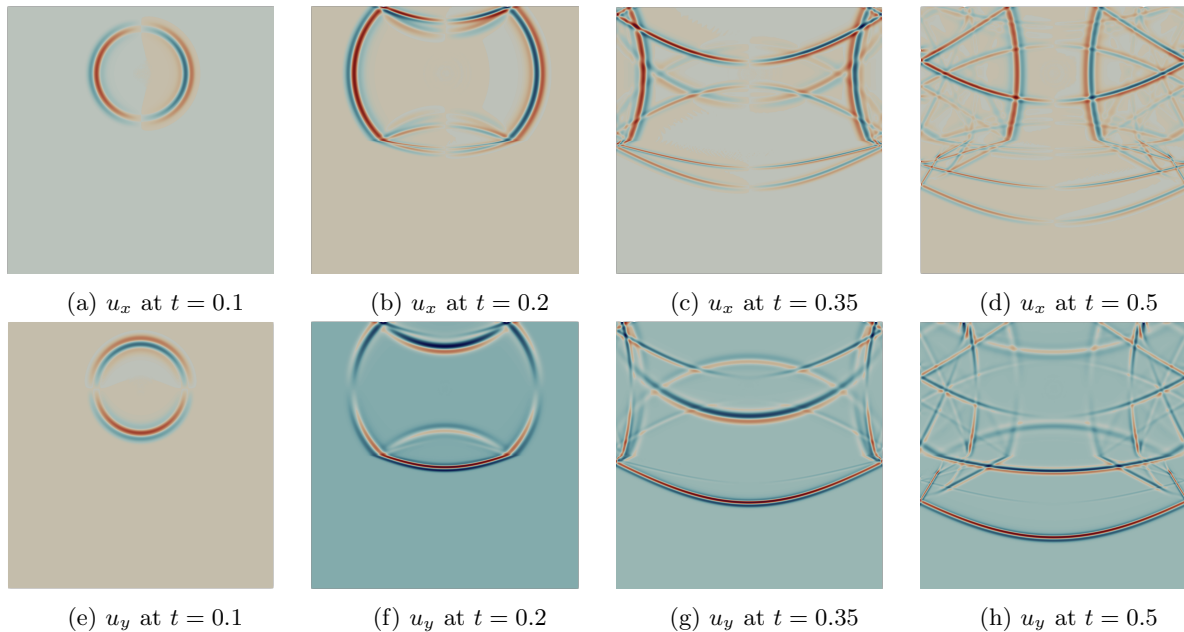


Figure 6.17: Ricker wavelet propagation in heterogeneous medium - Solution at different times, computed using  $P2$  finite elements with MANTA,  $x$ -component (top) and  $y$ -component (bottom) of the displacement.

Method	HHO(2,1)	HHO(3,2)	HHO(4,3)	P1 FEM	P2 FEM
Energy error	$1.7e^{-3}$	$2.2e^{-4}$	$9.6e^{-5}$	0.88	$3.9e^{-2}$

Table 6.9: Ricker wavelet propagation in heterogeneous medium -  $\ell^2$ -in-time energy-error on the displacement (maximal value over the 9 sensors) with mesh size  $h = 0.02$ , HHO with polynomial orders  $k \in \{1, 2, 3\}$  and mixed-order setting, P1 and P2 finite elements.

of the  $\ell^2$ -error in time is given in Table 6.11. In the coarse case, the HHO error is smaller than the finite element error and decreases with the polynomial order. In the fine case, the difference between HHO and finite elements is less pronounced. This is due to the fact that the time step is proportional to  $h$ . The time semi-discretization error converges as  $\Delta t^2$  and the space semi-discretization error converges as  $\mathcal{O}(h^{k+2})$  for HHO and as  $\mathcal{O}(h^{k+1})$  for finite elements. The convergence rate of the time semi-discretization error is thus slower than the convergence rate of the space semi-discretization for P2 finite elements and for HHO with polynomial orders  $k \in \{1, 2, 3\}$ . Nonetheless, employing the HHO method is more robust to dispersion in this heterogeneous setting and yields smaller errors than with finite elements. This confirms the attractive performances of the HHO method on this kind of problems.

Number of dofs	HHO(2,1)	HHO(3,2)	HHO(4,3)	P1 FEM	P2 FEM
36K	0.03	0.035	0.04	0.012	0.023
280K	0.01	0.012	0.015	0.004	0.008

Table 6.10: Ricker wavelet propagation in heterogeneous medium - Mesh size for HHO with polynomial orders  $k \in \{1, 2, 3\}$ , mixed-order setting, and for P1 and P2 finite elements, all corresponding to either 36K or 280K dofs.

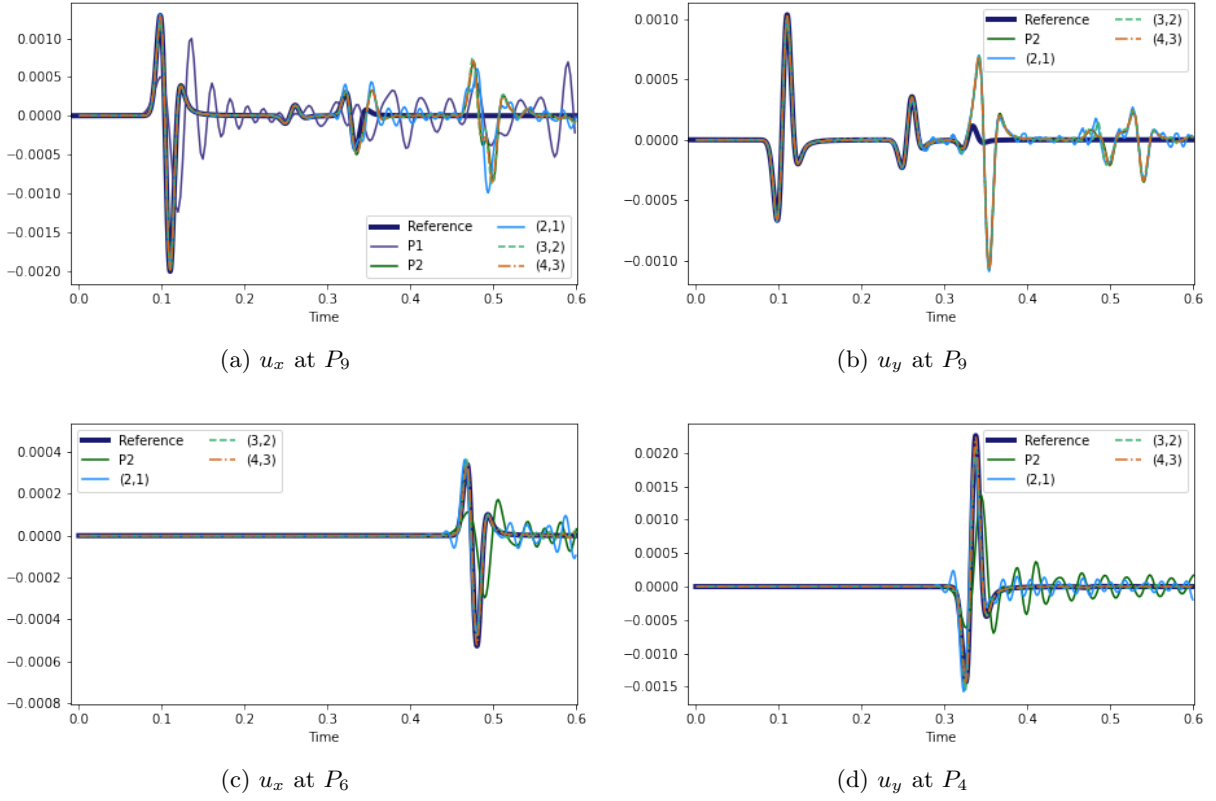


Figure 6.18: Ricker wavelet propagation in heterogeneous medium - Simulations on the same mesh with  $h = 0.02$ , HHO with polynomial orders  $k \in \{1, 2, 3\}$  in the mixed-order setting, P1 (only top left) and P2 finite elements, and semi-analytical solution (`Gar6more2D`).

$\ell^2$ -error	HHO(2,1)	HHO(3,2)	HHO(4,3)	P1 FEM	P2 FEM
36K	0.22	$8.9e^{-3}$	$7.2e^{-4}$	0.53	0.10
280K	$1.7e^{-4}$	$1.2e^{-4}$	$9.5e^{-5}$	$7.4e^{-3}$	$3.44e^{-3}$

Table 6.11: Ricker wavelet propagation in heterogeneous medium -  $\ell^2$ -in-time energy-error for HHO with polynomial orders  $k \in \{1, 2, 3\}$  in the mixed-order setting, and P1 and P2 finite elements, all corresponding to either 36K or 280K dofs.

## 6.4 Full discretization of structural dynamics with finite-strain plasticity

We directly consider finite-strain plasticity instead of small-strain plasticity since it is more relevant to industrial applications. The HHO discretization of static finite-strain plastic deformations has been developed in [2]. Here, we introduce the space semi-discretization of the structural dynamics equation using the same formalism. Interestingly, due to the incremental nature of plasticity, it is more natural to start with the semi-discretization in time. We then expose the HHO space semi-discretization and the associated splitting procedure, which is analogous to that devised for the nonlinear acoustic equation. Indeed, in the case of finite deformations, the HHO stabilization acts componentwise as the scalar stabilization used for acoustic waves. The convergence of the splitting procedure in the equal-order setting is hence

possible.

### 6.4.1 Plasticity model

Contrary to the linear elasticity model, the elastoplastic model considers that deformations are irreversible. We place ourselves within the framework of generalized standard materials [108, 126], which defines the thermodynamically admissible plasticity models. The plasticity model is assumed to be strain-hardening and rate-dependent, i.e. the speed of deformations has an impact on the solution.

**Kinematics.** We still consider a body occupying the reference domain  $\Omega$ . We define the deformation gradient  $\mathbf{F}(\mathbf{u}) = \mathbf{I}_d + \nabla \mathbf{u} \in \mathbb{R}_+^{d \times d}$ , where  $\mathbb{R}_+^{d \times d}$  is the set of  $\mathbb{R}^{d \times d}$ -matrices with positive determinant. We adopt the logarithmic strain framework [131] developed for anisotropic finite elastoplasticity, leading to the following (logarithmic) strain tensor:

$$\mathbf{E}(\mathbf{u}) := \frac{1}{2} \ln(\mathbf{F}(\mathbf{u})^\top \mathbf{F}(\mathbf{u})) \in \mathbb{R}_{\text{sym}}^{d \times d}, \quad (6.64)$$

where the logarithm is computed by means of an eigenvalue decomposition. The plastic deformations are measured using the (logarithmic) plastic strain tensor  $\mathbf{E}_p(\mathbf{u}) \in \mathbb{R}_{\text{sym}}^{d \times d}$ , and we assume that the (logarithmic) elastic strain tensor  $\mathbf{E}_e(\mathbf{u}) \in \mathbb{R}_{\text{sym}}^{d \times d}$  can be defined additively by the following decomposition:

$$\mathbf{E}_e(\mathbf{u}) := \mathbf{E}(\mathbf{u}) - \mathbf{E}_p(\mathbf{u}) \in \mathbb{R}_{\text{sym}}^{d \times d}. \quad (6.65)$$

**Plasticity modelling.** The local material state is described by the (logarithmic) plastic strain tensor and a finite collection of internal variables,  $\alpha := (\alpha_1, \dots, \alpha_m) \in \mathbb{R}^m$ , which typically contains at least the equivalent plastic strain  $p \geq 0$ . We define  $\mathcal{X}$ , the space of the generalized internal variables, as

$$\mathcal{X} := \{ \chi = (\mathbf{e}_p, \alpha) \in \mathbb{R}_{\text{sym}}^{d \times d} \times \mathbb{R}^m \text{ s.t. } \text{tr}(\mathbf{e}_p) = 0 \}. \quad (6.66)$$

The Helmholtz free energy  $\Psi : \mathbb{R}_{\text{sym}}^{d \times d} \times \mathbb{R}^m \rightarrow \mathbb{R}$  acts on a generic pair  $(\mathbf{e}_e, \alpha)$  representing the (logarithmic) elastic strain tensor and the internal variables. The Helmholtz free energy is assumed to satisfy the following hypothesis.

**Hypothesis 6.4.1** (Helmholtz free energy).  $\Psi$  can be decomposed additively into an elastic and a plastic part as follows:

$$\Psi(\mathbf{e}_e, \alpha) := \Psi^e(\mathbf{e}_e) + \Psi^p(\alpha), \quad \Psi^e(\mathbf{e}) := \frac{1}{2} \mathbf{e}_e : \mathbb{A} : \mathbf{e}_e, \quad (6.67)$$

where the elastic modulus  $\mathbb{A}$  is defined by (6.2) and  $\Psi^p : \mathbb{R}^m \rightarrow \mathbb{R}$  is strongly convex. We assume here that the Lamé coefficients  $\mu$  and  $\lambda$  are constant over  $\Omega$ .

Following the second principle of thermodynamics, the logarithmic stress tensor  $\mathbf{T} \in \mathbb{R}_{\text{sym}}^{d \times d}$  and the thermodynamic forces  $q \in \mathbb{R}^m$  are derived from  $\Psi$  as follows:

$$\mathbf{T}(\mathbf{e}_e) := \partial_{\mathbf{e}_e} \Psi^e(\mathbf{e}_e) := \mathbb{A} : \mathbf{e}_e \in \mathbb{R}_{\text{sym}}^{d \times d} \quad \text{and} \quad q(\alpha) = \partial_\alpha \Psi^p(\alpha) \in \mathbb{R}^m. \quad (6.68)$$

(Notice that  $\mathbf{T}(\boldsymbol{\epsilon}(\mathbf{u})) = \mathbb{A} : \boldsymbol{\epsilon}(\mathbf{u})$  coincides with the stress tensor in the case of the infinitesimal deformations considered in Section 6.1.)

The criterion to determine whether the deformations become plastic hinges on the scalar-valued yield function  $\Phi : \mathbb{R}_{\text{sym}}^{d \times d} \times \mathbb{R}^m \rightarrow \mathbb{R}$ , which is a continuous and convex function of the (logarithmic) stress tensor  $\mathbf{T}$  and the thermodynamic forces  $q$ . Letting  $\mathcal{A} := \{ (\mathbf{T}, q) \in \mathbb{R}_{\text{sym}}^{d \times d} \times \mathbb{R}^m \mid \Phi(\mathbf{T}, q) \leq 0 \}$  be the convex set of admissible states, the elastic domain  $\mathcal{A}_e$  is composed of all the pairs  $(\mathbf{T}, q)$  such that  $\Phi(\mathbf{T}, q) < 0$ , and the yield surface  $\partial \mathcal{A}$  of all the pairs  $(\mathbf{T}, q)$  such that  $\Phi(\mathbf{T}, q) = 0$ .

**Hypothesis 6.4.2** (Yield function). The yield function satisfies the following properties:

1.  $\Phi$  is piecewise analytical;
2. The point  $(\mathbf{0}, 0)$  lies in the elastic domain, i.e.,  $\Phi(\mathbf{0}, 0) < 0$ ;
3.  $\Phi$  is differentiable at all points on the yield surface  $\partial\mathcal{A}$ .

