

# Shape and topology optimization of multiphysics systems

Thèse de doctorat de l'Université Paris-Saclay  
préparée à l'École polytechnique

Ecole doctorale n°574 Ecole doctorale de mathématiques Hadamard (EDMH)  
Spécialité de doctorat : Mathématiques appliquées

Thèse présentée et soutenue à l'École polytechnique, le 16 décembre 2019, par

**FLORIAN FEPPON**

Composition du Jury :

Grégoire ALLAIRE Professeur, École polytechnique (CMAP)	Directeur de thèse
Samuel AMSTUTZ Professeur, École polytechnique (CMAP)	Examineur
Julien CORTIAL Ingénieur de recherche, Safran	Invité
Charles DAPOGNY Chargé de recherche, Université Grenoble Alpes (LJK)	Co-directeur de thèse
Sonia FLISS Enseignant chercheur, ENSTA (UMA)	Présidente du jury
Antoine HENROT Professeur, École des mines de Nancy	Rapporteur
Robert KOHN Professeur émérite, New York University (Courant Institute)	Rapporteur (absent)
Olivier PIRONNEAU Professeur émérite, Université Pierre et Marie Curie (LJLL)	Examineur
Yannick PRIVAT Professeur, Université de Strasbourg	Examineur
Nicole SPILLANE Chargé de recherche, École polytechnique (CMAP)	Examinatrice



# THÈSE DE DOCTORAT

de

L'UNIVERSITÉ PARIS-SACLAY

École doctorale de mathématiques Hadamard (EDMH, ED 574)

*Établissement d'inscription* : École polytechnique

*Laboratoire d'accueil* : Centre de mathématiques appliquées de Polytechnique, UMR 7641 CNRS

*Spécialité de doctorat* : Mathématiques appliquées

**Florian FEPPON**

Shape and topology optimization of multiphysics systems

*Date de soutenance* : 16 décembre 2019

*Après avis des rapporteurs* : ANTOINE HENROT (UNIVERSITÉ DE LORRAINE)  
ROBERT KOHN (COURANT INSTITUTE, NYU)

*Jury de soutenance* :

GRÉGOIRE ALLAIRE	(CMAP) Directeur de thèse
SAMUEL AMSTUTZ	(CMAP) Examineur
JULIEN CORTIAL	(Safran) Invité
CHARLES DAPOGNY	(LJK) Co-encadrant de thèse
SONIA FLISS	(ENSTA – UMA) Examinatrice
ANTOINE HENROT	(École des Mines de Nancy) Rapporteur
OLIVIER PIRONNEAU	(LJLL) Examineur
YANNICK PRIVAT	(Université de Strasbourg) Examineur
NICOLE SPILLANE	(CMAP) Examinatrice



## ABSTRACT

This work is devoted to shape and topology optimization of multiphysics systems motivated by aeronautic industrial applications. Shape derivatives of arbitrary objective functionals are computed for a weakly coupled thermal fluid-structure model. A novel gradient flow type algorithm is then developed for solving generic constrained shape optimization problems without the need for tuning non-physical metaparameters. Motivated by the need for enforcing non-mixing constraints in the design of liquid-liquid heat exchangers, a variational method is developed in order to simplify the numerical evaluation of geometric constraints: it allows to compute line integrals on a mesh by solving a variational problem without requiring the explicit knowledge of these lines on the spatial discretization. All these ingredients allowed us to implement a variety of 2-d and 3-d multiphysics shape optimization test cases: from single, double or three physics problems in 2-d, to moderately large-scale 3-d test cases for structural design, thermal conduction, aerodynamic design and a fluid-structure interacting system. A final opening chapter derives high order homogenized equations for the Stokes system in a porous medium. These high order equations encompass the three classical homogenized regimes—namely Stokes, Brinkman and Darcy—associated with different obstacle's size scalings. They could allow, in future works, to develop new topology optimization methods for the design of fluid systems characterized by multi-scale patterns such as industrial heat exchangers.

**Keywords:** Topology optimization, remeshing, convective heat transfer, fluid-structure interaction, geometric constraints, high order homogenization.

## RÉSUMÉ

Cette thèse est consacrée à l'optimisation de la topologie et de la forme de systèmes multiphysiques motivés par des applications de l'industrie aéronautique. Nous calculons les dérivées de forme de fonctions de coût arbitraires pour un modèle fluide, thermique et mécanique faiblement couplé. Nous développons ensuite un algorithme de type gradient adapté à la résolution de problèmes d'optimisation de formes sous contraintes qui ne requiert pas de réglage de paramètres non physiques. Nous introduisons ensuite une méthode variationnelle qui permet de calculer des intégrales le long de rayons sur un maillage par la résolution d'un problème variationnel qui ne requiert pas la détermination explicite de ces lignes sur la discrétisation spatiale. Cette technique nous a ainsi permis d'imposer une contrainte de non-mélange de phases pour une application à l'optimisation d'échangeurs de chaleur bi-tubes. Tous ces ingrédients ont été employés pour traiter une variété de cas tests d'optimisation de formes pour des systèmes multiphysiques 2-d ou 3-d. Nous avons considéré des problèmes à une seule, deux ou bien trois physiques couplées en 2-d, et des problèmes de tailles relativement élevées en 3-d pour la mécanique, la conduction thermique, l'optimisation de profils aérodynamiques, et de la forme de systèmes en interaction fluide-structure. Un dernier chapitre d'ouverture est consacré à l'étude de modèles homogénéisés d'ordres élevés pour les équations de Stokes en milieu poreux. Ces équations d'ordres élevés englobent les trois régimes homogénéisés classiques—Stokes, Brinkman et Darcy—associés à divers rapports d'échelles pour la taille des obstacles. Elles pourraient permettre, lors de futurs travaux, de développer de nouvelles méthodes d'optimisation pour la conception de systèmes fluides caractérisés par des motifs multiéchelles, tels que les échangeurs thermiques industriels.

**Mots clés :** Optimisation topologique, remaillage, transfert thermique convectif, interaction fluide-structure, contraintes géométriques, modèles homogénéisés d'ordres élevés.



## REMERCIEMENTS

ॐ भूर्भुवः स्वः    *om bhūr bhuvah suvah*  
तत्सवितुर्वरेण्यं    *tatsaviturvareṇyam*  
भर्गो देवस्य धीमहि    *bhargo devasyadhīmahi*  
धियो यो नः प्रचोदयात् ॥    *dhiyo yo naḥ prachodayāt*

*We choose the Supreme Light of the divine Sun; we aspire that it may impel our minds.  
(Gayatri Mantra, Trad. Sri Aurobindo)*

Mes remerciements vont en premier lieu à la providence qui a bien voulu créer les conditions propices à la réalisation de ce travail de par la réunion de bonnes volontés en un même endroit et en un même moment opportun. Certes, ces dispositions peuvent être comparées au fruit arrivé à maturation grâce à de nombreuses années de labeur; il suffit alors d'avoir la chance d'arriver au bon moment pour pouvoir le cueillir et en extraire la saveur. Cependant, malgré toutes ces conditions *nécessaires*, il s'en est fallu de peu pour que ce travail ne voie pas le jour sous sa forme actuelle. La recherche scientifique est une *écoute* alerte et vigilante permanente : il suffit d'une innocente autocensure de *rien du tout*, d'un *menu* attachement à une idée un peu trop personnelle, d'une *imperceptible* procrastination, d'une *anodine* discussion en moins avec un camarade doctorant éclairé pour que le fruit reste invisible au milieu du brouillard de l'opacité mentale. Je remercie donc cette grâce intérieure qui nourrit notre intuition (pour ne pas dire, nous souffle les preuves des théorèmes), nous guide vers des rencontres fructueuses, nous fait adorer la clarté, et nous réunit joyeusement à porter nos efforts vers une même direction.