**Example.** The hyperelasticity case corresponds to the case  $\Psi^P := 0$ . A simple plasticity model is the linear isotropic hardening model with a von Mises yield criterion. The only internal variable is the equivalent plastic strain  $p \geq 0$ . The plastic part of the free energy is

$$\Psi^P(p) := \sigma_{y,0}p + \frac{1}{2}Hp^2, \quad (6.69)$$

where  $H$  is the isotropic hardening modulus and  $\sigma_{y,0}$  the yield stress. The internal force is then  $q(p) := \sigma_{y,0} + Hp$ . The perfect plasticity model is retrieved by taking  $H := 0$ . Using the  $J_2$ -plasticity model with a von Mises criterion, the yield function writes

$$\Phi(\mathbf{T}, q) := \sqrt{\frac{3}{2}} \|\text{dev}(\mathbf{T})\|_{\ell^2} - q, \quad (6.70)$$

with  $\text{dev}(\mathbf{w}) := \mathbf{w} - \frac{1}{d}\text{tr}(\mathbf{w})\mathbf{I}_d$  the deviatoric operator and  $\|\cdot\|_{\ell^2}$  the Frobenius norm such that  $\|\mathbf{w}\|_{\ell^2} := \sqrt{\mathbf{w} : \mathbf{w}}$ , for all  $\mathbf{w} \in \mathbb{R}^{d \times d}$ .

**Finite elastoplasticity in incremental form.** We are interested in finding the dynamic evolution of the elastoplastic material body in the time interval  $J := [0, \mathfrak{T}]$ ,  $\mathfrak{T} > 0$ . Using the incremental form of the plasticity model (see [138, 130, 50]), we start by the time semi-discretization. We discretize the time interval into  $N$  discrete time intervals such that  $(t^n)_{n \in \{0:N\}}$  are the discrete time nodes with  $t^0 = 0$  and  $t^N := \mathfrak{T}$ . For simplicity, we consider a fixed time step  $\Delta t := \frac{\mathfrak{T}}{N}$ .

At each time step  $t^n$ , the body is subjected to a body force  $\mathbf{f}^n : \Omega \rightarrow \mathbb{R}^d$  and a prescribed displacement  $\mathbf{u}_D^n : \partial\Omega_D \rightarrow \mathbb{R}^d$  on  $\partial\Omega_D \subset \partial\Omega$ . For the sake of simplicity, we assume that  $\mathbf{u}_D^n = \mathbf{0}$  for all  $n \in \{0:N\}$ .

Let  $\mathbf{F}^n$ ,  $\mathbf{E}^n$ ,  $\mathbf{E}_e^n$ ,  $\mathbf{E}_p^n$ ,  $\mathbf{T}^n$  and  $\alpha^n$  be respectively the deformation gradient, the (logarithmic) strain tensor, the (logarithmic) elastic strain tensor, the (logarithmic) plastic strain tensor, the (logarithmic) stress tensor, and the internal variables at the discrete time  $t^n$ . We set  $\chi^n := (\mathbf{E}_p^n, \alpha^n)$ . In order to describe the evolution of the body, we introduce the first Piola–Kirchhoff stress tensor, which is defined at each discrete time  $t^n$  as follows:

$$\mathbf{\Pi}^n := \partial_{\mathbf{F}}\Psi(\mathbf{E}_e^n, \alpha^n) = \mathbf{T}^n(\mathbf{E}_e^n) : \partial_{\mathbf{F}}\mathbf{E}_e^n. \quad (6.71)$$

We denote by  $\mathbf{u}^n$  the unknown at time  $t^n$ . The leapfrog scheme consists in solving, for all  $n \in \{1:N-1\}$ , the following problem:

$$\frac{1}{\Delta t^2} (\rho(\mathbf{u}^{n+1} - 2\mathbf{u}^n + \mathbf{u}^{n-1}), \mathbf{w})_{\Omega} + (\mathbf{\Pi}^n, \nabla \mathbf{w})_{\Omega} = (\mathbf{f}^n, \mathbf{w})_{\Omega}, \quad \forall \mathbf{w} \in \mathbf{H}_0^1(\Omega), \quad (6.72a)$$

$$(\chi^n, \mathbf{\Pi}^n) := \mathcal{P}(\chi^{n-1}, \mathbf{F}^n, \mathbf{F}^{n-1}), \quad \text{a.e. in } \Omega, \quad (6.72b)$$

where  $\mathcal{P}$  is the function computing the generalized internal variables and the first Piola–Kirchhoff stress tensor at the current time  $t^n$ , using the generalized internal variables at the previous time  $t^{n-1}$ , and the deformation gradients at the discrete times  $t^n$  and  $t^{n-1}$ . The procedure consists in solving the following optimization problem:

$$\mathbf{E}_p^n - \mathbf{E}_p^{n-1} = \Lambda(\mathbf{T}^n, q^n) \partial_{\mathbf{T}}\Phi(\mathbf{T}^n, q^n), \quad \alpha^n - \alpha^{n-1} = -\Lambda(\mathbf{T}^n, q^n) \partial_q\Phi(\mathbf{T}^n, q^n), \quad (6.73a)$$

with the Lagrange multiplier  $\Lambda(\mathbf{T}^n, q^n)$  verifying the Karush–Kuhn–Tucker conditions

$$\Lambda(\mathbf{T}^n, q^n) \geq 0, \quad \Phi(\mathbf{T}^n, q^n) \leq 0, \quad \Lambda(\mathbf{T}^n, q^n)\Phi(\mathbf{T}^n, q^n) = 0. \quad (6.73b)$$

Problem (6.72) is solved in two steps:

1. Compute  $\mathbf{\Pi}^n$  by solving the optimization problem (6.73) where  $\chi^{n-1}$ ,  $\mathbf{F}^n$  and  $\mathbf{F}^{n-1}$  are known from the values at prior time steps or from the initial conditions.
2. Compute  $\mathbf{u}^{n+1}$  using (6.72a). This second step is fully explicit.

**Remark 6.4.3** (Initial conditions). Due to the incremental model used for plasticity, the initial conditions impose  $\mathbf{u}^0 = \mathbf{0}$  and  $\chi^0 = 0$ , i.e. the solid is at an equilibrium state at  $t = 0$ . Notice that we can still consider  $\mathbf{v}^0 \neq \mathbf{0}$ . This directly gives  $\mathbf{F}^0 = \mathbf{I}_d$  and  $\mathbf{\Pi}^0 = \mathbf{0}$ . The equivalent of equation (6.72a) at  $t = 0$  is then  $(\rho \mathbf{u}^1, \mathbf{w})_\Omega = \Delta t (\rho \mathbf{v}^0, \mathbf{w})_\Omega + \frac{\Delta t^2}{2} (\mathbf{f}^0, \mathbf{w})_\Omega$ , for all  $\mathbf{w} \in \mathbf{H}_0^1(\Omega)$ .  $\square$

## 6.4.2 Explicit HHO-leapfrog scheme

The HHO unknowns are defined in the same discrete spaces as in Section 6.1.2 for linear elasticity. However, in the finite-strain case, we consider directly a full gradient reconstruction operator. Its definition is similar to the one for the acoustic wave equation, but using vector-valued spaces. In this section, we detail the changes in the definition of the gradient reconstruction operator and the stabilization operator, before presenting the global discrete formulation. This formulation still leads to a static coupling between the cell and the face unknowns, that can be solved, once again, by using the splitting procedure introduced in Chapter 4.

### 6.4.2.1 HHO operators

We define the gradient reconstruction operator  $\mathbf{G}_T^k : \widehat{\mathbf{U}}_T^{l,k} \rightarrow \mathbf{Z}(T; \mathbb{R}^{d \times d})$ , where  $\mathbf{Z}(T; \mathbb{R}^{d \times d})$  is a reconstruction space composed of  $\mathbb{R}^{d \times d}$ -valued polynomials in the cell  $T$ . In the same manner as for the linear elasticity case, we consider the full-tensor (FT) polynomial space  $\mathbf{Z}(T; \mathbb{R}^{d \times d}) := \mathbb{P}_d^k(T; \mathbb{R}^{d \times d})$  and we define the gradient reconstruction operator  $\mathbf{G}_T^k : \widehat{\mathbf{U}}_T^{l,k} \rightarrow \mathbb{P}_d^k(T; \mathbb{R}^{d \times d})$  such that, for all  $\mathbf{q} \in \mathbb{P}_d^k(T; \mathbb{R}^{d \times d})$ ,

$$(\mathbf{G}_T^k(\widehat{\mathbf{v}}_T), \mathbf{q}_T)_T = (\nabla \mathbf{v}_T, \mathbf{q})_T + (\mathbf{v}_{\partial T} - \mathbf{v}_T, \mathbf{q} \cdot \mathbf{n}_T)_{\partial T}. \quad (6.74)$$

The stabilization operator in the mixed-order setting is still defined by (6.19), where we omit the superscript MO. In the equal-order setting, we use the displacement reconstruction operator  $\mathbf{R}_T^{k+1} := \mathbf{R}_T^{\text{FG}, k+1}$  defined by (6.20), and write (we omit the superscript FG)

$$\mathbf{S}_{TF}(\widehat{\mathbf{w}}_T) := \Pi_F^k(\delta_{TF}(\widehat{\mathbf{w}}_T) + (I - \Pi_T^k)\mathbf{R}_T^{k+1}(\widehat{\mathbf{v}}_T)(0, \delta_T(\widehat{\mathbf{w}}_T))|_F), \quad \forall \widehat{\mathbf{w}}_T \in \widehat{\mathbf{U}}_T^{k,k} := \widehat{\mathbf{U}}_T^k, \quad (6.75)$$

which is nothing but the vector-valued version of the equal-order stabilization operator considered in the acoustic case. The local stabilization bilinear form  $s_T$  is still defined by (6.23). For simplicity, we assume that the stabilization is scaled by a parameter  $\gamma\mu$  that is constant over all the mesh cells; a discussion for this choice can be found in [2].

**Remark 6.4.4.** (Gradient reconstruction) The full gradient reconstruction verifies

$$\nabla \mathbf{R}_T^{\text{FG}, k+1} = \Pi_{\nabla \mathbb{P}_d^{k+1}(T; \mathbb{R}^d)} \mathbf{G}_T^k.$$

As for linear elasticity, it is possible to define a full gradient reconstruction operator  $\mathbf{G}_T^{\text{FG}, k} := \nabla \mathbf{R}_T^{\text{FG}, k+1}$ , but, considering the consistency issues raised for linear elasticity, such a reconstruction is not considered for plasticity.  $\square$

### 6.4.2.2 Global discretization

On each cell  $T \in \mathcal{T}_h$ , we define the local deformation gradient operator  $\mathbf{F}_T^k : \widehat{\mathbf{U}}_T^{l,k} \rightarrow \mathbb{P}_d^k(T; \mathbb{R}^{d \times d})$  such that  $\mathbf{F}_T^k := \mathbf{I} + \mathbf{G}_T^k$ , where  $\mathbf{I}$  is the identity operator in  $\mathbb{P}_d^k(T; \mathbb{R}^{d \times d})$ . We define the global gradient

reconstruction and the global deformation gradient operators such that, for all  $T \in \mathcal{T}_h$  and all  $\hat{\mathbf{v}}_h \in \widehat{\mathcal{U}}_h^{l,k}$ ,  $\mathbf{G}_{\mathcal{T}}^k(\hat{\mathbf{v}}_h)|_T := \mathbf{G}_T^k(\hat{\mathbf{v}}_T)$  and  $\mathbf{F}_{\mathcal{T}}^k(\hat{\mathbf{v}}_h)|_T := \mathbf{F}_T^k(\hat{\mathbf{v}}_T)$ . We also define the global stabilization form by  $s_h(\hat{\mathbf{v}}_h, \hat{\mathbf{w}}_h) := \sum_{T \in \mathcal{T}_h} s_T(\hat{\mathbf{v}}_T, \hat{\mathbf{w}}_T)$ , for all  $\hat{\mathbf{v}}_h, \hat{\mathbf{w}}_h \in \widehat{\mathcal{U}}_h^{l,k}$ .

At the discrete level in space, the first Piola–Kirchhoff stress tensor is denoted  $\mathbf{\Pi}_{\mathcal{T}}^n$  and is defined at all the quadrature points  $\{\boldsymbol{\xi}_{T,i}\}_{i \in \{1:N_q\}}$  of any mesh cell  $T \in \mathcal{T}_h$ , with  $N_q$  quadrature points per cell. Thus, we write  $\mathbf{\Pi}_{\mathcal{T}}^n = (\mathbf{\Pi}_{T,i}^n)_{T \in \mathcal{T}_h, i \in \{1:N_q\}}$ . Similarly, the generalized internal variables  $\chi_{\mathcal{T}}^n$  are defined at all the quadrature points, and we write  $\chi_{\mathcal{T}}^n = (\chi_{T,i}^n)_{T \in \mathcal{T}_h, i \in \{1:N_q\}}$ . These two quantities are computed from the pointwise equation

$$(\chi_{T,i}^n, \mathbf{\Pi}_{T,i}^n) = \mathcal{P}(\chi_{T,i}^{n-1}, \mathbf{F}_{\mathcal{T}}^k(\hat{\mathbf{u}}_h^n)(\boldsymbol{\xi}_{T,i}), \mathbf{F}_{\mathcal{T}}^k(\hat{\mathbf{u}}_h^{n-1})(\boldsymbol{\xi}_{T,i})), \quad \forall T \in \mathcal{T}_h, \forall i \in \{1:N_q\}. \quad (6.76)$$

This equation is the fully discrete counterpart of (6.72b). The fully discrete counterpart of (6.72a) is

$$\frac{1}{\Delta t^2} (\mathbf{u}_{\mathcal{T}}^{n+1} - 2\mathbf{u}_{\mathcal{T}}^n + \mathbf{u}_{\mathcal{T}}^{n-1}, \mathbf{w}_{\mathcal{T}})_{\Omega} + (\mathbf{\Pi}_{\mathcal{T}}^n, \mathbf{G}_{\mathcal{T}}^k(\hat{\mathbf{w}}_h))_{\Omega, \mathcal{Q}} + s_h(\hat{\mathbf{u}}_h^n, \hat{\mathbf{w}}_h) = (\mathbf{f}^n, \mathbf{w}_{\mathcal{T}})_{\Omega}, \quad \forall \hat{\mathbf{w}}_h \in \widehat{\mathcal{U}}_{h,0}^{l,k}, \quad (6.77)$$

recalling that  $\mathbf{f}^n := \mathbf{f}(t^n)$  and where we have set

$$(\mathbf{\Pi}_{\mathcal{T}}^n, \mathbf{G}_{\mathcal{T}}^k(\hat{\mathbf{w}}_h))_{\Omega, \mathcal{Q}} := \sum_{T \in \mathcal{T}_h} \sum_{i \in \{1:N_q\}} \omega_{T,i} \mathbf{\Pi}_{T,i}^n : \mathbf{G}_T^k(\hat{\mathbf{w}}_T)(\boldsymbol{\xi}_{T,i}), \quad (6.78)$$

with  $\{\omega_{T,i}\}_{i \in \{1:N_q\}}$  the quadrature weights at every mesh cell  $T \in \mathcal{T}_h$ . The resolution of (6.77) is done in two steps:

1. Find  $\mathbf{u}_{\mathcal{F}}^n$  and  $\mathbf{\Pi}_{\mathcal{T}}^n$  using (6.76) and (6.77) with the test function  $(\mathbf{0}, \mathbf{w}_{\mathcal{F}})$ ,  $\forall \mathbf{w}_{\mathcal{F}} \in \mathcal{U}_{\mathcal{F},0}^k$ ,
2. Find  $\mathbf{u}_{\mathcal{T}}^{n+1}$  using (6.77) with the test function  $(\mathbf{w}_{\mathcal{T}}, \mathbf{0})$ ,  $\forall \mathbf{w}_{\mathcal{T}} \in \mathcal{U}_{\mathcal{T}}^l$ .

Step 1 requires to solve a nonlinear equation because both  $\mathbf{F}_{\mathcal{T}}^k(\hat{\mathbf{u}}_h^n)$  and  $\mathbf{\Pi}_{\mathcal{T}}^n$  are unknown. This can be done using a Newton algorithm or using the splitting method detailed below (and inspired from the method introduced in Chapter 4).

### 6.4.2.3 Algebraic formulation

We recall that  $N_{\mathcal{T}} := \dim(\mathcal{U}_{\mathcal{T}}^l)$ ,  $N_{\mathcal{F}} := \dim(\mathcal{U}_{\mathcal{F},0}^k)$  and that  $\{\boldsymbol{\phi}_i\}_{i \in \{1:N_{\mathcal{T}}\}}$ ,  $\{\boldsymbol{\psi}_i\}_{i \in \{1:N_{\mathcal{F}}\}}$  are bases of  $\mathcal{U}_{\mathcal{T}}^l$  and  $\mathcal{U}_{\mathcal{F},0}^k$ , respectively. We denote  $\hat{\mathbf{U}}_h^n := (\mathbf{U}_{\mathcal{T}}^n, \mathbf{U}_{\mathcal{F}}^n) \in \mathbb{R}^{N_{\mathcal{T}}} \times \mathbb{R}^{N_{\mathcal{F}}}$  the vector of degrees of freedom of the solution  $\hat{\mathbf{u}}_h^n$  on these bases, and  $\mathbf{F}_{\mathcal{T}}^n \in \mathbb{R}^{N_{\mathcal{T}}}$  the vector having components  $((\mathbf{f}^n, \boldsymbol{\phi}_i)_{\Omega})_{i \in \{1:N_{\mathcal{T}}\}}$ .  $\mathcal{M}$  denotes the mass matrix associated with the vector-valued displacement on the cells, and  $\mathcal{S}$  denotes the stabilization matrix associated with the bilinear form  $s_h$ . We define the vectors

$$\mathbf{P}_{\mathcal{T}}^n := ((\mathbf{\Pi}_{\mathcal{T}}^n, \mathbf{G}_{\mathcal{T}}^k(\boldsymbol{\phi}_i, \mathbf{0}))_{\Omega, \mathcal{Q}})_{i \in \{1:N_{\mathcal{T}}\}}, \quad \mathbf{P}_{\mathcal{F}}^n := ((\mathbf{\Pi}_{\mathcal{T}}^n, \mathbf{G}_{\mathcal{T}}^k(\mathbf{0}, \boldsymbol{\psi}_i))_{\Omega, \mathcal{Q}})_{i \in \{1:N_{\mathcal{F}}\}}. \quad (6.79)$$

The algebraic formulation reads

$$\frac{1}{\Delta t^2} \begin{bmatrix} \mathcal{M} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{pmatrix} \mathbf{U}_{\mathcal{T}}^{n+1} - 2\mathbf{U}_{\mathcal{T}}^n + \mathbf{U}_{\mathcal{T}}^{n-1} \\ \cdot \end{pmatrix} + \begin{pmatrix} \mathbf{P}_{\mathcal{T}}^n \\ \mathbf{P}_{\mathcal{F}}^n \end{pmatrix} + \begin{bmatrix} \mathcal{S}_{\mathcal{T}\mathcal{T}} & \mathcal{S}_{\mathcal{T}\mathcal{F}} \\ \mathcal{S}_{\mathcal{F}\mathcal{T}} & \mathcal{S}_{\mathcal{F}\mathcal{F}} \end{bmatrix} \begin{pmatrix} \mathbf{U}_{\mathcal{T}}^n \\ \mathbf{U}_{\mathcal{F}}^n \end{pmatrix} = \begin{bmatrix} \mathbf{F}_{\mathcal{T}}^n \\ \mathbf{0} \end{bmatrix}. \quad (6.80)$$

As a consequence of the structure of the global mass matrix, the face component is replaced by a “.” in the acceleration term. The submatrix  $\mathcal{S}_{\mathcal{T}\mathcal{T}}$  is block-diagonal, since  $s_h$  does not couple cell degrees of freedom from different cells, and, in the mixed-order setting,  $\mathcal{S}_{\mathcal{F}\mathcal{F}}$  is block-diagonal as well. Finally, as in the linear case, the cell mass matrix  $\mathcal{M}$  is also block-diagonal.

#### 6.4.2.4 Splitting procedure in the mixed-order setting

The mixed-order setting offers a convenient operator splitting, which consists, for  $m \geq 0$ , in finding  $\mathbf{u}_{\mathcal{F}}^{n,m+1}$  such that

$$s_h((\mathbf{0}, \mathbf{u}_{\mathcal{F}}^{n,m+1}), (\mathbf{0}, \mathbf{w}_{\mathcal{F}})) = -(\mathbf{\Pi}_{\mathcal{T}}^{n,m}, \mathbf{G}_{\mathcal{T}}^k((\mathbf{0}, \mathbf{w}_{\mathcal{F}})))_{\Omega, \mathcal{Q}} - s_h((\mathbf{u}_{\mathcal{T}}^n, \mathbf{0}), (\mathbf{0}, \mathbf{w}_{\mathcal{F}})), \quad (6.81)$$

for all  $\mathbf{w}_{\mathcal{F}} \in \mathcal{U}_{\mathcal{F},0}^k$ , where  $\mathbf{\Pi}_{\mathcal{T}}^{n,m}$  is the first Piola–Kirchhoff stress tensor obtained by solving

$$(\chi_{T,i}^{n,m}, \mathbf{\Pi}_{T,i}^{n,m}) = \mathcal{P}(\chi_{T,i}^{n-1}, \mathbf{F}_{\mathcal{T}}^k(\mathbf{u}_{\mathcal{T}}^n, \mathbf{u}_{\mathcal{F}}^{n,m})(\boldsymbol{\xi}_{T,i}), \mathbf{F}_{\mathcal{T}}^k(\hat{\mathbf{u}}_h^{n-1})(\boldsymbol{\xi}_{T,i})), \quad \forall T \in \mathcal{T}_h, \forall i \in \{1:N_q\}. \quad (6.82)$$

The algebraic form of (6.81) is as follows: Setting  $\mathbf{U}_{\mathcal{F}}^{n,0} := \mathbf{U}_{\mathcal{F}}^{n-1}$ , we seek  $\mathbf{U}_{\mathcal{F}}^{n,m+1} \in \mathbb{R}^{N_{\mathcal{F}}}$  such that

$$\mathcal{S}_{\mathcal{F}\mathcal{F}} \mathbf{U}_{\mathcal{F}}^{n,m+1} = -\mathbf{P}_{\mathcal{F}}^{n,m} - \mathcal{S}_{\mathcal{F}\mathcal{T}} \mathbf{U}_{\mathcal{T}}^n, \quad (6.83)$$

with  $\mathbf{P}_{\mathcal{F}}^{n,m} := ((\mathbf{\Pi}_{\mathcal{T}}^{n,m}, \mathbf{G}_{\mathcal{T}}^k(\mathbf{0}, \boldsymbol{\psi}_i))_{\Omega, \mathcal{Q}})_{i \in \{1:N_{\mathcal{F}}\}}$ . The above splitting procedure is computationally effective since, as mentioned above, the face-face stabilization submatrix  $\mathcal{S}_{\mathcal{F}\mathcal{F}}$  is block-diagonal.

#### 6.4.2.5 Splitting procedure in the equal-order setting

Let  $\zeta_T$  be the local bilinear form such that, for all  $T \in \mathcal{T}_h$  and all  $\mathbf{v}_{\partial T}, \mathbf{w}_{\partial T} \in \mathcal{U}_{\partial T}^k$ ,

$$\begin{aligned} \zeta_T((\mathbf{0}, \mathbf{v}_{\partial T}), (\mathbf{0}, \mathbf{w}_{\partial T})) &= \gamma \mu \sum_{F \in \mathcal{F}_T} \eta_{TF}^{-1} \left\{ ((I - \mathbf{\Pi}_T^k) \mathbf{R}_T^{k+1}(\mathbf{0}, \mathbf{v}_{\partial T})|_F, \mathbf{w}_F)_F \right. \\ &\quad + (\mathbf{v}_F, (I - \mathbf{\Pi}_T^k) \mathbf{R}_T^{k+1}(\mathbf{0}, \mathbf{w}_{\partial T})|_F)_F \\ &\quad \left. + (\mathbf{\Pi}_F^k (I - \mathbf{\Pi}_T^k) \mathbf{R}_T^{k+1}(\mathbf{0}, \mathbf{v}_{\partial T})|_F, \mathbf{\Pi}_F^k (I - \mathbf{\Pi}_T^k) \mathbf{R}_T^{k+1}(\mathbf{0}, \mathbf{w}_{\partial T})|_F)_F \right\}, \end{aligned} \quad (6.84)$$

and let  $s_T^*$  be the local bilinear form defined by (6.51). Then the equal-order stabilization form writes

$$s_T = s_T^* + \zeta_T. \quad (6.85)$$

Let us introduce the global bilinear forms  $s_h^*((\mathbf{0}, \mathbf{v}_{\mathcal{F}}), (\mathbf{0}, \mathbf{w}_{\mathcal{F}})) := \sum_{T \in \mathcal{T}_h} s_T^*((\mathbf{0}, \mathbf{v}_{\partial T}), (\mathbf{0}, \mathbf{w}_{\partial T}))$  and  $\zeta_h((\mathbf{0}, \mathbf{v}_{\mathcal{F}}), (\mathbf{0}, \mathbf{w}_{\mathcal{F}})) := \sum_{T \in \mathcal{T}_h} \zeta_T((\mathbf{0}, \mathbf{v}_{\partial T}), (\mathbf{0}, \mathbf{w}_{\partial T}))$ , so that  $s_h = s_h^* + \zeta_h$ . This leads to the following iterative procedure, with the same initial condition as for the mixed-order setting: For all  $m \geq 0$ , find  $\mathbf{u}_{\mathcal{F}}^{n,m+1} \in \mathcal{U}_{\mathcal{F},0}^k$  such that

$$s_h^*((\mathbf{0}, \mathbf{u}_{\mathcal{F}}^{n,m+1}), (\mathbf{0}, \mathbf{w}_{\mathcal{F}})) = -(\mathbf{\Pi}_{\mathcal{T}}^{n,m}, \mathbf{G}_{\mathcal{T}}^k(\mathbf{0}, \mathbf{w}_{\mathcal{F}}))_{\Omega, \mathcal{Q}} - \zeta_h((\mathbf{0}, \mathbf{u}_{\mathcal{F}}^{n,m}), (\mathbf{0}, \mathbf{w}_{\mathcal{F}})) - s_h((\mathbf{u}_{\mathcal{T}}^n, \mathbf{0}), (\mathbf{0}, \mathbf{w}_{\mathcal{F}})), \quad (6.86)$$

for all  $\mathbf{w}_{\mathcal{F}} \in \mathcal{U}_{\mathcal{F},0}^k$ , where  $\mathbf{\Pi}_{\mathcal{T}}^{n,m}$  is obtained by solving (6.82). At the algebraic level, we define two matrices  $\mathcal{S}_{\mathcal{F}\mathcal{F}}^*$  and  $\mathcal{Z}_{\mathcal{F}\mathcal{F}}$  such that  $\mathcal{S}_{\mathcal{F}\mathcal{F}} = \mathcal{S}_{\mathcal{F}\mathcal{F}}^* + \mathcal{Z}_{\mathcal{F}\mathcal{F}}$ ,  $\mathcal{Z}_{\mathcal{F}\mathcal{F}}$  corresponds to the bilinear form  $\zeta_h$  and  $\mathcal{S}_{\mathcal{F}\mathcal{F}}^*$  to  $s_h^*$ . Then the splitting procedure (6.86) translates into the following iterative algorithm: For all  $m \geq 0$ , find  $\mathbf{U}_{\mathcal{F}}^{n,m+1} \in \mathbb{R}^{N_{\mathcal{F}}}$  such that

$$\mathcal{S}_{\mathcal{F}\mathcal{F}}^* \mathbf{U}_{\mathcal{F}}^{n,m+1} = -\mathbf{P}_{\mathcal{F}}^{n,m} - \mathcal{Z}_{\mathcal{F}\mathcal{F}} \mathbf{U}_{\mathcal{F}}^{n,m} - \mathcal{S}_{\mathcal{F}\mathcal{T}} \mathbf{U}_{\mathcal{T}}^n. \quad (6.87)$$

**Remark 6.4.5** (Convergence of the splitting in the equal-order setting). Since the stabilization is the same as the one used for linear elasticity, the condition  $\alpha < 1$  is verified for all the polynomial orders  $k \in \{1, 2, 3\}$  and all the cell shapes considered in this chapter.  $\square$



## 6.5 Numerical experiments with finite-strain plasticity

We investigate in this section the feasibility and the computational efficiency of the splitting procedure, compared to the semi-implicit method with a Newton algorithm. We first study the impact of plasticity on the splitting parameter  $\gamma$ , on the CFL condition and on the number of splitting iterations. In particular, we compare the values of  $\gamma^*$  (the smallest value of  $\gamma$  such that the splitting procedure converges) to the values obtained for linear elasticity. Then, we compare the execution time and the precision of the splitting procedure with those for the semi-implicit method. The numerical experiments are carried out with MANTA, which uses the opensource software `Mfront` [111] for the integration of the behavior laws (documentation and code available on the dedicated web page). More precisely, at each integration point, given the deformation gradient and the internal variables, `Mfront` returns the chosen stress measure, here the first Piola–Kirchhoff stress tensor and the updated internal variables. We use the same quadrature order for both the plasticity and the inertial terms. The quadrature order is chosen so that there is no quadrature error for the inertial term, i.e.  $2k$  in the equal-order setting and  $2k + 2$  in the mixed-order setting.