Il convient ensuite de témoigner ma gratitude à l'ensemble des personnes qui m'ont accompagné et porté sur le chemin de cette thèse. Je remercie Grégoire pour avoir déroulé un long tapis rouge jonché de belles opportunités devant mes pieds depuis la fin 2014. Merci pour le dévouement, l'humanité, et la bienveillance avec laquelle tu encadres tes thésards. Un grand merci à Charles pour la sincérité de ton encadrement, pour l'abondante assistance que tu as pu m'apporter durant ces trois ans, pour ton accueil à Grenoble, pour l'effort de tes relectures et commentaires attentifs, et enfin plus encore pour les précieuses discussions et conseils à un thésard s'interrogeant sur sa vie. Je suis certain que la fin de cette thèse n'est qu'un début vers de nouvelles aventures communes, et si possible bien évidemment près des montagnes...

Je tiens à remercier sincèrement l'encadrement dont j'ai bénéficié à Safran Tech : Julien, pour son attachement caractéristique à la mise en place et au maintien de toutes les *conditions d'optimalité* nécessaires au bon déroulement de la thèse, puis Felipe et Christian pour leurs encouragements soutenus tout au long du travail<sup>1</sup>. Bien sûr, la vie à Safran Tech aurait été moins agréable sans la présence des permanents : je remercie entre autres Oana, Nissrine, Sophie, Noémie, Tonya, Wafa, Emma, Sébastien, Augustin, Fabien, Franck, Mohammed, Guy, Nicolas et Frédéric. L'environnement Safran Tech, à mi-chemin entre l'industrie et le monde académique, m'a paru fort enrichissant, aussi je recommanderais sans hésiter cette expérience CIFRE à de futurs doctorants.

Merci ensuite aux personnes du monde académique qui se sont intéressées de près à ce travail : notamment Yannick Privat pour son invitation à Strasbourg et ses encouragements, et Georgios pour avoir continué à dédier un peu de ton temps et de ton expérience pour nous. Je suis également reconnaissant envers Frédéric Hecht et Olivier Pironneau pour les échanges autour de mes travaux et les invitations au Laboratoire Jacques Louis Lions (LJLL). Je remercie chaleureusement Algiane, Pierre Jolivet et Vincent Moureau pour l'aide précieuse qu'ils m'ont apportée pour la réalisation des cas tests tridimensionnels.

Un grand merci à Antoine Henrot et Robert Kohn qui m'ont fait l'honneur de lire et rapporter mon travail. Je remercie de même tous les membres de mon jury de thèse qui ont accepté sans hésiter ce rôle et qui ont bien voulu offrir leur présence et leur attention le 16 décembre 2019.

Du côté du Centre de Mathématiques Appliquées (CMAP), je suis bien sûr reconnaissant envers Nasséra, Alex et Maud pour leur accueil, leur dévouement et leur patience infinies envers les doctorants trop peu soucieux des règles administratives (occupés qu'ils sont à faire des mathématiques, il leur faut environ trois ans pour assimiler toutes les règles...). Merci ensuite à toutes les autres personnes avec lesquelles j'ai eu l'opportunité d'échanger : entre autres Nicole Spillane, François Alouges, Samuel Amstutz, Lucas Chennel. Merci à Pierre et Sylvain pour leur disponibilité à toute épreuve en ce qui concerne les demandes informatiques. Merci à tous les doctorants du laboratoire avec lesquels j'ai eu

<sup>1</sup>Notons que la thèse, *in fine*, est aussi un travail soutenu.

le plaisir de partager ces trois années (en vous adressant mes plus sincères excuses si par hasard votre nom n'est pas cité). Merci tout particulièrement à Léa au moins pour une balade ensoleillée en vélo en Bretagne, Martin A. pour les moments de partages autour de nos passions artistiques, à Kevish pour son enthousiasme et sa simplicité d'être et à qui j'espère ne pas avoir (trop) abimé les genoux. Merci à Céline, Pamela, Juliette, Geneviève, Pierre C., Florian B., Fedor, Heythem, Nicolas, Corentin H. et Corentin C., Frédéric, toujours enthousiastes pour se retrouver, discuter, ou bien mettre de la vie dans les couloirs du laboratoire.

Merci à toutes les personnes du bureau 2016 qui ont subi ma présence : plus particulièrement Jean-Bernard, illustre maître du bureau, pour le partage de tes réflexions avisées sur l'actualité politique, Hadrien en tant qu'exemple vivant que la réussite peut naître de la désorganisation, Raphaël et Cédric. Aux personnes qui ont quitté ce bureau moins tôt : Mathieu toujours très réceptif à mon humour et qui nous a enseigné un moyen très astucieux pour maintenir en place une fenêtre de bureau récalcitrante. Merci à Julie, partante pour toutes les sorties sportives et qui parvient à nous convaincre (non sans usage d'une certaine insistance...) de venir l'accompagner par nuit noire courir sur des chemins boueux de forêt. Merci à ce cher Cheikh jamais de mauvaise humeur. Merci à Paul T. qui est sans aucun doute la personne m'ayant le mieux accompagné pendant cette thèse et ayant eu la plus grande indulgence devant mes erreurs... lorsqu'il jouait du violoncelle. Merci Paul J. pour m'avoir fait redécouvrir la beauté de la langue Française avec Cyrano de Bergerac. Aux autres doctorants qui bien assez vite deviendront à leur tour les plus anciens du bureau : Bowen, Louis, Alexis, Naoufal, Alann.

Je remercie les ami(e)s du groupe de Grégoire qui eux aussi dédient une partie de leur vie à l'optimisation de formes. Merci Mathilde, dont on envie le rayonnement permanent, pour tout ce débordement de bonne humeur que tu sais communiquer aux autres. Merci Perle pour les nombreux moments passés en école d'été, en conférence, (aux jardins botaniques), et d'avoir maintes fois œuvré pour nous rassembler entre doctorants de Safran ou bien du CMAP. Merci pour les figues (et les guêpes !), les pommes, et les analyses sur la psychologie de la personnalité féline. Merci Lalaina, Martin, Jeet, Alex, Beni, Alexis, Mathias pour les nombreuses expériences communes (notamment, de soupes, de fromages, mais par pitié sans crème fraîche ni gâteaux liquides) et d'échanges mathématiques enthousiastes.

Je suis bien également reconnaissant envers mes camarades doctorants de Safran. Je remercie Perle (pour m'avoir cité trois fois dans ses remerciements), Camille (pour la bienveillance), Adrien (qui reconnaîtra que l'Abondance est bien le meilleur des fromages), Maxence W. (pour son alignement et Baahubali), Yannis (pour son vocabulaire délicat), Clément B (pour avoir épargné mes genoux), Antony (pour une leçon d'anglais à ECCM), Thomas (pour son amour pour les ROM), Clément O. (pour ses recommandations de courses dans la vallée de la mère en thèse), Moubine (pour l'enthousiasme), et bien sûr Loïc J. (pour nous avoir transmis ton intérêt pour les aliments et boissons fermentées!). Un sincère remerciement pour l'ambiance empreinte de légèreté et d'humour qui a perduré durant ces trois années grâce à vos personnalités caractéristiques et colorées, et grâce aux nombreuses occasions qui nous ont réunis en dehors de Châteaufort.

J'adresse ensuite des remerciements sincères envers tous mes ami(e)s qui à leur manière ont contribué à maintenir la bonne santé mentale du thésard. Un *énorme*<sup>2</sup> merci à Maxence E. pour la fidélité et la sincérité de ton amitié qui me font te pardonner volontiers les cinq petites minutes d'attente systématiques. Merci Frédérique et Thibault (pour les déjeuners à l'INRIA après la piscine) Nicole et Benjamin (notamment pour une descente mémorable du pic rouge de Bassiès), Cécile, Sylvain (entre autres pour ce qui nous relie mutuellement à la France et l'Inde), Benjamin (merci pour ton invitation à une table ronde bien singulière), Christophe (au moins pour quelques bonnes rasades d'Orchata avec Anouk entre deux minisymposia, mais aussi pour la main gauche de Ravel), Amaury (avec qui il est toujours un plaisir de discuter), Chenzhang (pour ses conseils pianistiques et encouragements avisés), Guillaume (pour une visite à Zurich), Camille (pour m'avoir par exemple fait découvrir que les courgettes peuvent se croquer crues), et Marc (pour un beau week-end en Picardie).

Je remercie les précieux amis qui pensent toujours à moi depuis les U.S.A. (ou bien depuis Ferney-Voltaire International Airport) : *thank you* Viral, Dragosh, Wei Wei, Suhyoun, Parnika, Saviz, Salman, Ravi, Corbin, Rohit, Rameech, Chinmay, Arko ...

Enfin, je n'oublie pas l'éternelle *dream team* du lycée : merci Alex, Paul, David, Raphaël, Laurent, Cédric, Martin, Olivier de (presque) toujours parvenir à trouver le moyen d'être présent lors de mes venues parcimonieuses en contrées savoyardes.

Je termine par les remerciements à mes deux familles qui m'ont permis à leurs contacts de recharger et d'équilibrer régulièrement mes batteries d'énergie physique, mentale et émotionnelle. Ma famille

<sup>2</sup>énorme comment?



spirituelle tout d'abord, spécifiquement mes ami(e)s de l'École André Van Lysebeth auquel(le)s le  
reliement me nourrit d'instant en instant. En particulier, merci Joumana pour la confiance que tu me  
témoignes depuis mon arrivée à Orsay. Et puis, je remercie bien affectueusement mes racines terrestres,  
toutes les personnes de ma famille de chair qui m'estiment, m'encouragent et colorent ma vie depuis son  
commencement. J'adresse en particulier une pensée sincère et profonde pour Kévin, qui me rappelle à  
l'humilité.

Et avant de commencer, je Te remercie Toi, ô cher(ère) Lecteur(trice), pour ce qui porte Tes yeux à  
parcourir ce travail!



# CONTENTS

<b>Introduction – Résumé de la thèse</b>	<b>13</b>
<b>Introduction – Summary of the thesis</b>	<b>25</b>
<b>1 Shape and topology optimization based on Hadamard’s boundary variation method</b>	<b>37</b>
1.1 Notation . . . . .	37
1.2 The boundary variation method of Hadamard for shape sensitivity of PDE constrained problems . . . . .	38
1.3 On the signed distance function and its main properties . . . . .	48
1.4 A classical shape optimization numerical workflow using a level-set based mesh evolution method . . . . .	53
<b>2 Hadamard’s shape derivatives for a coupled thermal fluid structure problem</b>	<b>61</b>
2.1 Introduction . . . . .	61
2.2 Setting of the three-physic problem . . . . .	64
2.3 Shape derivatives for a simplified scalar fluid structure interaction problem . . . . .	68
2.4 Shape derivatives for the three-physic problem . . . . .	76
2.5 Numerical test cases . . . . .	80
2.6 Appendix . . . . .	100
<b>3 Null space gradient flows for constrained shape optimization</b>	<b>107</b>
3.1 Introduction . . . . .	107
3.2 Null space gradient flows for equality-constrained optimization in Hilbert spaces . . . . .	111
3.3 Extension to equality and inequality constraints . . . . .	117
3.4 Numerical discretization and time-stepping schemes for the null space ODE . . . . .	125
3.5 Comparisons with other methods and illustrations on academic test cases . . . . .	127
3.6 Optimization within the set of Lipschitz subdomains: applications to shape optimization . . . . .	138
<b>4 A variational method for computing shape derivatives of geometric constraints along rays</b>	<b>149</b>
4.1 Introduction . . . . .	149
4.2 Weighted graph space of the advection operator $\beta \cdot \nabla$ for velocity fields of class $\mathcal{C}^1$ . . . . .	152
4.3 Numerical methods for integration along normal rays . . . . .	163
4.4 Applications to maximum and minimum thickness constraints in shape optimization . . . . .	178
<b>5 Topology optimization of 2-d heat exchangers</b>	<b>189</b>
5.1 Design optimization of 2-d liquid-liquid heat exchangers with a non-mixing constraint . . . . .	189
5.2 Topology optimization of a 2-d air-oil heat exchanger . . . . .	194
<b>6 Towards 3-d and industrial applications: implementation recipes for a variety of numerical test cases</b>	<b>209</b>
6.1 Implementation recipes for 3-d constrained topology optimization of multiphysics system . . . . .	209
6.2 A few (moderately) large-scale three dimensional multiphysics applications . . . . .	219
<b>7 High order homogenized equations for perforated problems: towards fluid topology optimization by the homogenization method</b>	<b>239</b>
7.1 Introduction . . . . .	239
7.2 Motivations from shape optimization and summary of results . . . . .	240
7.3 High order homogenization for the perforated Poisson problem . . . . .	249
7.4 High order homogenization for the perforated elasticity system . . . . .	271
7.5 High order homogenization for the Stokes system in a porous medium . . . . .	287
<b>References</b>	<b>311</b>



## INTRODUCTION – RÉSUMÉ DE LA THÈSE

Cette thèse est consacrée à l’optimisation de la forme et de la topologie de systèmes multi-physiques motivés par des applications de l’industrie aéronautique. Nous employons principalement la méthode de Hadamard qui est particulièrement adaptée à la prise en compte de diverses contraintes industrielles détaillées ci-après. Cependant, nous effectuons dans un dernier chapitre indépendant des développements théoriques en ce qui concerne l’applicabilité future de la méthode d’optimisation topologique de systèmes fluides par homogénéisation dans un dernier chapitre indépendant.

L’optimisation des formes est l’art mathématique de générer des formes “optimales” qui satisfont au mieux un objectif proposé. Dans le contexte industriel, il s’agit de concevoir des systèmes physiques qui atteignent des performances optimales. Quelle est la forme de la structure la plus rigide utilisant une quantité donnée de matière ? Quel profil aérodynamique choisir pour générer une force de portance voulue ? Quelle distribution de fluides permet de réaliser le meilleur transfert de chaleur ? Comment concevoir la forme d’une coque de bateau ? Ces questions sont des problèmes industriels très classiques.

### 1 DE L’OPTIMISATION DE FORMES À L’OPTIMISATION TOPOLOGIQUE : LA MÉTHODE DE HADAMARD, LES MÉTHODES PAR DENSITÉ ET PAR HOMOGENÉISATION

La méthode de Hadamard tire ses origines du mémoire de 1908 [180]. Elle consiste à calculer la sensibilité du problème considéré à de petites déformations de la frontière (Figure 1). Ceci permet d’optimiser la forme via la détermination de déformations particulières qui donnent lieu à de meilleurs designs. De part

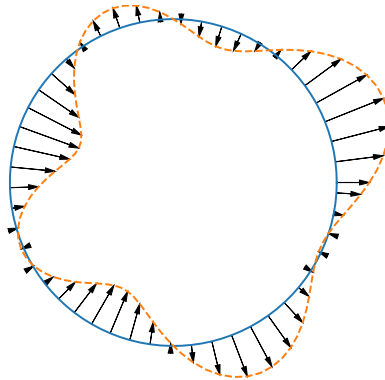


Figure 1: La méthode de variation de frontière de Hadamard : il s’agit d’évaluer la sensibilité du problème d’optimisation par rapport à de petites déformations de la forme.

l’essor de l’informatique qui a rendu possible la simulation numérique de systèmes physiques toujours plus complexes, de nombreux développements autour de la méthode ont été effectués dans les années 1970 avec les travaux fondateurs de Céa [86, 87], Murat, Simon [241, 242], Dervieux et Palmerio [251, 252] et Pironneau [259, 260]. Ces auteurs ont été parmi les premiers à proposer des cadres théoriques et numériques pour la résolution de problèmes d’optimisation de formes contraints par des équations aux dérivées partielles (EDPs). Ces techniques ont très vite suscité un engouement dans l’industrie et ont été suivies de nombreux articles envisageant des applications physiques variées : pour la conception optimale des structures mécaniques [290, 182, 57], en théorie des coques [89], pour des problèmes de conduction thermique [121, 83, 91], et pour des applications en mécanique des fluides [90, 172, 191, 192, 226, 123].