### 6.5.1 Numerical study of the splitting parameters

**Impact of plasticity on the splitting.** In a general nonlinear case, there is no analytical value for  $\gamma^*$ . It is only possible to determine a value which depends on the test case, and which is the smallest value such that the splitting procedure converges on this precise test case. The value of  $\gamma^*$  depends on the polynomial order, the type of mesh, but also the material behavior and the loading. We start by computing  $\gamma^*$  on a given mesh for different values of  $k$ . We use a hyperelastic behavior as well as a linear isotropic hardening behavior with a von Mises criterion. This plasticity model involves two parameters: the isotropic hardening,  $H$ , and the yield stress,  $\sigma_{y,0}$ . We consider a 2D square domain  $\Omega := (0, 1)^2$  meshed with squares of size  $h = 0.1\text{m}$  and make the hypothesis of plane strain, i.e. there is no strain along the  $z$ -axis. The body is made of the reference material ( $\lambda = \mu = 1$ ,  $\rho = 1$ ) and is subject to a surface force  $\mathbf{f}$  uniformly applied on the top edge and pulling upwards, homogeneous Dirichlet boundary conditions on the bottom edge, the other edges being constraint-free. The surface force  $\mathbf{f}$  is a function of time: in the time interval  $t \in (0, 1)\text{s}$ , the force linearly increases, i.e.  $\mathbf{f}(t) := [0, f_{\max}t]^\top$  with  $f_{\max} := 0.2$  N and in the time interval  $t \in (1, 5)\text{s}$ , the surface force is constant,  $\mathbf{f}(t) := [0, f_{\max}]^\top$ . Figure 6.19a shows the 2D HHO solution at  $t = 1.42\text{s}$  for  $k = 2$  in the equal-order setting. Figures 6.19b and 6.19c show the  $x$ - and  $y$ -displacement of the point  $A$  of initial coordinates  $(1, 1)$  over time, for five sets of plasticity parameters  $(H, \sigma_{y,0})$ . These figures illustrate the impact of both parameters on the observed displacement: smaller values of  $H$  represent a softer plastic material and larger values of  $\sigma_{y,0}$  reduce the plastic strain.

We are interested in the value of  $\gamma^*$  for the five pairs of plastic parameters considered in Figure 6.19. Table 6.12 displays the smallest integer such that the splitting procedure converges on the mesh of size  $h = 0.1\text{m}$  for  $k \in \{1, 2, 3\}$ , equal- and mixed-order settings, for the five parameter pairs. We consider that the splitting converges if less than 1000 iterations are needed at each time step. If more iterations are required, the splitting procedure is deemed to be computationally inefficient. The smallest value of  $\gamma^*$  is obtained for the hyperelastic behavior. Moreover, all the reported values of  $\gamma^*$  are equivalent or smaller than the values of  $\gamma^*$  for the linear elastic problem in the mixed-order case (see Table 4.1). Another interesting observation is that, unlike what was observed for the acoustic wave equation, there is a large difference between the value of  $\gamma^*$  for mixed- and equal-order settings with the same  $k$ . The mixed-order setting yields values of  $\gamma^*$  that are 1.5 to 2 times larger than the value of  $\gamma^*$  for the equal-order case.

Regarding the impact of the plastic behavior, both  $H$  and  $\sigma_{y,0}$  influence the value of  $\gamma^*$ . The larger the plastic strain (and hence the displacement observed in Figure 6.19), the larger  $\gamma^*$  (i.e. the larger the scaling of the stabilization needed for the splitting to converge). This is in agreement with the behavior observed for the nonlinear acoustic wave equation: the more nonlinear the behavior, the larger  $\gamma$  needs to be. The values of  $\gamma^*$  reported in Table 6.12 remain, however, of the same order of magnitude as for the

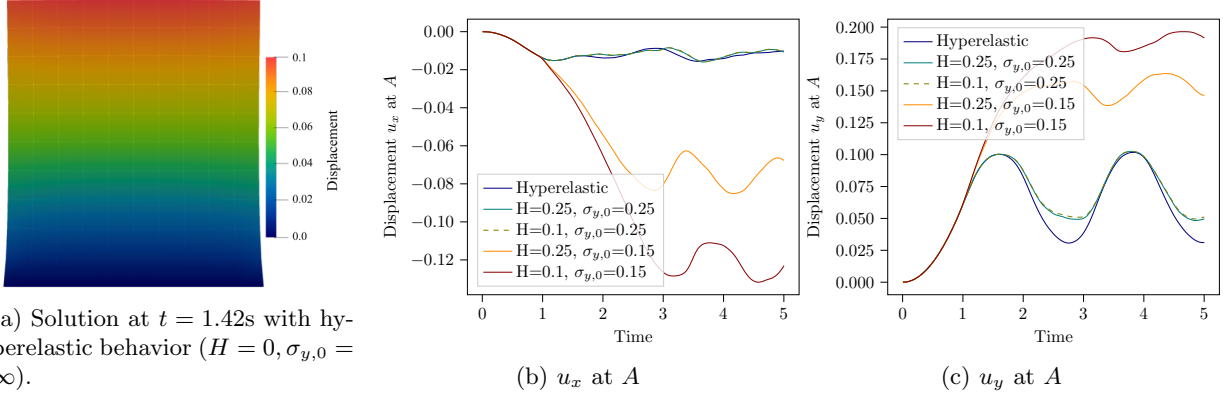


Figure 6.19: 2D traction finite-strain experiment (reference material) - Setting and solution on a  $10 \times 10$  mesh,  $(l, k) = (1, 1)$  for various values of the plasticity parameters  $H$  and  $\sigma_{y,0}$ .

hyperelastic behavior. Thus, a practical conclusion is that, when considering a new test case, experiments on a small mesh with hyperelastic behavior may be useful in order to determine a lower bound for  $\gamma^*$ .

Polynomial order	(1,1)	(2,2)	(3,3)	(2,1)	(3,2)	(4,3)
Hyperelasticity	13	27	50	25	49	81
$H = 0.25, \sigma_{y,0} = 0.25$	13	27	52	25	49	82
$H = 0.1, \sigma_{y,0} = 0.25$	13	27	52	25	49	82
$H = 0.25, \sigma_{y,0} = 0.15$	14	29	57	28	56	94
$H = 0.1, \sigma_{y,0} = 0.15$	18	32	63	31	62	103

Table 6.12: 2D traction with finite-strain (reference material) - Value of  $\gamma^*$  depending on the plasticity parameters (in GPa), polynomial orders  $k \in \{1, 2, 3\}$ , equal- and mixed-order settings.

It is interesting to look at the number of splitting iterations at each time step for the previous experiment. Figure 6.20 reports this number in the time interval  $(0, 5)$ s for the hyperelastic model and the plasticity model with  $H = 0.1$ GPa and  $\sigma_{y,0} = 0.15$ GPa. In the hyperelastic case, all the polynomial orders lead to the same behavior: the number of splitting iterations is maximal when the deformation at point  $A$  is maximal ( $t \approx 1.5$ s and  $t \approx 3.8$ s). The largest number of splitting iterations is between 100 and 300 depending on the polynomial order. This is a relatively large number of iterations, which indicates that, in the hyperelastic case, the value of  $\gamma^*$  depends on the maximal strain. In the plastic case, the number of splitting iterations does not display the same behavior. For all polynomial orders except  $(l, k) = (1, 1)$ , the number of splitting iterations increases swiftly in the time interval  $t \in (2, 3)$ s, and then stays relatively high. This interval corresponds to the time when plastic deformation occurs, see Figure 6.20c which displays statistics of the equivalent cumulative plastic strain in the entire body over the time interval  $(0, 5)$ s. During the time interval  $(2, 3)$ s, the maximal plastic strain increases, as well as the variations of the plastic stain in the entire body, measured by the interquartile range (interval between the 25th and the 75th percentiles). Thus, the splitting procedure requires more iterations during large plastic deformations. The results reported in Figure 6.20 also indicate that some computational efficiency could be gained if the value of  $\gamma$  could be tuned on the fly during the simulation, increasing it during large plastic deformations and reducing it otherwise. How to realize this in an automated fashion is left to future work.

**Impact of mesh type.** We observed in the case of linear acoustic and elastic waves that the shape of the mesh cells influences the value of  $\gamma^*$ . It is thus expected that the value of  $\gamma^*$  should change with the shape of the cell here as well. To study this, we keep the same setting as in the previous test case

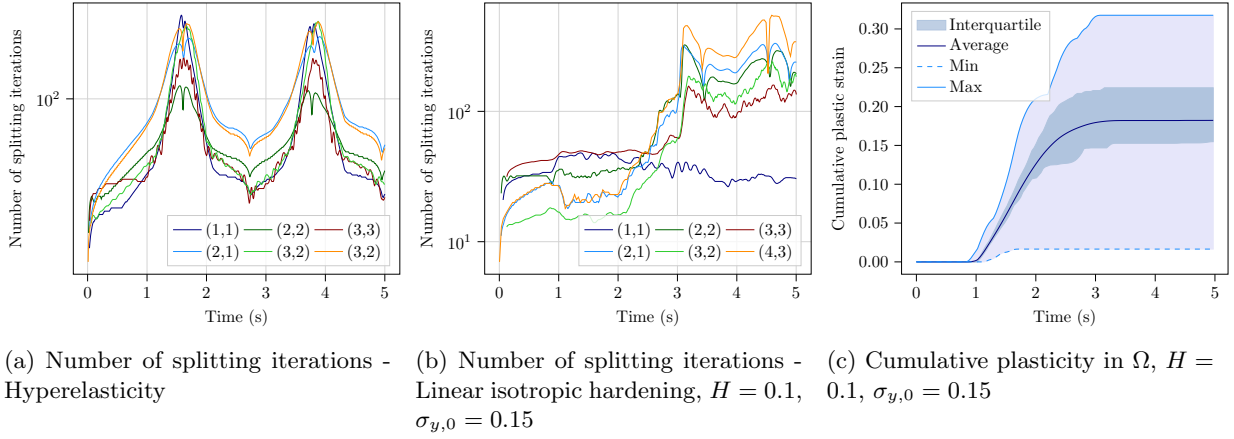


Figure 6.20: 2D traction with finite-strain - Number of splitting iterations at each time step for  $\gamma = \gamma^*$  and cumulative plastic strain, polynomial orders  $k \in \{1, 2, 3\}$ , mixed- and equal-order settings.

and focus on linear isotropic hardening with  $H = 0.1\text{GPa}$  and  $\sigma_{y,0} = 0.15\text{GPa}$ . Table 6.13 displays the value of  $\gamma^*$  for square cells (the same values as in Table 6.12, these values are recalled for convenience), right-triangular cells, and unstructured triangular and quadrangular cells. For all the meshes, the mesh size is  $h = 0.1\text{m}$ . We observe that, as in the linear case, the value of  $\gamma^*$  is larger for triangular meshes, and increases for all meshes with the polynomial order. We also observe that, for all meshes, the difference between the value of  $\gamma^*$  for mixed- and equal-order settings is more pronounced than for the acoustic wave equation.

Polynomial order	(1,1)	(2,2)	(3,3)	(2,1)	(3,2)	(4,3)
Squares	18	32	63	31	62	103
Right triangles	33	55	69	105	133	175
Unstructured quadrangles	29	45	60	89	104	145
Unstructured triangles	23	35	48	68	84	116

Table 6.13: 2D traction with finite-strain (reference material, isotropic linear hardening  $H = 0.1\text{GPa}$ ,  $\sigma_{y,0} = 0.15\text{GPa}$ ) - Value of  $\gamma^*$  for different mesh shapes,  $h = 0.1\text{m}$ , polynomial orders  $k \in \{1, 2, 3\}$ , equal- and mixed-order settings.

**Impact of material parameters.** We compute the value of  $\gamma^*$  in the same manner as before for the materials of Table 6.1 (we exclude rubber since it is nearly incompressible). We consider the domain  $\Omega := (0,1)^2$ , meshed with squares of size  $h := 0.1$ , and the same boundary conditions as in the previous test case. We now take  $f_{\max} := 2 \times 10^5\text{N}$  and  $\mathfrak{T} := 10\text{ms}$ . The loading is increased linearly in the time interval  $(0,5)\text{ms}$  and remains constant in the time interval  $(5,10)\text{ms}$ . We consider isotropic linear strain-hardening plasticity with the parameters  $(H, \sigma_{y,0})$  reported in Table 6.14, the other material properties (density, Young modulus and Poisson ratio) are being reported in Table 6.1. Figure 6.21 displays the

Material	Steel	Copper	Diamond	Gold
Hardening slope $H$ (GPa)	0.13	0.1	20	0.2
Yield strength $\sigma_{y,0}$ (GPa)	0.45	0.4	2.8	0.02

Table 6.14: Plastic properties for a selection of metals (linear isotropic hardening): hardening slope and yield strength.

displacement of point  $A$  of initial coordinates  $(1, 1)$  as a function of time for the four materials considered. The displacement  $u_y$  is displayed in logarithmic scale, so that the result for diamond is visible. The behavior of steel and copper is, as expected, very similar, since both materials have similar material parameters, with a slightly larger displacement for copper. Gold plastifies for a much smaller strain, and diamond for a much larger strain, with a very small displacement. The values of  $\gamma^*$  for HHO with polynomial orders  $k \in \{1, 2, 3\}$  reported in Table 6.15 show that the value of  $\gamma^*$  is directly linked to the largest displacement. The value of  $\gamma^*$  for diamond is two times smaller than that of steel, which is also about two times smaller than that of gold. The value of  $\gamma^*$  for copper is about 50% higher than that of steel. This illustrates that the value of  $\gamma^*$  is directly impacted by the material coefficients and the amount of deformation in the body.

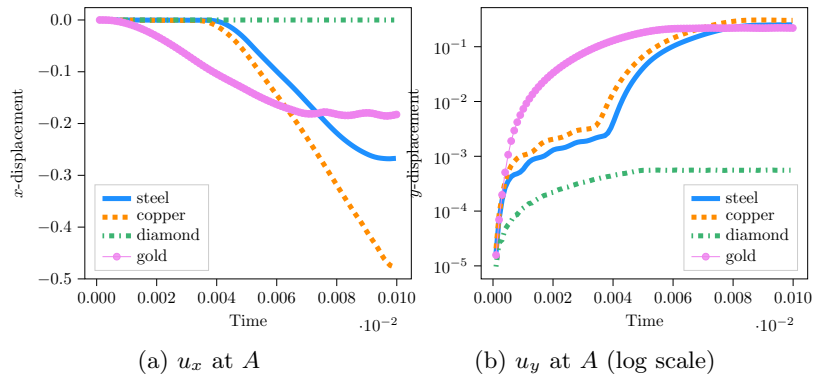


Figure 6.21: 2D traction with finite-strain, linear isotropic hardening - Displacement of point  $A = (1, 1)$  at initial time for different materials, HHO with polynomial orders  $k \in \{1, 2, 3\}$ , mixed- and equal-order settings.

Material	Steel	Copper	Diamond	Gold
<b>(1,1)</b>	20	27	9	37
<b>(2,2)</b>	40	54	20	75
<b>(3,3)</b>	78	112	34	149
<b>(2,1)</b>	39	61	17	81
<b>(3,2)</b>	77	105	34	141
<b>(4,3)</b>	130	175	56	239

Table 6.15: 2D traction finite-strain experiment, linear isotropic hardening - Value of  $\gamma^*$  for plastic behavior on different alloys,  $k \in \{1, 2, 3\}$ , mixed- and equal-order settings.

## 6.5.2 Comparison between the splitting procedure and a Newton scheme

We compare the execution time of the splitting procedure and the Newton scheme. We chose as test case Cook's membrane, since it is known for generating volumetric locking when using a finite element space semi-discretization. This is a well known bending-dominated test case [139, 95, 9]. It consists of a tapered panel, clamped on one side, and subjected to a vertical load applied uniformly on the opposite side. Figure 6.22a illustrates the setting. This test case is usually used in a quasi-static setting. We adapt the test case to a dynamic setting. Let  $\mathfrak{T} := 0.5\text{ms}$ . In the time interval  $(0, 0.25)\text{ms}$ , the force applied on the right side is increased linearly, up to  $\mathbf{f} := [0, f_{\max}]^T$ , with  $f_{\max} := 7\text{kN}$ . In the time interval  $(0.25, 0.5)\text{ms}$ , the applied force has a constant value. The material parameters are those of steel with isotropic linear-hardening plasticity, density and Lamé coefficients from Table 6.1 and plasticity

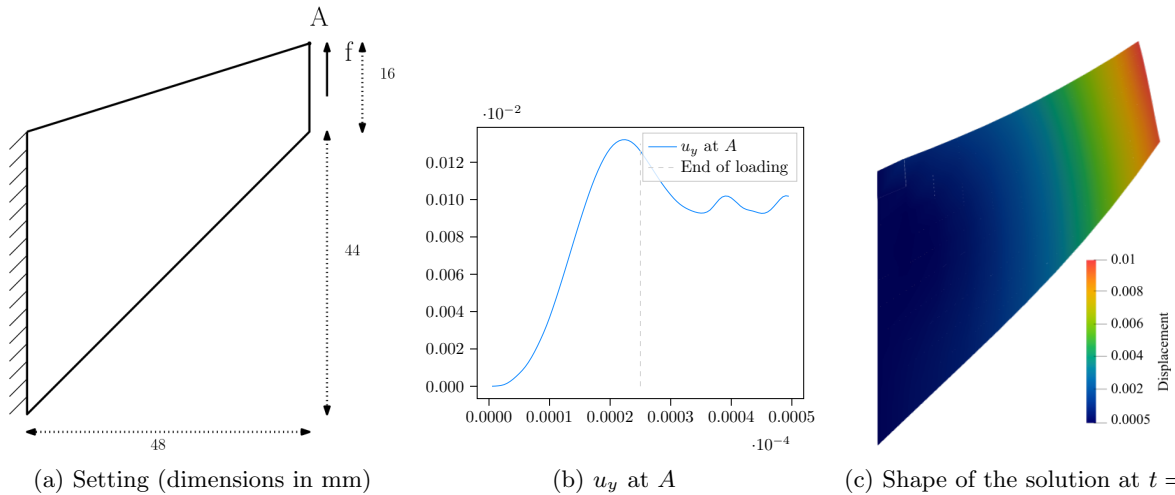


Figure 6.22: Cook's membrane, steel, linear isotropic hardening - Setting, displacement of point  $A$  and 2D solution at  $\mathfrak{T}$ , HHO with polynomial order  $k = 2$  mixed-order setting,  $h = 0.05\text{m}$ .

parameters from Table 6.14. Figure 6.22b displays the  $y$ -displacement of point  $A$  with initial coordinates  $(0.048, 0.06)$ , and Figure 6.22c displays the shape of the membrane at the final time  $\mathfrak{T}$  with the colormap representing the  $y$ -displacement. We first verify that the HHO method does not present volumetric

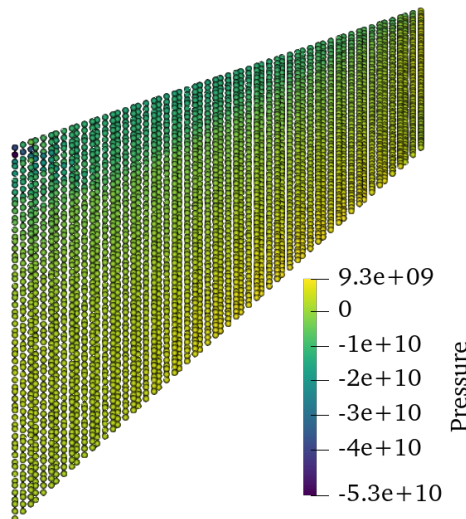


Figure 6.23: Cook's membrane, steel, linear isotropic hardening - Hydrostatic pressure at integration points at final time  $t = 0.5\text{ms}$ , HHO with polynomial order  $k = 2$ , mixed-order setting,  $h = 0.05\text{m}$ .

locking. Figure 6.23 displays the trace of the Cauchy stress tensor (hydrostatic pressure) at the final time  $\mathfrak{T} = 50\text{ms}$ , at each integration point of a mesh with size  $h = 0.05\text{m}$ , for the HHO method with polynomial order  $(l, k) = (3, 2)$ . No oscillations are observed, which is the expected result. The other polynomial orders lead to the same result.

We now focus on the comparison of the splitting procedure with the semi-implicit scheme with a Newton algorithm. Since the setting is different from that of Section 6.5.1, we recompute the splitting parameter  $\gamma^*$  in the same fashion for this test case. We use values of  $\gamma$  at least 30% larger than the value of  $\gamma^*$  to avoid having to perform too many splitting iterations, and consider a convergence criterion

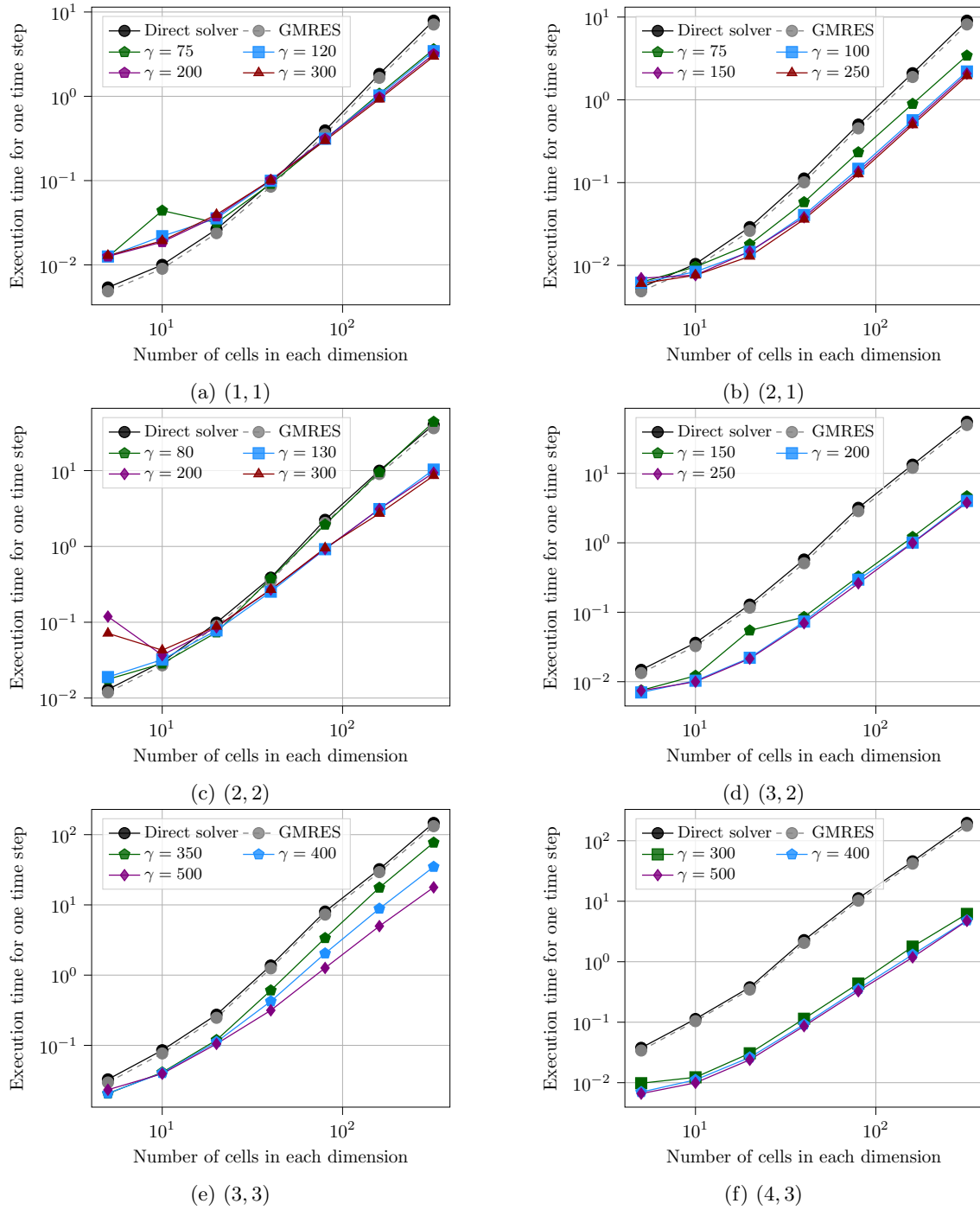


Figure 6.24: Cook's membrane, steel, linear isotropic hardening - Average execution time for one time step, HHO with polynomial orders  $k \in \{1, 2, 3\}$ , mixed- and equal-order settings, using either the semi-implicit scheme or the splitting procedure.

$\epsilon = 10^{-10}$ . Figure 6.24 displays the execution time for one time step using the semi-implicit scheme or the splitting procedure, for a sequence of refined meshes. In the Newton algorithm, the linear system can be solved using either a direct solver or an iterative solver. We test a direct solver and a GMRES

iterative solver with a Jacobi preconditioner. For the splitting procedure, different values of  $\gamma$  are tested. The displayed execution time is the average over the simulation. The behavior is similar to the one observed for the nonlinear acoustic wave equation: on the coarsest meshes, the splitting procedure can be slower than the semi-implicit method with a Newton algorithm. But on the finer meshes, the splitting procedure is always faster. We notice that, for the semi-implicit scheme, the GMRES iterative solver and the direct solver yield similar execution times. The gain in execution time for the splitting on the finest mesh depends on the polynomial order and the value of  $\gamma$ . The largest gain on the finest mesh is reported in Table 6.16. In the equal-order setting, the impact of the value of  $\gamma$  is more pronounced than for the mixed-order setting. The smallest values of  $\gamma$  considered for the splitting procedure yield execution times comparable to that of the semi-implicit procedure for  $k = 2$  in the equal-order setting. In the mixed-order setting for  $k \in \{2, 3\}$ , the three considered values of  $\gamma$  yield similar execution times. As the polynomial order increases, the gain in computation time on the finer meshes increases with the splitting procedure. This experiment illustrates the efficiency of the splitting procedure compared to the

Polynomial order	(1,1)	(2,2)	(3,3)	(2,1)	(3,2)	(4,3)
Gain in execution time	2.3	4.7	8.2	4.6	14.5	30.5

Table 6.16: Cook’s membrane, steel, linear isotropic hardening - Largest gain in computational time using the splitting procedure rather than the Newton algorithm, HHO with polynomial orders  $k \in \{1, 2, 3\}$  mixed- and equal-order settings,  $h = 0.313\text{mm}$ .

semi-implicit method with a Newton algorithm. Since different values of  $\gamma$  are used, we verify that they all yield the same solution. Figure 6.25 displays the displacement of the point  $A$  over time on a (coarse) mesh having 5 cells in each direction, polynomial orders  $k \in \{1, 2, 3\}$ , equal- and mixed-order settings, for all the values of  $\gamma$  used in Figure 6.24. The reference solution of Figure 6.22b is displayed using a dashed gray line. All curves of the same color essentially superimpose, which justifies that, for a given polynomial order, the value of  $\gamma$  does not impact strongly the quality of the solution. Only for polynomial order  $k = 1$  and equal-order setting, the curves do not perfectly align. There is also a larger difference between the mixed- and the equal-order predictions for  $k \in \{1, 2\}$  than for  $k = 3$ . This can be expected since the considered mesh is quite coarse, so that the solutions differ depending on the polynomial order. Higher orders yield a solution closer to the reference.

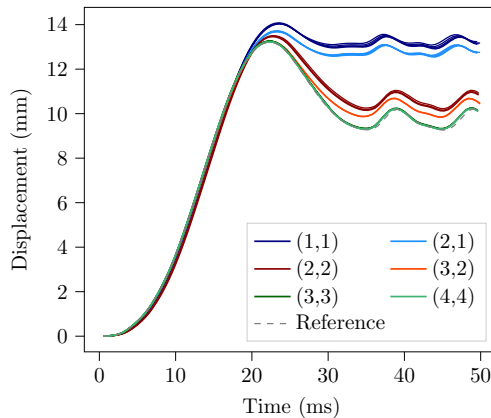


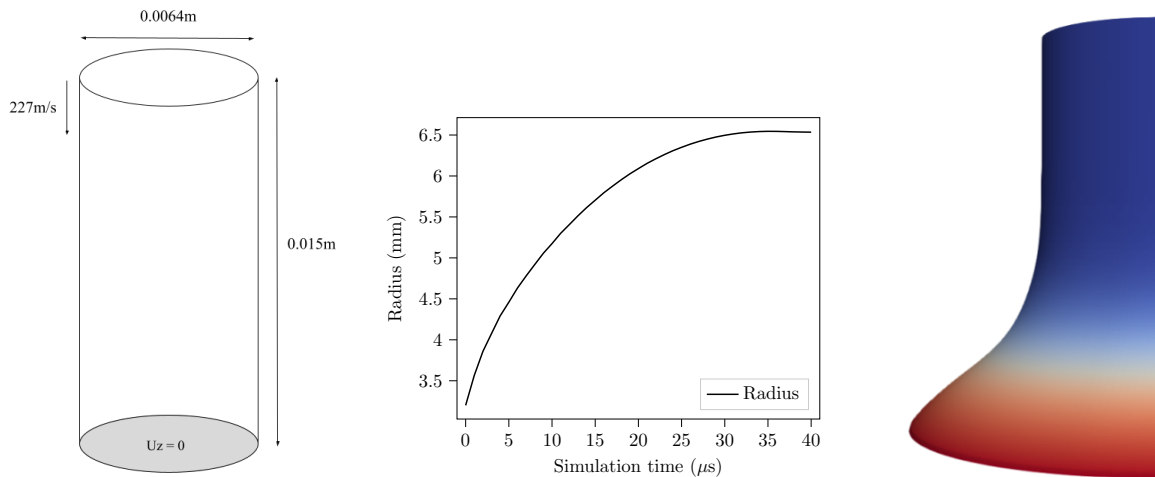
Figure 6.25: Cook’s membrane, steel, linear isotropic hardening - Displacement of point  $A$  over time for polynomial orders  $k \in \{1, 2, 3\}$ , equal- and mixed-order settings, and reference solution.

The numerical experiments of this section illustrated how the splitting procedure can be used when dealing with plastic behaviors. Since this is a nonlinear problem, the value of  $\gamma^*$  depends on the test case, and some preliminary computations are required. But, contrary to the linear elastic case, both

mixed- and equal-order settings yield converging splitting procedures. Moreover, the splitting procedure accelerates the computations by a factor of at least 2, and up to 30 (for the mixed-order setting (4,3)), compared to the semi-implicit scheme with a Newton algorithm. Finally, as in the quasi-static case, the HHO space semi-discretization does not suffer from volumetric locking. These advantages suggest that the fully explicit HHO scheme with the splitting procedure could be a good candidate to replace the classical leapfrog-finite element discretization.