Depuis lors, la discipline de l’optimisation de formes n’a cessé de se développer et de s’améliorer quant aux aspects théoriques [78, 79, 186] ou numériques [179, 229]. Les premiers algorithmes de calcul de formes optimales étaient basés sur la technique de déformation de maillages [260, 291, 34] : une forme candidate est discrétisée en un maillage, qui est déformé de manière itérative afin d’obtenir une géométrie améliorée. (Figure 2). Une avancée majeure a eu lieu au début des années 2000 avec l’introduction de la méthode des lignes de niveaux dans les algorithmes d’optimisation de formes [32, 311, 250] : la forme à optimiser n’est pas maillée explicitement mais plutôt capturée de manière implicite comme l’ensemble des valeurs négatives d’une fonction “*lignes de niveaux*”, ce qui permet de capturer l’évolution de la forme sur un maillage fixe. Cette méthode permet l’apparition de changements topologiques complexes de la forme au cours des itérations, tels que la fusion de trous ou la fusion de frontières (Figure 3). Cette capacité de ces types d’algorithmes à gérer les changements topologiques permet de s’affranchir de tout *a priori* sur la topologie et sur la géométrie de la forme à optimiser, tel que le nombre ou la localisation des

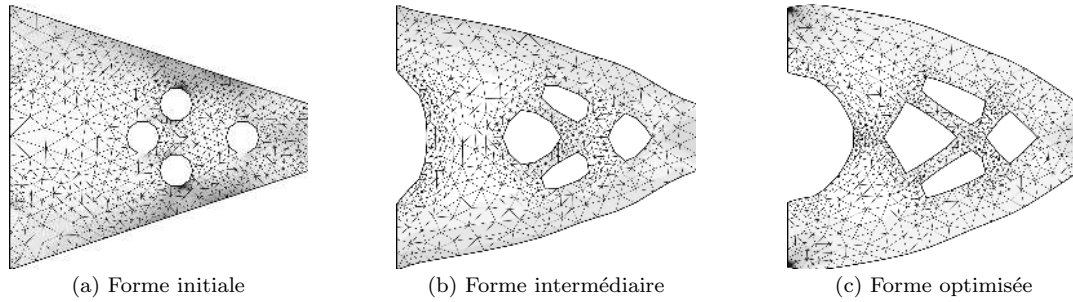


Figure 2: La méthode de variation de frontière de Hadamard utilisant la déformation d’un maillage pour l’optimisation de la forme d’une console 2-d en flexion (figure extraite de [17]). Les changements topologiques sont difficiles à traiter numériquement : les formes initiales et finales ont le même nombre de trous.

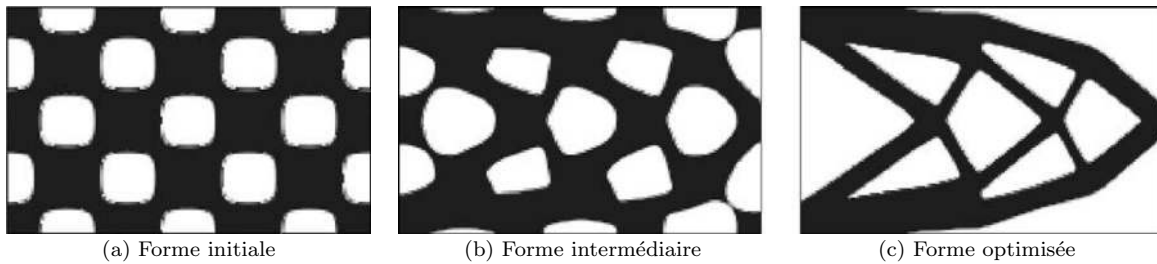


Figure 3: La méthode de variation de frontières de Hadamard utilisant la méthode des lignes de niveaux pour l’optimisation de la forme d’une console 2-d en flexion (figures extraites de [31]). Les changements topologiques sont possibles : certains trous ont fusionné entre l’initialisation et la forme finale.

trous (ou bien des arches en 3-d). Pour cette raison, les algorithmes basés sur la méthode des lignes de niveaux sont souvent considérés comme des techniques d’*optimisation topologique*, tandis que ceux basés sur la déformation d’un maillage de la forme est traditionnellement appelée *optimisation géométrique* [17]. Une tendance récente suggère d’intégrer des techniques de remaillage à la démarche d’optimisation topologique de façon à combiner les avantages des deux méthodes [24, 107]. Il est ainsi possible de garder une discrétisation explicite des formes au cours de l’optimisation tout en permettant à des changements topologiques d’avoir lieu (Figure 4).

L’optimisation topologique est une discipline très large qui bénéficie d’un large panel de techniques mathématiques et numériques. Nous pouvons mentionner, de manière non exhaustive, la méthode du gradient topologique [85, 145, 165, 289, 40], les méthodes par lignes de niveaux qui ne sont pas spécifiquement basées sur la méthode de Hadamard [133, 316], les méthodes de changement de phase [73, 299], les méthodes d’optimisations topologiques dites “évolutionnaires” [193], et enfin les méthodes par densité [66] et par homogénéisation [18, 93] (nous renvoyons à [284] pour une revue plus complète du domaine). Parmi celles-ci, les méthodes par densité n’optimisent pas explicitement la frontière de la forme mais plutôt une fonction de densité  $\rho$  qui modélise la présence de matériau ( $\rho = 1$ ) ou bien son absence ( $\rho = 0$ ), ou bien des états “intermédiaires” ( $0 < \rho < 1$ ). Des schémas numériques de pénalisation doivent être utilisés de façon à obtenir la convergence vers des formes “admissibles”, c’est-à-dire des états purement “noirs et blancs” ( $\rho = 0$  ou bien  $\rho = 1$  comme sur la Figure 5). Les méthodes par homogénéisation remplacent quant à elles le problème d’optimisation de formes par la recherche d’une microstructure optimale pour un matériau composite effectif. Une première étape consiste à optimiser les paramètres de la microstructure tels que la densité de matière ou bien l’orientation des cellules. Dans un second temps, si aucune pénalisation n’a été utilisée, ces champs de paramètres sont interprétés afin de déterminer une forme qui approche convenablement la microstructure optimale (Figure 6) [254, 27, 166, 177].

## 2 QUELQUES ENJEUX ACTUELS DE L’OPTIMISATION TOPOLOGIQUE POUR LES APPLICATIONS DE L’INDUSTRIE AÉRONAUTIQUE

Traditionnellement, les ingénieurs conçoivent les systèmes industriels grâce à l’aide de logiciels de Conception Assistée par Ordinateur (CAO). Les industries sont souvent dépendantes des formats CAO du fait de leur compatibilité avec toutes les étapes de conception, de la simulation numérique des pro-

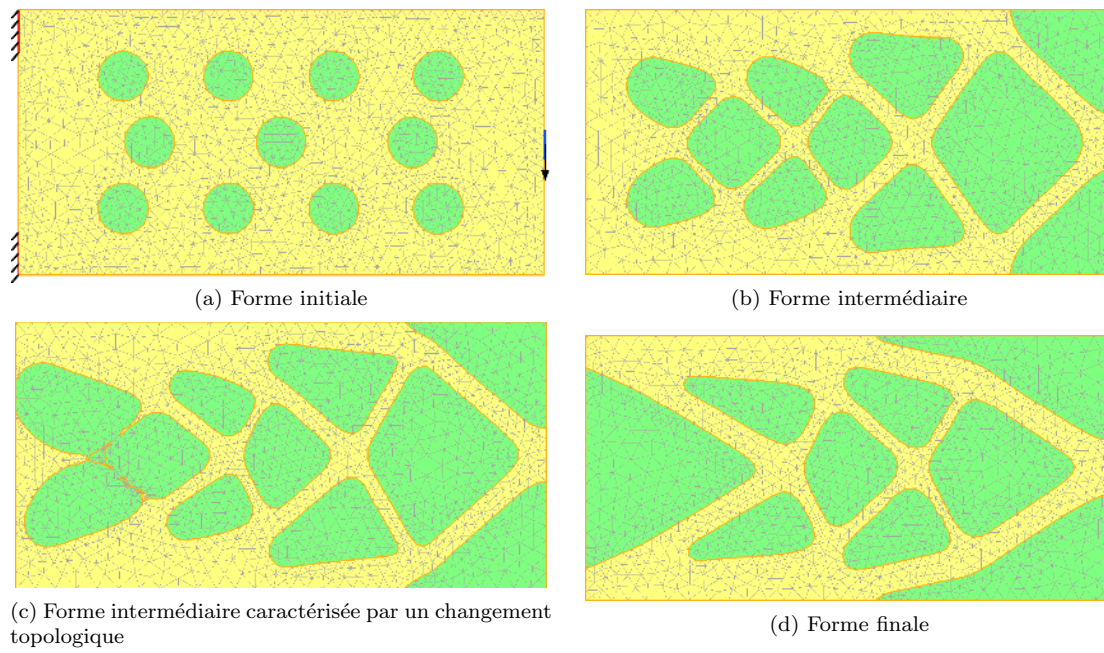


Figure 4: La méthode de variation de frontière de Hadamard implémentée en utilisant la méthode des lignes de niveaux combinée à une technique de remaillage pour l'optimisation d'une console 2-d en flexion (figures extraites de [107]). Les changements topologiques sont permis tout en conservant une discrétisation explicite et conforme des formes.