### 6.5.3 Taylor rod test case

The goal of this section is to consider a 3D test case and to compare the splitting procedure for the HHO method with the standard finite element method (FEM). We chose an adaptation of the 3D Taylor rod test case, since it is widely used to assess the efficiency of various devices to counter volumetric locking, see for example [28] for a comparison of tetrahedral and hexahedral elements using reduced integration. In this test case, we consider the impact of a small cylindrical copper bar against a rigid wall. The bar has initial length  $h := 15\text{mm}$  and initial radius  $r := 3.2\text{mm}$ . The initial velocity of the bar is  $v_0 := 227\text{m/s}$  and the final simulation time is  $\mathfrak{T} := 40\mu\text{s}$ . We consider that the bar is made of copper, with the material properties of Table 6.1, considering a linear isotropic plastic behavior. We do not simulate contact between the cylinder and the wall. Instead we consider that the cylinder is blocked by the wall in the  $z$ -direction, see Figure 6.26a. Since the problem has rotation symmetries, we only



(a) Setting (dimensions in mm) (b) Radius of the bottom of the rod over time (c) Final shape of a quarter rod

Figure 6.26: Taylor's rod, linear isotropic hardening - Complete 3D setting, radius of the bottom of the rod over time and 3D shape of the solution at  $t = \mathfrak{T}$ , linear finite elements with reduced integration and anti-hour glass (Europlexus software),  $h = 70\mu\text{m}$ .

simulate a quarter of the rod and apply symmetry boundary conditions (normal component of the face unknowns is set to zero). Since the standard finite element method is subject to volumetric locking, due to the strong plastic deformations, we use as reference a linear finite element method with reduced integration and anti hourglass corrections implemented in the Europlexus software developed at CEA. Figure 6.26c illustrates the final shape of the full rod on a fine mesh of hexahedra with  $h = 70\mu\text{m}$ . We are interested in radius of the bottom of the rod at the final time  $t = \mathfrak{T}$ . Figure 6.26b shows the radius of the bottom of the rod over time. The final radius of the rod on our fine mesh is  $r_f = 6.58\text{mm}$ . This is the reference value for further simulations.

Figure 6.27a displays the solution for a quarter of the rod on our finest mesh with  $h = 70\mu\text{m}$  for the standard finite element implementation, and Figure 6.27b displays the solution obtained with the



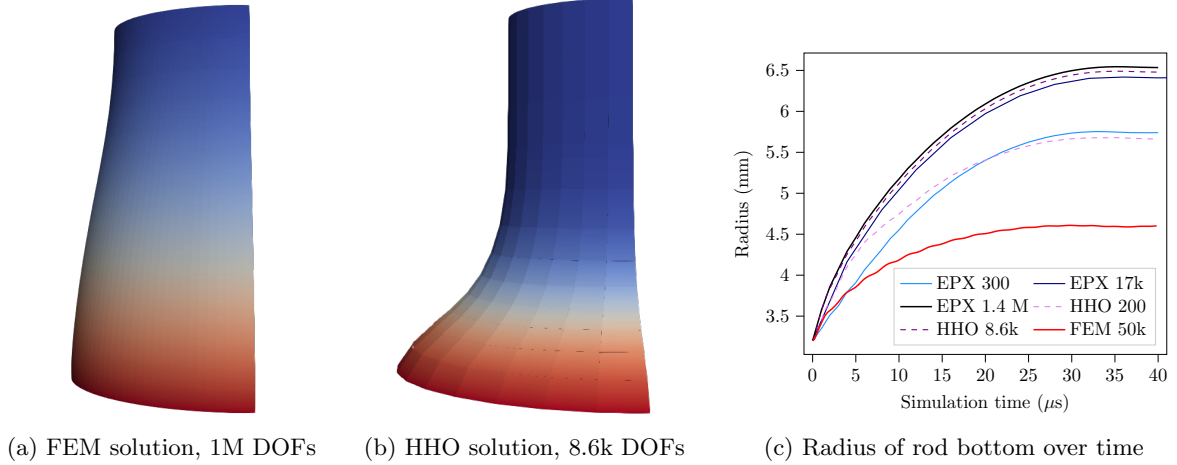
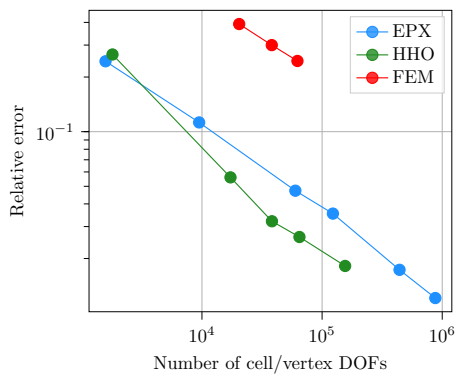


Figure 6.27: Taylor's rod, linear isotropic hardening - Shape of the FEM solution on a mesh containing 50k DOFs, shape of the HHO solution for polynomial order  $k = 1$ , mixed-order setting on a mesh containing 8.6k DOFs and radius over time for different number of DOFs, standard FEM, HHO and FEM with reduced integration.

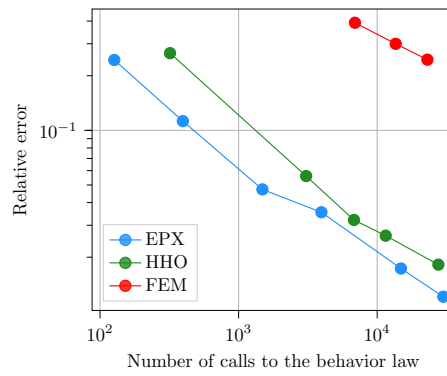
HHO method with the splitting procedure, mixed-order setting,  $k = 1$ ,  $\gamma = 200$  for a much coarser mesh with  $h = 0.5\text{mm}$ . The convergence criterion for the splitting procedure is set to  $\epsilon = 10^{-8}$ . A qualitative comparison shows that the HHO method on a coarse mesh yields a solution much closer to the reference value than the finite element solution. Figure 6.27c displays the radius of the rod bottom over time for different meshes for the HHO method, the standard FEM and the FEM with reduced integration. The HHO solution obtained with 8.6k DOFs is closer to the reference than the solution obtained with the FEM with reduced integration with 17k DOFs.

Figure 6.28a displays the error on the final radius of the rod bottom for the linear standard FEM, the linear FEM with reduced integration and the HHO method in the mixed-order setting with  $k = 1$ , on a series of refined meshes, as a function of the number of DOFs. We consider the number of cell DOFs for the HHO method and the number of vertex DOFs for the FEM method. This illustrates the fast convergence of the HHO method: the same error is obtained with about 5 times less cells for the HHO method than for the FEM with reduced integration.

Comparisons of execution times between the HHO method with the splitting procedure and the finite element method with reduced integration are under study. The HHO method requires quite long execution times due to two factors: a CFL condition about 2 times smaller than that of the finite element method, and an increased number of calls to the behavior law, due to the repeated calls during the splitting iterations and the larger number of integration points. The Aitken acceleration can reduce the number of splitting iterations down to 3 or 4, but the other two factors still lead to at least 10 times more calls to the behavior law. Figure 6.28b displays the relative error as a function of the number of calls to the behavior law. For the FEM with reduced integration, only one integration point per cell is used, whereas 8 integration points are used for the standard FEM and the HHO method. Using this metric, for a given error, the HHO method is slightly slower than the FEM with reduced integration. It could be possible to reduce the number of calls to the behavior law by using less quadrature points for the HHO method (4 points can deliver an exact quadrature for  $\mathbb{P}_3^1$  in a tetrahedron, the 8 points deliver for exact integration of polynomials in  $\mathbb{Q}_3^1$  in a tetrahedron). In this case, the number of calls to the behavior law between the HHO method and FEM with reduced integration for a given error would be comparable.



(a) Relative error vs. of the number of DOFs



(b) Relative error vs. of the number of calls to the behavior law

Figure 6.28: Relative error on the final radius as a function of the number of DOFs and of the number of calls to the behavior law. HHO method with polynomial order  $k = 1$ , mixed-order setting, standard linear FEM and linear FEM with reduced integration.

## Chapter 7

# Conclusion and perspectives

In this Thesis, we introduced an explicit time-scheme for the discretization of wave propagation problems discretized with the Hybrid High-Order method (HHO) in space and the leapfrog (central finite difference) scheme in time. The challenge laid in the existence of a static coupling between the cell and the face unknowns at each time step. This static coupling leads to a semi-implicit time stepping, requiring to solve, at each time step, a linear system in the linear case, and a nonlinear system in the nonlinear case. In order to avoid the cost entailed by this resolution, we devised an explicit HHO method for the acoustic wave equation, both linear and nonlinear, and for the structural dynamics equation, both in the elastic regime and with finite-strain plasticity. This explicit time scheme hinges on an operator splitting and can be viewed as a fixed-point algorithm. Convergence was proved in the linear case and numerically verified for linear and nonlinear equations. The computational efficiency of the proposed method was compared with the standard finite element method, using an implementation in the open-source simulation software MANTA.

Let us summarize the main results of this Thesis. Firstly, a proof of convergence of the fully discrete scheme for the linear acoustic wave equation using the HHO method for the space discretization and the leapfrog scheme for the time discretization was provided. The proof relied on energy and duality arguments. It was presented first in the time-continuous case and then in the fully discrete case.

Secondly, the explicit time-stepping scheme was introduced for the linear acoustic wave equation. The idea was to remove the static coupling existing between the cell and the face unknowns at each time step by means of a splitting of the stiffness operator. For the mixed-order setting, the splitting relied on the face-face part of the stabilization, which is a block-diagonal matrix. In the equal-order setting, the splitting required an additional splitting of the stabilization operator. We showed that the resulting fixed-point algorithm converges if the stabilization bilinear form is scaled by a parameter that is taken to be large enough. Quite importantly, our numerical tests revealed that increasing the weight of the stabilization only mildly affected the solution accuracy. Moreover, the impact on the CFL stability restriction remained moderate.

Thirdly, the splitting procedure was adapted to the nonlinear acoustic wave equation by exploiting the fact that the stabilization is linear. Experiments with a p-structure potential and a vibrating membrane showed good performances compared to the finite element method in terms of accuracy and computational efficiency.

Finally, the splitting procedure was extended to the elastic wave equation and the structural dynamics equation with finite-strain plasticity. The discretization of the elastic wave equation hinged on a symmetric gradient reconstruction operator instead of a plain gradient reconstruction operator. This impacted unfavorably the splitting procedure only in the equal-order case. In the mixed-order case, numerical experiments illustrated the efficiency of the HHO method compared to the finite element method when considering heterogeneous media. For the structural dynamics equation with finite-strain plasticity, the splitting procedure could be used in both equal- and mixed-order settings. Numerical examples demon-

strated the computational efficiency of the splitting procedure compared to a semi-implicit scheme with a Newton algorithm.

The work presented in this Thesis can be pursued in several directions. From a mathematical and numerical point of view, the following directions can be investigated:

- Adaptive mesh refinement (AMR). The HHO method is well-suited for nonconforming mesh refinement. AMR is paramount in the design of efficient high-performance computing. Leveraging the explicit HHO time-stepping developed in this Thesis, adaptive mesh refinement could lead to a very efficient numerical setting for industrial applications.
- *hp*-adaptivity. While AMR deals with *h*-refinement, one can also consider refining the polynomial order *p*. In *hp*-adaptivity, one can refine the polynomial order and the mesh size in different zones of the computational domain. This is another approach to optimize the amount of computations for a given precision, see [75] for an application to  $H^1$ -conforming methods.
- Study of the scheme obtained with the scaling of the stabilization going to infinity. We observed that, for the diffusion equation, increasing the scaling of the stabilization leads to a stable solution and did not increase the error. In the mixed-order setting, we were able to express the face unknowns in terms of the cell unknowns. One perspective is to rewrite the scheme with an infinite scaling of the stabilization as a discontinuous Galerkin method and investigate its mathematical and numerical properties.
- Tuning of the splitting parameter during the simulation. We observed in Chapter 6 that in the plastic case, the number of splitting iterations varies strongly during the simulation. Automatically tuning the splitting parameter  $\gamma$  on the fly during the simulation could reduce the number of splitting iterations and hence the computational cost of the splitting procedure.
- Fully discrete  $L^2$ -error estimation. Chapter 3 only delivered energy-norm estimates for the HHO-leapfrog scheme. Drawing on the time-continuous case, an  $L^2$ -error estimate could possibly be derived.
- Higher-order and local time integration. This Thesis focused on the second-order leapfrog scheme for simplicity. However, higher-order time integration schemes would allow to leverage the high order of convergence of the space semi-discretization. Possible schemes are higher-order finite difference schemes as in [51], leapfrog-Chebyshev schemes [40], or Adams–Bashforth schemes for instance. Another approach to improve computational efficiency is the use of local time stepping, see e.g. [90, 106].

Concerning problems in computational mechanics, we can mention the following directions to continue this work:

- Contact and damage. Considering the final goal of application of accident simulations, contact between parts of the structure or several structures and damages need to be taken into account. Indeed, above a given threshold in large plastic deformations, the material can be damaged, or bend and collide with other parts of the structure. The HHO method has already been applied to the simulation of contact with linear elasticity and Tresca friction for quasi-static problems [61].
- Industrial applications. The present Thesis focused on the mathematical design and the numerical validation on simple cases of the splitting algorithm. Since MANTA is dedicated to industrial applications, it is interesting to pursue further validation test cases on more challenging problems and with more realistic physical assumptions such as large strain. This can lead to the use of finer meshes requiring large scale parallelization or GPU programming.

# Bibliography

- [1] M. Abbas, A. Ern, and N. Pignet. Hybrid high-order methods for finite deformations of hyperelastic materials. *Comput. Mech.*, 62(4):909–928, 2018.
- [2] M. Abbas, A. Ern, and N. Pignet. A hybrid high-order method for finite elastoplastic deformations within a logarithmic strain framework. *Int. J. Numer. Methods Eng.*, 120(3):303–327, 2019.
- [3] M. Abbas, A. Ern, and N. Pignet. A hybrid high-order method for incremental associative plasticity with small deformations. *Comp. Meth. Appl. Mech. Eng.*, 346:891–912, 2019.
- [4] C. Agelet de Saracibar, M. Chiumenti, Q. Valverde, and M. Cervera. On the orthogonal subgrid scale pressure stabilization of finite deformation j2 plasticity. *Comput. Methods Appl. Mech. Eng.*, 195(9-12):1224–1251, 2006.
- [5] J. Aghili, S. Boyaval, and D. A. Di Pietro. Hybridization of mixed high-order methods on general meshes and application to the Stokes equations. *Comput. Methods Appl. Math.*, 15(2):111–134, 2015.
- [6] J. Aghili, D. A. Di Pietro, and B. Ruffini. An hp-hybrid high-order method for variable diffusion on general meshes. *Comput. Methods Appl. Math.*, 17(3):359–376, 2017.
- [7] M. Ainsworth and C. Parker. Unlocking the secrets of locking: Finite element analysis in planar linear elasticity. *Comput. Methods Appl. Mech. Eng.*, 395:115034, 2022.
- [8] A. C. Aitken. IV.—studies in practical mathematics. V. on the iterative solution of a system of linear equations. *Proc. R. Soc. Edinburgh Sec. A: Mathematics*, 63(1):52–60, 1950.
- [9] D. Al Akhrass, J. Bruchon, S. Drapier, and S. Fayolle. Integrating a logarithmic-strain based hyperelastic formulation into a three-field mixed finite element formulation to deal with incompressibility in finite-strain elastoplasticity. *Finite Elem. Anal. Des.*, 86:61–70, 2014.
- [10] P. F. Antonietti, M. Botti, and I. Mazzieri. On mathematical and numerical modelling of multi-physics wave propagation with polytopal discontinuous Galerkin methods: a review. *Vietnam J. Math*, 50(4):997–1028, 2022.
- [11] P. F. Antonietti, I. Mazzieri, and F. Migliorini. A discontinuous Galerkin time integration scheme for second order differential equations with applications to seismic wave propagation problems. *Comput. Math. Appl.*, 134:87–100, 2023.
- [12] P. F. Antonietti, I. Mazzieri, M. Muhr, V. Nikolić, and B. Wohlmuth. A high-order discontinuous Galerkin method for nonlinear sound waves. *J. Comput. Phys.*, 415:109484, 2020.
- [13] D. Appelö and T. Hagstrom. A new discontinuous Galerkin formulation for wave equations in second-order form. *SIAM J. Numer. Anal.*, 53(6):2705–2726, 2015.
- [14] D. N. Arnold, F. Brezzi, B. Cockburn, and L. D. Marini. Unified analysis of discontinuous Galerkin methods for elliptic problems. *SIAM J. Numer. Anal.*, 39:1749–1779, 2002.

- [15] B. Ayuso de Dios, K. Lipnikov, and G. Manzini. The nonconforming virtual element method. *ESAIM Math. Mod. Numer. Anal.*, 50(3):879–904, 2016.
- [16] I. Babuška and M. Suri. Locking effects in the finite element approximation of elasticity problems. *Numer. Math.*, 62(1):439–463, 1992.
- [17] I. Babuška and M. Suri. On locking and robustness in the finite element method. *SIAM J. Numer. Anal.*, 29(5):1261–1293, 1992.
- [18] G. A. Baker. Error estimates for finite element methods for second order hyperbolic equations. *SIAM J. Numer. Anal.*, 13(4):564–576, 1976.
- [19] H. Barucq, H. Calandra, J. Diaz, and E. Shishenina. Space–time Trefftz-DG approximation for elasto-acoustics. *Appl. Anal.*, 99(5):747–760, 2020.
- [20] Y. Bazilevs, L. Beirão da Veiga, J. A. Cottrell, T. J. R. Hughes, and G. Sangalli. Isogeometric analysis: approximation, stability and error estimates for h-refined meshes. *Math. Models Methods Appl. Sci.*, 16(07):1031–1090, 2006.
- [21] L. Beirão da Veiga, Franco Brezzi, and L. D. Marini. Virtual elements for linear elasticity problems. *SIAM J. Numer. Anal.*, 51(2):794–812, 2013.
- [22] L. Beirão da Veiga, C. Lovadina, and D. Mora. A virtual element method for elastic and inelastic problems on polytope meshes. *Comput. Methods Appl. Mech. Eng.*, 295:327–346, 2015.
- [23] L. Beirão da Veiga, F. Brezzi, L. D. Marini, and A. Russo. *Virtual Element Implementation for General Elliptic Equations*, pages 39–71. Lecture Notes in Computational Science and Engineering. Springer International Publishing, 2016.
- [24] T. Belytschko, W. K. Liu, B. Moran, and K. I. Elkhodary. *Nonlinear Finite Elements for Continua and Structures, 2nd Edition*. Wiley, 2014.
- [25] F. Bertrand, C. Carstensen, B. Gräßle, and N. T. Tran. Stabilization-free HHO a posteriori error control. *arXiv.org*, 2207.01038, 2022. Preprint.
- [26] D. Boffi, M. Botti, and D. A. Di Pietro. A nonconforming high-order method for the Biot problem on general meshes. *SIAM J. Sci. Comput.*, 38(3):A1508–A1537, 2016.
- [27] F. Bonaldi, D. A. Di Pietro, G. Geymonat, and F. Krasucki. A hybrid high-order method for Kirchhoff-Love plate bending problems. *ESAIM Math. Model. Numer. Anal.*, 52(2):393–421, 2018.
- [28] J. Bonet and A. Burton. A simple average nodal pressure tetrahedral element for incompressible and nearly incompressible dynamic explicit applications. *Commun. Numer. Methods Eng.*, 14(5):437–449, 1998.
- [29] S. Bonetti, M. Botti, I. Mazzieri, and P. F. Antonietti. Numerical modelling of wave propagation phenomena in thermo-poroelastic media via discontinuous Galerkin methods. *J. Comput. Phys.*, page 112275, 2023.
- [30] L. Botti, D. A. Di Pietro, and J. Droniou. A hybrid high-order method for the incompressible Navier-Stokes equations based on Temam’s device. *J. Comput. Phys.*, 376:786–816, 2019.
- [31] M. Botti, D. A. Di Pietro, and P. Sochala. A hybrid high-order method for nonlinear elasticity. *SIAM J. Numer. Anal.*, 55(6):2687–2717, 2017.
- [32] E. Burman, M. Cicuttin, G. Delay, and A. Ern. An unfitted hybrid high-order method with cell agglomeration for elliptic interface problems. *SIAM J. Sci. Comput.*, 43(2):A859, 2021.

- [33] E. Burman, G. Delay, and A. Ern. An unfitted hybrid high-order method for the Stokes interface problem. *IMA J. Numer. Anal.*, 41(4):2362–2387, 2021.
- [34] E. Burman, O. Duran, and A. Ern. Hybrid high-order methods for the acoustic wave equation in the time domain. *Commun. Appl. Math. Comput. (CAMC)*, 4(2):597–633, 2022.
- [35] E. Burman, O. Duran, and A. Ern. Unfitted hybrid high-order methods for the wave equation. *Comp. Meth. Appl. Mech. Eng.*, 389:114366, 2022.
- [36] E. Burman, O. Duran, A. Ern, and M. Steins. Convergence analysis of hybrid high-order methods for the wave equation. *J. Sci. Comput.*, 87(3):91, 2021.
- [37] E. Burman and A. Ern. An unfitted hybrid high-order method for elliptic interface problems. *SIAM J. Numer. Anal.*, 56:1525–1546, 2018.
- [38] V. Calo, M. Cicuttin, Q. Deng, and A. Ern. Spectral approximation of elliptic operators by the hybrid high-order method. *Math. Comp.*, 88(318):1559–1586, 2019.
- [39] A. Cangiani, Z. Dong, E. H. Georgoulis, and P. Houston. *hp-version discontinuous Galerkin methods on polygonal and polyhedral meshes*. SpringerBriefs in Mathematics. Springer, Cham, 2017.
- [40] C. Carle, M. Hochbruck, and A. Sturm. On leapfrog-Chebyshev schemes. *SIAM J. Numer. Anal.*, 58(4):2404–2433, 2020.
- [41] C. Carstensen, A. Ern, and S. Puttkammer. Guaranteed lower bounds on eigenvalues of elliptic operators with a hybrid high-order method. *Numer. Math.*, 149(2):273–304, 2021.
- [42] C. Carstensen and N. T. Tran. Convergent adaptive hybrid higher-order schemes for convex minimization. *Numer. Math.*, 151(2):329–367, 2022.
- [43] C. Carstensen and T. Tran. Unstabilized hybrid high-order method for a class of degenerate convex minimization problems. *SIAM J. Numer. Anal.*, 59(3):1348–1373, 2021.
- [44] K. L. Cascavita, J. Bleyer, X. Chateau, and A. Ern. Hybrid discretization methods with adaptive yield surface detection for Bingham pipe flows. *J. Sci. Comput.*, 77(3):1424–1443, 2018.
- [45] K. L. Cascavita, F. Chouly, and A. Ern. Hybrid high-order discretizations combined with Nitsche’s method for Dirichlet and Signorini boundary conditions. *IMA J. Numer. Anal.*, 40(4):2189–2226, 2020.
- [46] D. Castanon Quiroz, D. A. Di Pietro, and A. Harnist. A hybrid high-order method for incompressible flows of non-Newtonian fluids with power-like convective behaviour. *IMA J. Numer. Anal.*, 43(1):144–186, 2023.
- [47] M. Cervera, M. Chiumenti, L. Benedetti, and R. Codina. Mixed stabilized finite element methods in nonlinear solid mechanics. part iii: Compressible and incompressible plasticity. *Comput. Methods Appl. Mech. Eng.*, 285:752–775, 2015.
- [48] M. Cervera, M. Chiumenti, and R. Codina. Mixed stabilized finite element methods in nonlinear solid mechanics: Part i: Formulation. *Comput. Methods Appl. Mech. Eng.*, 199(37-40):2559–2570, 2010.
- [49] M. Cervera, M. Chiumenti, Q. Valverde, and C. Agelet de Saracibar. Mixed linear/linear simplicial elements for incompressible elasticity and plasticity. *Comput. Methods Appl. Mech. Eng.*, 192(49-50):5249–5263, 2003.
- [50] M. Cervera, N. Lafontaine, R. Rossi, and M. Chiumenti. Explicit mixed strain-displacement finite elements for compressible and quasi-incompressible elasticity and plasticity. *Comput. Mech.*, 58(3):511–532, 2016.

- [51] J. Chabassier and S. Imperiale. Introduction and study of fourth order theta schemes for linear wave equations. *J. Comput. Appl. Math.*, 245:194–212, 2013.
- [52] J. Chabassier and S. Imperiale. Space/time convergence analysis of a class of conservative schemes for linear wave equations. *CRAS*, 355(3):282–289, 2017.
- [53] J. Chabassier and P. Joly. Energy preserving schemes for nonlinear Hamiltonian systems of wave equations: Application to the vibrating piano string. *Comp. Meth. Appl. Mech. Eng.*, 199(45):2779–2795, 2010.
- [54] T. Chaumont-Frelet, A. Ern, S. Lemaire, and F. Valentin. Bridging the multiscale hybrid-mixed and multiscale hybrid high-order methods. *ESAIM Math. Model. Numer. Anal.*, 56(1):261, 2022.
- [55] F. Chave, D. A. Di Pietro, and L. Formaggia. A hybrid high-order method for Darcy flows in fractured porous media. *SIAM J. Sci. Comput.*, 40(2):1063–1094, 2018.
- [56] F. Chave, D. A. Di Pietro, and L. Formaggia. A hybrid high-order method for passive transport in fractured porous media. *Int. J. Geomath.*, 10(12):1–34, 2019.
- [57] F. Chave, D. A. Di Pietro, F. Marche, and F. Pigeonneau. A hybrid high-order method for the Cahn-Hilliard problem in mixed form. *SIAM J. Numer. Anal.*, 54(3):1873–1898, 2016.
- [58] P. Chellapandi, K. Velusamy, S. C. Chetal, S. B. Bhoje, H. Lal, and V. S. Sethi. Analysis for mechanical consequences of a core disruptive accident in prototype fast breeder reactor. In *Transactions of the 17th International Conference on Structural Mechanics in Reactor Technology (SMiRT 17)*, 2003.
- [59] M. Chiumenti, M. Cervera, and R. Codina. A mixed three-field FE formulation for stress accurate analysis including the incompressible limit. *Comput. Methods Appl. Mech. Eng.*, 283:1095–1116, 2015.
- [60] M. Chiumenti, Q. Valverde, C. Agelet de Saracibar, and M. Cervera. A stabilized formulation for incompressible plasticity using linear triangles and tetrahedra. *Int. J. Plast.*, 20(8-9):1487–1504, 2004.
- [61] F. Chouly, A. Ern, and N. Pignet. A hybrid high-order discretization combined with Nitsche’s method for contact and Tresca friction in small strain elasticity. *SIAM J. Sci. Comput.*, 42(4):2300–2324, 2020.
- [62] P. G. Ciarlet. *The Finite Element Method for Elliptic Problems*. North-Holland, Amsterdam, 1978.
- [63] M. Cicuttin, D. A. Di Pietro, and A. Ern. Implementation of discontinuous skeletal methods. *J. Comput. Appl. Math.*, 344:852–874, 2018.
- [64] M. Cicuttin, A. Ern, and T. Gudi. Hybrid high-order methods for the elliptic obstacle problem. *J. Sci. Comput.*, 83(1):8, 2020.
- [65] M. Cicuttin, A. Ern, and S. Lemaire. A hybrid high-order method for highly oscillatory elliptic problems. *Comput. Methods Appl. Math.*, 19(4):723–748, 2019.
- [66] M. Cicuttin, A. Ern, and N. Pignet. *Hybrid High-Order Methods. A Primer with Application to Solid Mechanics*. Springer, 2021. SpringerBriefs in Mathematics.
- [67] B. Cockburn. Discontinuous Galerkin methods. *ZAMM - J. Appl. Math. Mech*, 83(11):731–754, 2003.
- [68] B. Cockburn. Static condensation, hybridization, and the devising of the HDG methods. In *Building bridges: connections and challenges in modern approaches to numerical partial differential equations*, pages 129–177. Springer, 2016.



- [69] B. Cockburn, Z. Fu, A. Hungria, L. Ji, M. A. Sánchez, and F.-J. Sayas. Störmer-Numerov HDG methods for acoustic waves. *J. Sci. Comput.*, 75(2):597–624, 2018.
- [70] B. Cockburn, J. Gopalakrishnan, and R. Lazarov. Unified hybridization of discontinuous Galerkin, mixed, and continuous Galerkin methods for second order elliptic problems. *SIAM J. Numer. Anal.*, 47(2):1319–1365, 2009.
- [71] B. Cockburn, D. A. Di Pietro, and A. Ern. Bridging the hybrid high-order and hybridizable discontinuous Galerkin methods. *ESAIM Math. Mod. Numer. Anal.*, 50(3):635–650, 2016.
- [72] B. Cockburn and V. Quenneville-Bélair. Uniform-in-time superconvergence of the HDG methods for the acoustic wave equation. *Math. Comp.*, 83(285):65–85, 2014.
- [73] G. Cohen, P. Joly, J. E. Roberts, and N. Tordjman. Higher order triangular finite elements with mass lumping for the wave equation. *SIAM J. Numer. Anal.*, 38(6):2047–2078, 2001.
- [74] G. Cohen, P. Joly, and N. Tordjman. Higher-order finite elements with mass-lumping for the 1d wave equation. *Finite Elem. Anal. Des.*, 16(3-4):329–336, 1994.
- [75] P. Daniel, A. Ern, I. Smears, and M. Vohralík. An adaptive hp-refinement strategy with computable guaranteed bound on the error reduction factor. *Comput. Math. Appl.*, 76(5):967–983, 2018.
- [76] D. A. Di Pietro and J. Droniou. A hybrid high-order method for Leray-Lions elliptic equations on general meshes. *Math. Comp.*, 86(307):2159–2191, 2017.
- [77] 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 discretisation of Leray-Lions problems. *Math. Models Methods Appl. Sci.*, 27(5):879–908, 2017.
- [78] D. A. Di Pietro and J. Droniou. *The Hybrid High-Order Method for Polytopal Meshes*. Modeling, Simulation and Applications series. Springer International Publishing, 2020.
- [79] D. A. Di Pietro, J. Droniou, and A. Ern. A discontinuous-skeletal method for advection-diffusion-reaction on general meshes. *SIAM J. Numer. Anal.*, 53(5):2135–2157, 2015.
- [80] D. A. Di Pietro, J. Droniou, and G. Manzini. Discontinuous skeletal gradient discretisation methods on polytopal meshes. *J. Comput. Phys.*, 355:397–425, 2018.
- [81] D. A. Di Pietro and A. Ern. *Mathematical Aspects of Discontinuous Galerkin Methods*, volume 69. Springer, 2012.
- [82] D. A. Di Pietro and A. Ern. A hybrid high-order locking-free method for linear elasticity on general meshes. *Comp. Meth. Appl. Mech. Eng.*, 283:1–21, 2015.
- [83] D. A. Di Pietro and A. Ern. Hybrid high-order methods for variable-diffusion problems on general meshes. *C. R. Math. Acad. Sci. Paris*, 353(1):31–34, 2015.
- [84] 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. *Comput. Methods Appl. Math.*, pages 461–472, 2014.
- [85] D. A. Di Pietro, A. Ern, and S. Lemaire. A review of hybrid high-order methods: Formulations, computational aspects, comparison with other methods. In *Building Bridges: Connections and Challenges in Modern Approaches to Numerical Partial Differential Equations*, vol 114 of Lect. Notes Comput. Sci. Eng., pages 205–236. Springer, 2016.
- [86] D. A. Di Pietro, A. Ern, A. Linke, and F. Schieweck. A discontinuous skeletal method for the viscosity-dependent Stokes problem. *Comput. Methods Appl. Mech. Engrg.*, 306:175–195, 2016.