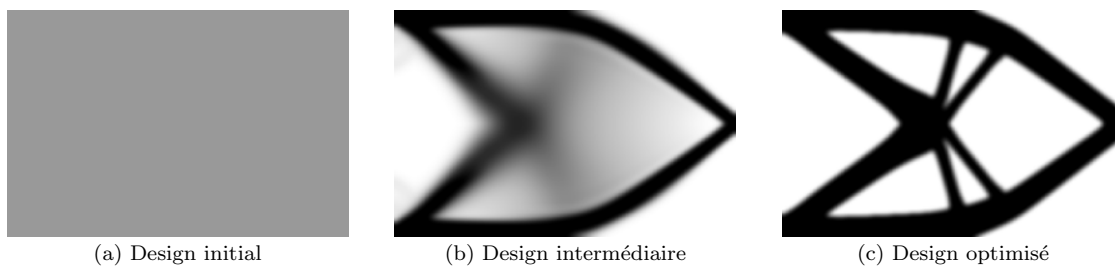


Figure 5: Optimisation topologique d'une console 2-d en flexion par une méthode de densité (ici la méthode SIMP—Solid Isotropic Material Penalization—). Figure obtenue à partir du code source décrit dans [43].

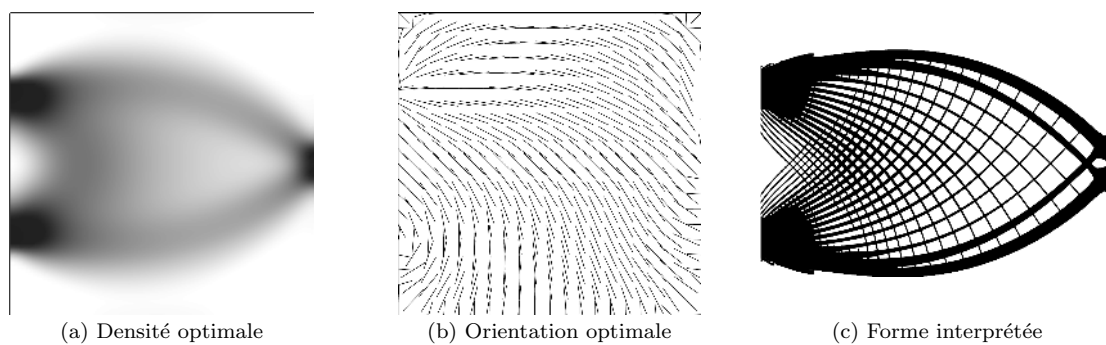


Figure 6: Optimisation topologique d'une console 2-d en flexion par une méthode d'homogénéisation. Figure extraite de [166].

cessus physiques par des codes industriels à la fabrication effective par des machines automatisées. Lors de l’optimisation d’une forme au moyen d’outils CAO, la topologie du système à concevoir est généralement proposée *a priori* par l’ingénieur. Elle est par exemple déterminée par un ensemble de points de contrôles de courbes de Bézier, par des paramètres de splines, ou bien par d’autres types de paramètres géométriques. Une optimisation paramétrique de ces paramètres de conception est alors effectuée pour améliorer le design initial [181]. Habituellement, la sensibilité du problème par rapport à la géométrie n’est pas connue de manière analytique. Celle-ci est parfois estimée par différentiation automatique (un outil très populaire en conception aérodynamique [94, 139, 130]) [235, 288].

Cependant, dans de nombreux cas, l’optimisation est effectuée par l’exploration d’un sous-ensemble suffisamment vaste de l’espace des paramètres appelé “plan d’expériences” (cela est fait par exemple avec le logiciel commercial Optimus [245]), voir par exemple [142, 286]. Puisque ces méthodes reposent sur

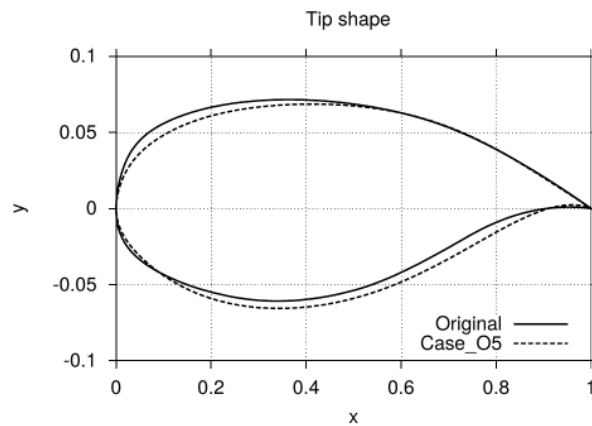


Figure 37. Tip wing section. Original vs optimized ones.

Figure 7: Optimisation de la CAO d’un profil aérodynamique pour le couple portance–traînée. Figure extraite de [142].

des choix importants de paramétrisation de la géométrie, elles ne conduisent généralement qu’à de très faibles modifications de la forme initialement proposée (Figure 7). En principe, l’optimisation topologique pourrait permettre d’obtenir des designs aux performances sensiblement meilleures, parmi un ensemble beaucoup moins restreint de formes admissibles.

Les développements de la fabrication additive depuis les années 1990 permettent aujourd’hui à l’industrie de fabriquer des systèmes à la géométrie toujours plus complexe, difficilement décrite par une paramétrisation CAO ; ceci suscite aujourd’hui un enthousiasme renouvelé pour l’optimisation topologique. Un grand nombre d’applications envisagées dans les travaux de recherche actuels sont issus de l’industrie aéronautique. Celles-ci constituent quelques-unes des motivations principales à l’origine de ce travail, par exemple lorsqu’il est question de la conception des composants de moteurs d’avion. Un des enjeux majeurs, en ce qui concerne l’application de l’optimisation topologique aux systèmes aéronautiques, est la nécessité de prendre en compte leur caractère intrinsèquement multi-physique : il est bien souvent souhaitable d’optimiser de tels systèmes au regard de critères impliquant des propriétés thermiques, mécaniques et hydrauliques couplées.

Un problème d’actualité est celui de l’optimisation topologique appliquée à la conception d’échangeurs de chaleur [255, 258, 189, 271, 275]. Les échangeurs de chaleurs sont des dispositifs utilisés dans les moteurs pour refroidir des fluides chauds en les transportant à proximité de certains gaz ou liquides réfrigérants. Les modèles industriels présentent généralement un assemblage de nombreux tubes et ailettes afin de maximiser la surface d’échange entre les phases chaudes et froides (Figure 8). Naturellement, il est aussi nécessaire de prendre en compte un ensemble de contraintes multi-physiques lors de la conception, telles que la perte de charge (la différence de pression entre l’entrée et la sortie du dispositif), ou bien la résistance des structures mécaniques en présence d’une charge thermique élevée.

De nombreux autres composants des moteurs aéronautiques pourraient certainement bénéficier de l’optimisation de formes, à commencer par le système de refroidissement par canaux internes des aubes de turbines [62, 162] (Figure 9).



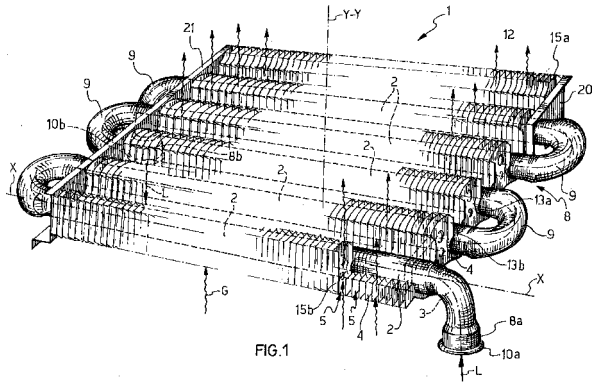


Figure 8: Échangeur de chaleur industriel liquide-gaz. Figure extraite de [169].

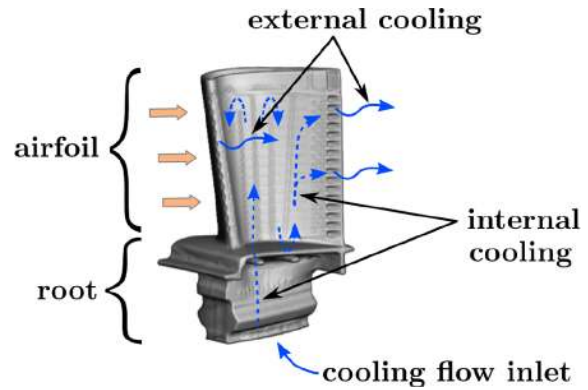


Figure 9: Système de refroidissement d'aubes de turbines avec présence de canaux internes réfrigérants. Figure extraite de [162].

### 3 CADRE DE LA THÈSE

À ce jour, la méthode de Hadamard utilisée à des fins d'optimisation topologique (dans laquelle on autorise des déformations libres de la forme optimisée) n'est pas encore suffisamment mature pour être intégrée à de véritables applications industrielles et multi-physiques. La plupart des cas tests traités dans la littérature se placent le plus souvent dans le cadre de l'élasticité linéaire dans [25, 108], ou dans celui de la conduction thermique [9, 325], et bien peu d'exemples numériques tridimensionnels, à résolution élevée, sont présentés [202]. Seulement récemment, la méthode a été appliquée à des physique plus complexes telles que la plasticité mécanique [230] ou bien l'électromagnétisme [160, 213]. Cette observation s'applique également à la communauté d'optimisation par méthodes de densités qui présente un nombre croissant de récents travaux portant sur des applications multidisciplinaires [116, 129, 128, 318, 319, 240, 288].

L'objectif de ce travail est de développer des outils théoriques et numériques pour l'optimisation de la forme et de la topologie de systèmes multi-physiques, guidé par la prise en compte d'un certain nombre d'exigences industrielles à long terme. À cet effet, nous utilisons principalement la méthode de variation de frontière de Hadamard ; un chapitre d'ouverture (le chapitre 7) est dédié à des développements mathématiques sur l'homogénéisation des milieux poreux qui pourraient permettre d'envisager des applications nouvelles pour l'optimisation de systèmes impliquant des fluides.

Par "exigences industrielles", nous entendons plusieurs besoins industriels qui sont autant d'axes majeurs pour notre travail et que nous décrivons dans les cinq prochains paragraphes. Les contributions de la thèse sont ensuite résumées chapitre par chapitre dans la section qui suit.

#### Systèmes multi-physiques

Comme en témoignent les applications industrielles évoquées dans la la section précédente, il existe une demande croissante en ce qui concerne l'optimisation topologique pour des applications mêlant plusieurs physiques en interaction. Notre étude se concentre sur des systèmes fluides présentant des propriétés hydrauliques, thermiques et mécaniques couplées. Ces derniers sont caractérisés par les variables physiques suivantes:

- $\mathbf{v}$  et  $p$  désignent les champs de vitesse et de pression associés à un ou plusieurs fluides évoluant à travers le système ;
- $T$  pour le champs de température dans les phases solides et liquides ;
- $\mathbf{u}$  pour le champs de déplacement élastique des structures mécaniques solides mises en jeu dans le système.

Mathématiquement, ces variables sont déterminées par la résolution d'un système d'équations aux dérivées partielles, ces dernières provenant elles-mêmes du choix de la modélisation physique. Un cadre académique suffisamment représentatif pour nos applications, est décrit dans le chapitre 2.

En termes d’applications, une partie de notre travail est guidée par la conception d’échangeurs de chaleur. Une étude numérique dans ce contexte est proposée dans le [chapitre 5](#), où nous optimisons la forme et la topologie de deux modèles d’échangeurs de chaleur bidimensionnels.

### Optimisation sous contraintes

Lors de la conception d’un système industriel, il est d’usage de prendre en compte un cahier des charges dont les spécifications doivent être satisfaites dans les conditions d’utilisation réelles. Par exemple, les contraintes mécaniques ou thermiques d’une structure solide ne doivent pas dépasser un seuil donné afin d’éviter un endommagement prématuré. En d’autres termes, un problème de conception optimale consiste à déterminer la forme du système qui réalise la meilleure performance pour un ensemble donné de contraintes physiques à respecter.

De tels problèmes sont modélisés de manière générique par des programmes mathématiques de la forme :

$$\begin{aligned} \min_{\Gamma} \quad & J(\Gamma, \mathbf{v}(\Gamma), p(\Gamma), T(\Gamma), \mathbf{u}(\Gamma)) \\ \text{s.t.} \quad & \begin{cases} g_i(\Gamma, \mathbf{v}(\Gamma), p(\Gamma), T(\Gamma), \mathbf{u}(\Gamma)) = 0, & 1 \leq i \leq p, \\ h_j(\Gamma, \mathbf{v}(\Gamma), p(\Gamma), T(\Gamma), \mathbf{u}(\Gamma)) \leq 0, & 1 \leq j \leq q, \end{cases} \end{aligned} \quad (1)$$

où  $\Gamma$  est la frontière du système à optimiser, typiquement l’interface entre une structure solide et un fluide (ce qui laisse la possibilité de rendre compte de plusieurs phases liquides).  $J$  est une fonction de coût à minimiser qui quantifie la performance du système. Par exemple, la compliance d’une structure mécanique, la force de traînée générée par un obstacle, ou encore la chaleur emmagasinée par un système thermique sont des fonctions de coût classiques couramment rencontrées en optimisation de la forme. Les fonctions  $g_i$  et  $h_j$  modélisent des contraintes d’égalité et d’inégalité: elles renvoient à des valeurs cibles que certaines quantités physiques doivent atteindre (par exemple, un emplacement souhaité pour le centre de gravité de la structure) ou bien ne doivent pas dépasser (par exemple, une limite imposée sur la température du système, ou bien une taille caractéristique minimale). La fonction objectif  $J$  et les contraintes  $g_i, h_j$  dépendent à la fois de la forme  $\Gamma$  et des variables physiques  $\mathbf{v}(\Gamma), p(\Gamma), \mathbf{u}(\Gamma), T(\Gamma)$ , qui dépendent elles-mêmes de  $\Gamma$  via les équations d’états.

Dans le contexte des méthodes de densité, la variable optimisée  $\rho$  appartient à un espace vectoriel de type  $\mathbb{R}^n$  après discrétisation, ce qui permet de résoudre des programmes mathématiques contraints par des méthodes très classiques d’optimisation du premier ordre [[244](#), [298](#), [310](#)]. Dans le cadre de la méthode de Hadamard, la variable d’optimisation  $\Gamma$  n’appartient pas à un espace vectoriel mais plutôt à une “variété” de dimension infinie ce qui rend la résolution de (1) beaucoup plus délicate. La plupart des travaux de la littérature s’appuyant sur la méthode de Hadamard convertissent (1) en un problème de minimisation sans contrainte via l’ajout de termes de pénalisation à la fonction objectif  $J$ ; une technique couramment employée à cet égard est la méthode du Lagrangien augmenté [[234](#), [107](#)]. Des méthodes plus complexes, telles que l’optimisation linéaire successive (*Sequential Linear Programming, SLP*) ou la méthode des directions admissibles (*Method of Feasible Directions, MFD*) ont été utilisées avec succès dans un nombre restreint de travaux [[135](#), [150](#)]. Cependant, ces méthodes ne sont pas complètement satisfaisantes pour une utilisation industrielle car elles nécessitent un réglage fin de paramètres d’optimisation qui dépendent fortement de la situation considérée et qui peuvent être difficiles à ajuster, surtout en présence d’un grand nombre de contraintes.