- [87] D. A. Di Pietro and S. Krell. A hybrid high-order method for the steady incompressible Navier-Stokes problem. *J. Sci. Comput.*, 74(3):1677–1705, 2018.
- [88] D. A. Di Pietro and R. Tittarelli. *An Introduction to Hybrid High-Order Methods*, pages 75–128. volume 15 of SEMA SIMAI Springer Ser. Springer, 2018.
- [89] J. Diaz and A. Ezziani. Garmore2D software. <http://www.spice-rtn.org/library/software/Gar6more2D/>, 2008.
- [90] J. Diaz and M. J. Grote. Multi-level explicit local time-stepping methods for second-order wave equations. *Comput. Methods Appl. Mech. Engrg.*, 291:240–265, 2015.
- [91] W. P. Doherty, E. L. Wilson, and R. L. Taylor. Stress analysis of axisymmetric solids utilizing higher order quadrilateral finite elements, SESM Report 69-3, Department of Civil Engineering, *University of California, Berkeley*, 1969.
- [92] Z. Dong and A. Ern. Hybrid high-order method for singularly perturbed fourth-order problems on curved domains. *ESAIM Math. Model. Numer. Anal.*, 55(66):3091–3114, 2021.
- [93] Z. Dong and A. Ern. Hybrid high-order and weak Galerkin methods for the biharmonic problem. *SIAM J. Numer. Anal.*, 60(5):2626–2656, 2022.
- [94] T. Dupont.  $l^2$ -estimates for Galerkin methods for second order hyperbolic equations. *SIAM J. Numer. Anal.*, 10(5):880–889, 1973.
- [95] T. Elguedj, Y. Bazilevs, V. M. Calo, and T. J. R. Hughes.  $\bar{B}$  and  $\bar{F}$  projection methods for nearly incompressible linear and non-linear elasticity and plasticity using higher-order nurbs elements. *Comput. Methods Appl. Mech. Eng.*, 197(33):2732–2762, 2008.
- [96] T. Elguedj and T. J. R. Hughes. Isogeometric analysis of nearly incompressible large strain plasticity. *Comput. Methods Appl. Mech. Eng.*, 268:388–416, 2014.
- [97] A. Ern and J.-L. Guermond. *Finite Elements I: Approximation and Interpolation*, volume 72 of *Texts in Applied Mathematics*. Springer, Cham, 2021.
- [98] A. Ern and J.-L. Guermond. *Finite elements. II. Galerkin approximation, elliptic and mixed PDEs*, volume 73 of *Texts in Applied Mathematics*. Springer, Cham, 2021.
- [99] A. Ern and J.-L. Guermond. Quasi-optimal nonconforming approximation of elliptic PDEs with contrasted coefficients and  $H^{1+r}$ ,  $r > 0$ , regularity. *Found. Comput. Math.*, 22(5):1273–1308, 2022.
- [100] A. Ern, F. Hédin, G. Pichot, and N. Pignet. Hybrid high-order methods for flow simulations in extremely large discrete fracture networks. *SMAI J. Comput. Math.*, 8:375–398, 2022.
- [101] A. Ern and M. Steins. Convergence analysis for the wave equation discretized with hybrid methods in space (HHO, HDG and WG) and the leapfrog scheme in time. *submitted*, 2023.
- [102] V. Faucher and T. Gautier. Simulation de l’accident de dimensionnement du confinement dans un réacteur de 1<sup>re</sup> génération avec une approche parallèle hybride dans Europlexus. In *11e colloque national en calcul des structures*, 2013.
- [103] S. Geevers, W. A. Mulder, and J. J. W. van der Vegt. New higher-order mass-lumped tetrahedral elements for wave propagation modelling. *SIAM J. Sci. Comput.*, 40(5):A2830–A2857, 2018.
- [104] F. X. Giraldo and M. A. Taylor. A diagonal-mass-matrix triangular-spectral-element method based on cubature points. *J. Engrg. Math.*, 56(3):307–322, 2006.
- [105] S. Glaser and F. Armero. On the formulation of enhanced strain finite elements in finite deformations. *Eng. Comput.*, 14(7):759–791, 1997.

- [106] M. J. Grote, S. Michel, and S. A. Sauter. Stabilized leapfrog based local time-stepping method for the wave equation. *Math. Comp.*, 90(332):2603–2643, 2021.
- [107] M. J. Grote, A. Schneebeli, and D. Schötzau. Discontinuous Galerkin finite element method for the wave equation. *SIAM J. Numer. Anal.*, 44(6):2408–2431, 2006.
- [108] B. Halphen and Q. S. Nguyen. Sur les matériaux standard généralisés. *Journal de mécanique*, 14(1):39–63, 1975.
- [109] P. Hansbo and M. G. Larson. Discontinuous Galerkin methods for incompressible and nearly incompressible elasticity by Nitsche’s method. *Comput. Methods Appl. Mech. Eng.*, 191(17-18):1895–1908, 2002.
- [110] C. Harder, D. Paredes, and F. Valentin. A family of multiscale hybrid-mixed finite element methods for the darcy equation with rough coefficients. *J. Comput. Phys.*, 245:107–130, 2013.
- [111] Thomas. Helfer, J.-M. Proix, and O. Fandeur. Implantation de lois de comportement mécanique à l’aide de MFront: simplicité, efficacité, robustesse et portabilité. In *12e Colloque national en calcul des structures*. CSMA, 2015.
- [112] Y. Huang, J. Li, and D. Li. Developing weak Galerkin finite element methods for the wave equation. *Numer. Methods Partial Differ. Equations*, 33(3):868–884, 2017.
- [113] B. Hudobivnik, F. Aldakheel, and P. Wriggers. A low order 3d virtual element formulation for finite elasto-plastic deformations. *Comput. Mech.*, 63:253–269, 2019.
- [114] T. J. R. Hughes. Generalization of selective integration procedures to anisotropic and nonlinear media. *Int. J. Numer. Methods Eng.*, 15(9):1413–1418, 1980.
- [115] T. J. R. Hughes. *The Finite Element Method*. Linear Static and Dynamic Finite Element Analysis. Prentice-Hall, Englewood Cliffs, New Jersey, 1987.
- [116] P. Jana, N. Kumar, and B. Deka. A systematic study on weak Galerkin finite-element method for second-order wave equation. *Comput. Appl. Math.*, 41(8):Paper No. 359, 25, 2022.
- [117] P. Joly. Variational methods for time-dependent wave propagation problems. In *Topics in computational wave propagation*, volume 31 of *Lect. Notes Comput. Sci. Eng.*, pages 201–264. Springer, Berlin, 2003.
- [118] E. P. Kasper and R. L. Taylor. A mixed-enhanced strain method: Part i: Geometrically linear problems. *Comput. Struct.*, 75(3):237–250, 2000.
- [119] E. P. Kasper and R. L. Taylor. A mixed-enhanced strain method: Part ii: Geometrically nonlinear problems. *Comput. Struct.*, 75(3):251–260, 2000.
- [120] M. Kronbichler, S. Schoeder, C. Müller, and W. A. Wall. Comparison of implicit and explicit hybridizable discontinuous Galerkin methods for the acoustic wave equation. *Int. J. Numer. Methods Eng.*, 106(9):712–739, 2016.
- [121] P. Krysl. Mean-strain eight-node hexahedron with optimized energy-sampling stabilization for large-strain deformation. *Int. J. Numer. Methods Eng.*, 103(9):650–670, 2015.
- [122] P. Krysl. Mean-strain 8-node hexahedron with optimized energy-sampling stabilization. *Finite Elem. Anal. Des.*, 108:41–53, 2016.
- [123] P. Le Tallec, C. Rahier, and A. Kaiss. Three-dimensional incompressible viscoelasticity in large strains: formulation and numerical approximation. *Comput. Methods Appl. Mech. Eng.*, 109(3-4):233–258, 1993.

- [124] C. Lehrenfeld. *Hybrid Discontinuous Galerkin methods for solving incompressible flow problems*. PhD thesis, Rheinisch-Westfälischen Technischen Hochschule (RWTH), Aachen, 2010.
- [125] C. Lehrenfeld and J. Schöberl. High order exactly divergence-free hybrid discontinuous Galerkin methods for unsteady incompressible flows. *Comp. Meth. Appl. Mech. Eng.*, 307:339–361, 2016.
- [126] J. Lemaitre and J.-L. Chaboche. *Mechanics of solid materials*. Cambridge university press, 1994.
- [127] R. Liu, M. F. Wheeler, C. N. Dawson, and R. H. Dean. A fast convergent rate preserving discontinuous Galerkin framework for rate-independent plasticity problems. *Comput. Methods Appl. Mech. Eng.*, 199(49-52):3213–3226, 2010.
- [128] R. Liu, M. F. Wheeler, and I. Yotov. On the spatial formulation of discontinuous Galerkin methods for finite elastoplasticity. *Comput. Methods Appl. Mech. Eng.*, 253:219–236, 2013.
- [129] D. S. Malkus and T. J. R. Hughes. Mixed finite element methods — reduced and selective integration techniques: A unification of concepts. *Comput. Methods Appl. Mech. Eng.*, 15(1):63–81, 1978.
- [130] C. Miehe, N. Apel, and M. Lambrecht. Anisotropic additive plasticity in the logarithmic strain space: modular kinematic formulation and implementation based on incremental minimization principles for standard materials. *Comput. Methods Appl. Mech. Eng.*, 191(47):5383–5425, 2002.
- [131] Christian. Miehe, N. Apel, and M. Lambrecht. Anisotropic additive plasticity in the logarithmic strain space: modular kinematic formulation and implementation based on incremental minimization principles for standard materials. *Comput. Methods Appl. Mech. Eng.*, 191(47-48):5383–5425, 2002.
- [132] N. C. Nguyen and J. Peraire. Hybridizable discontinuous Galerkin methods for partial differential equations in continuum mechanics. *J. Comput. Phys.*, 231(18):5955–5988, 2012.
- [133] N. C. Nguyen, J. Peraire, and B. Cockburn. High-order implicit hybridizable discontinuous Galerkin methods for acoustics and elastodynamics. *J. Comput. Phys.*, 230, 2011.
- [134] I. Ramière and T. Helfer. Iterative residual-based vector methods to accelerate fixed point iterations. *Comput. Math. Appl.*, 70(9):2210–2226, 2015.
- [135] S. Reese. On a consistent hourglass stabilization technique to treat large inelastic deformations and thermo-mechanical coupling in plane strain problems. *Int. J. Numer. Methods Eng.*, 57(8):1095–1127, 2003.
- [136] M. A. Sánchez, C. Ciuca, N. C. Nguyen, J. Peraire, and B. Cockburn. Symplectic Hamiltonian HDG methods for wave propagation phenomena. *J. Comput. Phys.*, 350:951–973, 2017.
- [137] L. R. Scott and M. Vogelius. Norm estimates for a maximal right inverse of the divergence operator in spaces of piecewise polynomials. *ESAIM Math. Model. Numer. Anal.*, 19(1):111–143, 1985.
- [138] J. C. Simo. Algorithms for static and dynamic multiplicative plasticity that preserve the classical return mapping schemes of the infinitesimal theory. *Comput. Methods Appl. Mech. Eng.*, 99(1):61–112, 1992.
- [139] J. C. Simo and F. Armero. Geometrically non-linear enhanced strain mixed methods and the method of incompatible modes. *Int. J. Numer. Methods Eng.*, 33(7):1413–1449, 1992.
- [140] J. C. Simo, F. Armero, and R. L. Taylor. Improved versions of assumed enhanced strain tri-linear elements for 3d finite deformation problems. *Comput. Methods Appl. Mech. Eng.*, 110(3-4):359–386, 1993.
- [141] J. C. Simo and M. S. Rifai. A class of mixed assumed strain methods and the method of incompatible modes. *Int. J. Numer. Methods Eng.*, 29(8):1595–1638, 1990.

- [142] J. C. Simo, R. L. Taylor, and K. Pister. Variational and projection methods for the volume constraint in finite deformation elasto-plasticity. *Comput. Methods Appl. Mech. Eng.*, 51(1-3):177–208, 1985.
- [143] M. Stanglmeier, N. C. Nguyen, J. Peraire, and B. Cockburn. An explicit hybridizable discontinuous Galerkin method for the acoustic wave equation. *Comp. Meth. Appl. Mech. Eng.*, 300:748–769, 2016.
- [144] M. Steins, A. Ern, O. Jamond, and F. Drui. Time-explicit hybrid high-order method for the nonlinear acoustic wave equation. *ESAIM Math. Mod. Numer. Anal.*, 57(5):2977–3006, 2023.
- [145] R. Stevenson. Optimality of a standard adaptive finite element method. *Found. Comput. Math.*, 7:245–269, 2007.
- [146] J. S. Sun, K. H. Lee, and H. P. Lee. Comparison of implicit and explicit finite element methods for dynamic problems. *J. Mater. Process. Technol.*, 105(1):110–118, 2000.
- [147] R. L. Taylor. A mixed-enhanced formulation tetrahedral finite elements. *Int. J. Numer. Methods Eng.*, 47(1-3):205–227, 2000.
- [148] M. Vogelius. A right-inverse for the divergence operator in spaces of piecewise polynomials: Application to the p-version of the finite element method. *Numer. Math.*, 41(1):19–37, 1983.
- [149] J. Wang and X. Ye. A weak Galerkin finite element method for second-order elliptic problems. *J. Comput. Appl. Math.*, 241:103–115, 2013.
- [150] M. F. Wheeler. A priori  $L_2$  error estimates for Galerkin approximations to parabolic partial differential equations. *SIAM J. Numer. Anal.*, 10:723–759, 1973.
- [151] P. Wriggers and B. Hudobivnik. A low order virtual element formulation for finite elasto-plastic deformations. *Comput. Methods Appl. Mech. Eng.*, 327:459–477, 2017.
- [152] O. C. Zienkiewicz, R. L. Taylor, and J. M. Too. Reduced integration technique in general analysis of plates and shells. *Int. J. Numer. Methods Eng.*, 3(2):275–290, 1971.