Dans le [chapitre 3](#), nous proposons un nouvel algorithme de type gradient pour la résolution de problèmes d’optimisation sous contraintes adapté au contexte de la méthode de Hadamard.

### Algorithmes non-intrusifs

Dans un souci de compatibilité avec les outils industriels, le processus d’optimisation de formes doit être entièrement compatible avec la chaîne de conception du système. Par exemple, les outils d’optimisation topologique doivent idéalement pouvoir évaluer les variables d’état en faisant appel aux mêmes solveurs physiques que ceux utilisés lors de la phase de validation industrielle. De nombreuses méthodes d’optimisation topologique, notamment les méthodes de densité ou de lignes de niveaux, ne remplissent pas cette condition car elles reposent sur des modifications artificielles des équations d’état pour des raisons purement numériques.

Ce besoin constitue l’une des raisons de notre utilisation systématique pour tous nos cas tests numériques de l’algorithme d’évolution de maillages proposé par [[25](#)]. Cet algorithme illustré sur la [Figure 4](#) et ses grands principes sont rappelés dans le [chapitre 1](#). Cette technique, basée sur la méthode

des lignes de niveaux, permet de conserver une représentation maillée explicite de la forme optimisée durant toutes les itérations du processus d'optimisation, ce qui permet de résoudre les équations d'état sans aucune modification ou approximation de la physique considérée. Cette étape de résolution pourrait en principe, être effectuée par un solveur tiers. Pour nos applications, nous avons utilisé de manière intensive le logiciel open source `mmg` [108] pour les étapes de remaillage 2-d ou 3-d.

Une nécessité supplémentaire imposée par le souhait de non-intrusivité des algorithmes d'optimisation de formes est la minimisation de la quantité d'information demandée à l'utilisateur. En particulier, un utilisateur formulant un problème du type (1) devrait pouvoir renseigner des fonctions de coûts et des contraintes arbitraires sans avoir à connaître les détails mathématiques (plutôt techniques) du calcul des dérivées de forme au sens de Hadamard. Dans le chapitre 2, nous établissons des formules pour les dérivées de formes de fonctionnelles de coût *arbitraires*, qui ont l'avantage de pouvoir être assemblées numériquement de manière automatique à partir de de la seule connaissance des dérivées partielles (qui classiquement, sont très simples à calculer).

### Contraintes géométriques

La prise en compte de contraintes géométriques quant aux formes optimisées est un besoin industriel classique en optimisation topologique du fait de la précision limitée des processus de fabrication. Elles peuvent par exemple se manifester par des contraintes sur l'épaisseur minimale d'une structure ou bien sur l'angle maximal de ses parties en porte-à-faux avec la direction verticale. Elles se présentent également lors de la conception des échangeurs de chaleur bi-tubes : dans le chapitre 5, nous modélisons une condition de non-mélange de deux phases fluides par une contrainte de distance géométrique.

De nombreux travaux ont analysé l'intégration de telles contraintes en optimisation topologique, voir par exemple [234], chapitre 3, pour une revue. Dans le contexte de l'optimisation de formes par la méthode de Hadamard, il est commode de modéliser les contraintes géométriques au moyen de la fonction de distance signée [30]. Cependant, l'implémentation pratique des dérivées de forme associées est très délicate et dépend sensiblement du type de discrétisation envisagée pour la forme optimisée. Dans le chapitre 4, nous proposons une nouvelle méthode variationnelle qui facilite l'évaluation numérique des dérivées de forme pour les contraintes géométriques de distance, et qui peut être implémentée de manière très commode par la méthode des éléments finis.

### Calcul haute performance

Une dernière exigence industrielle importante concerne l'application de l'optimisation topologique à des problèmes physiques 3-d et discrétisé par des maillages ayant des résolutions élevées. Lors d'un processus d'optimisation, les équations d'état doivent être résolues de l'ordre d'une à plusieurs centaines de fois, ce qui est très coûteux en temps de calcul pour les problèmes dits de "grande échelle" (pour lesquels l'étape de résolution requiert la détermination d'un très grand nombre de variables). Ce coût est d'autant plus élevé lorsque les physiques mises en jeu comprennent des écoulements de fluides. Cette question a été assez largement traitée par la communauté utilisant les méthodes de densité [1, 10, 71, 232, 39], et elle commence à être également abordée dans le cadre des algorithmes d'optimisation topologiques dits "évolutionnaires" [227] ainsi que de ceux basés sur la méthode des lignes de niveaux [202]. Dans le chapitre 6, nous présentons et discutons notre implémentation en `FreeFEM` [183] de cas tests 3-d présentant des tailles de maillages relativement importantes (moyennant la résolution de systèmes linéaires jusqu'à 2 millions de degrés de liberté). À cet effet, nous exploitons le calcul parallèle et des techniques récentes de décomposition de domaines associées à des solveurs itératifs préconditionnés [238].

Dans un autre registre, les modèles industriels classiques de conceptions des échangeurs de chaleur (tel que celui illustré sur la Figure 8) suggèrent le besoin de générer des designs présentant des motifs multi-échelles. Dans le chapitre 7, nous proposons des modèles homogénéisés d'ordres supérieurs pour les écoulement de fluides en milieux poreux, qui pourraient permettre d'envisager de traiter ces applications par les méthodes d'optimisation topologique par homogénéisation dans de tels contextes.

## 4 RÉSUMÉ PAR CHAPITRE

### Chapitre 1 : optimisation de formes par la méthode de variation de frontière de Hadamard

Ce chapitre préliminaire présente le contexte et le matériel technique de base concernant l'optimisation de formes par la méthode de Hadamard. Nous décrivons les ingrédients classiques requis pour la mise en œuvre d'un processus d'optimisation topologique via la méthode d'évolution de maillages de [25].

Plus précisément, nous commençons par un survol de la méthode de Hadamard et des techniques mathématiques classiques qui permettent de calculer les dérivées de forme de fonctions de coût dépendant de solutions d'équations aux dérivées partielles. Nous discutons ensuite du problème classique de l'adaptation de l'algorithme d'optimisation par la méthode du gradient au contexte de la dimension infinie. Une partie indépendante est dédiée à des rappels sur les définitions et les principales propriétés de la fonction de distance signée. Enfin nous résumons les différentes étapes de la méthode d'évolution de maillage de Allaire, Dapogny et Frey [25] pour l'optimisation de formes.

## Chapitre 2 : Dérivées de formes au sens de Hadamard pour un problème hydraulique thermique mécanique faiblement couplé

Ce chapitre présente un modèle thermique fluide-structure simplifié qui fait office de référence pour toutes nos applications numériques: les variables  $\mathbf{v}, p, T, \mathbf{u}$  sont définies comme les solutions de trois équations d'état qui sont faiblement couplées dans le sens où ces dernières peuvent être résolues successivement et indépendamment. Nous considérons les équations de Navier-Stokes incompressibles et stationnaires pour la vitesse et la pression du fluide, l'équation de convection-diffusion pour le champ de température au sein des phases fluides et solides, et le système de la thermo-élasticité linéaire pour le déplacement élastique de la structure solide. Ce modèle est une généralisation de plusieurs situations impliquant une seule physique qui sont couramment étudiées dans la littérature ; ils permettent en outre de traiter de nouveaux cas test couplés tels que l'interaction fluide-structure dans un régime de petites déformations, ou bien le transfert de chaleur par convection.

Une nouveauté importante proposée par cette partie est la détermination de formules de dérivées de forme pour des fonctions de coût *arbitraires*. Dans la littérature consacrée à la méthode de Hadamard, de telles formules reposent en effet sur des hypothèses sur la structure de la fonction de coût considérée. Ainsi, il s'avère nécessaire de refaire le calcul des dérivées de forme lors de chaque changement de fonction objectif ou contrainte, ce qui va à l'encontre d'une intégration systématique dans des logiciels industriels. Nos formules permettent au contraire d'automatiser le calcul des dérivées de forme : celles-ci sont assemblées après la résolution d'équations adjointes à partir de la seule connaissance des dérivées *partielles* de la fonction de coût, une information qui peut en principe être fournie facilement par un utilisateur externe.

Une autre contribution de ce chapitre concerne la prise en compte spécifique des problèmes d'interaction fluide-structure. Une propriété surprenante réside dans le couplage du système adjoint : l'égalité des contraintes normales sur l'interface fluide-structure (une condition de type "Neumann") se traduit par une égalité des valeurs des variables adjointes sur cette frontière (une condition de type "Dirichlet").

La validité de nos formules est enfin vérifiée numériquement sur plusieurs cas tests à deux dimensions impliquant simultanément une, deux ou trois des physiques mentionnées.

Le contenu de ce chapitre a été publié pour l'essentiel dans l'article [153]:

F. FEPPON, G. ALLAIRE, F. BORDEU, J. CORTIAL, AND C. DAPOGNY, *Shape optimization of a coupled thermal fluid-structure problem in a level set mesh evolution framework*, SeMA Journal, (2019), pp. 1–46.

Mentionnons cependant que nous avons ajouté de nombreux exemples numériques supplémentaires et amélioré certains cas tests de la publication [153].

## Chapitre 3 : flots de gradient à noyaux pour l'optimisation de formes sous contrainte

Ce chapitre présente l'algorithme d'optimisation que nous avons spécifiquement développé dans la perspective de son application en optimisation de formes. Nous proposons un schéma d'optimisation qui peut être interprété comme la discrétisation d'un flot de gradient capable de "voir" les contraintes d'égalité et d'inégalité (il s'agit d'une généralisation de certains algorithmes d'optimisation par systèmes dynamiques [317, 300]). Ce flot de gradient est appelé à *noyau* car la direction d'optimisation s'écrit comme la somme d'un vecteur appartenant à l'espace tangent aux contraintes (*i.e.* au noyau de leur différentielle) qui vise à faire décroître les valeurs de la fonction objectif, et d'un vecteur orthogonal au premier et qui ramène le chemin d'optimisation à l'intérieur du domaine admissible. Le calcul de ces vecteurs dépend de la résolution d'un sous-problème quadratique qui indique quand décoller de la frontière des contraintes pour revenir à l'intérieur du domaine d'optimisation (Figure 10). Une caractéristique de l'algorithme est son aptitude à diminuer les valeurs de la fonction objectif tout en maintenant les contraintes satisfaites.

Nous établissons des propriétés de convergence de notre algorithme pour les trajectoires d'optimisation continues ; celles-ci garantissent l'amélioration de solutions optimisées jusqu'à ce qu'un minimum local

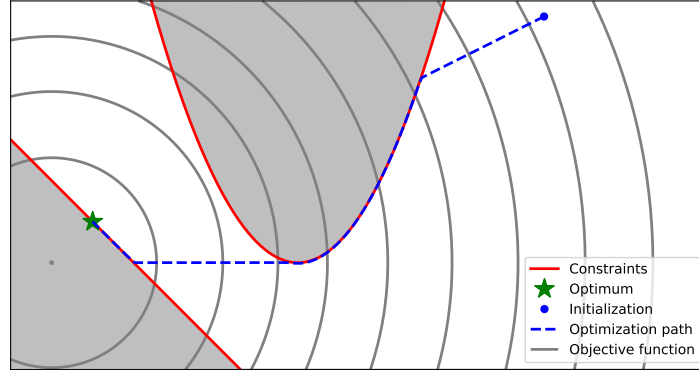


Figure 10: Illustration de notre flot de gradient à noyau pour la résolution des problèmes d'optimisation sous contraintes. Les trajectoires d'optimisation suivent toujours la meilleure direction de descente possible.

soit trouvé, à une discrétisation suffisamment fine près. Nous illustrons la robustesse et l'efficacité de notre méthode d'optimisation de formes sur un cas test de pont en chargement multiple présentant 10 contraintes d'optimisation.

La plupart du contenu de ce chapitre est à paraître dans la prépublication soumise [155]:

F. FEPPON, G. ALLAIRE, AND C. DAPOGNY, *Null space gradient flows for constrained optimization with applications to shape optimization*, submitted, (2019).

Le chapitre contient en outre quelques comparaisons de notre méthode avec des algorithmes plus classiques sur des exemples académiques.

#### Chapitre 4 : une méthode variationnelle pour calculer les dérivées de formes de contraintes géométriques le long des rayons

Les applications envisagées de l'optimisation topologique aux échangeurs de chaleur bi-tubes nécessitent de prendre en compte une contrainte de non-mélange des phases fluides misent en jeu lors de la résolution du problème d'optimisation. Une approche très naturelle pour ce faire consiste à prescrire une distance minimale entre les deux phases non miscibles, ce qui peut être formulé mathématiquement grâce à la fonction de distance signée. Plus généralement, de nombreuses autres contraintes géométriques peuvent être formulées de manière analogue, telles que les contraintes d'épaisseurs minimales ou maximales.

Des travaux précédents [30, 234] ont proposé des expressions mathématiques pour calculer les dérivées de formes de fonctionnelles de coût faisant intervenir la fonction de distance signée, afin d'imposer des contraintes géométriques telles que l'épaisseur maximale ou minimale de la forme optimisée. Cependant, la mise en œuvre directe de ces formules est difficile car elle nécessite une intégration numérique le long des rayons normaux au bord de la forme (Figure 11), ainsi que des estimations précises du squelette et des courbures principales. L'implémentation de ces opérations est notoirement difficile et est très dépendante de la dimension (2-d ou 3-d) ainsi que du type de discrétisation envisagé pour capturer la forme optimisée (implicite ou explicite, sur maillage structuré ou non structuré).

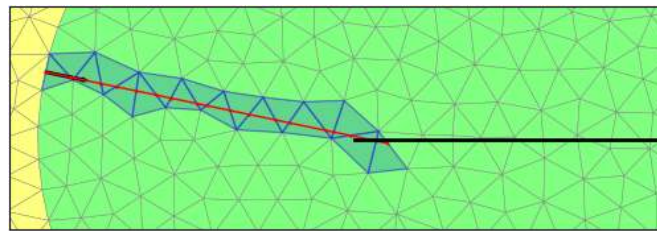


Figure 11: Les méthodes classiques pour calculer les dérivées de formes des contraintes géométriques (impliquant la fonction de distance signée) requièrent le calcul des rayons normaux à la forme (en rouge), ce qui nécessite de parcourir la discrétisation de maille en maille.

Dans ce chapitre, nous montrons que les dérivées de forme de contraintes géométriques peuvent être calculées au moyen de la résolution d'un problème variationnel, qui peut être aisément effectuée par la

méthode des éléments finis : la valeur des intégrales le long des rayons est obtenue à partir des valeurs de cette solution variationnelle aux sommets de la frontière de la forme. En toute généralité, notre méthode permet de calculer des intégrales le long des courbes caractéristiques d'un champ de vitesse  $\beta$  sans avoir à calculer explicitement ces courbes sur la discrétisation spatiale. L'implémentation s'en trouve considérablement simplifiée car elle ne nécessite pas les calculs de rayons, ni celui des courbures de la forme, mais uniquement la détermination d'un poids relativement arbitraire s'annulant approximativement sur le squelette.

Nous établissons le caractère bien posé de la formulation variationnelle proposée grâce à une analyse détaillée d'espaces à poids associés à l'opérateur d'advection  $\beta \cdot \nabla$ . Dans le contexte de l'optimisation de formes,  $\beta$  est le gradient de la fonction de distance signée au domaine sur lequel on souhaite imposer des contraintes géométriques, mais notre analyse permet de traiter également des champs de vecteurs plus généraux. Notre approche permet notamment de traiter des champs de vitesse ayant une divergence non bornée: les travaux classiques effectuant des analyses similaires [144] supposent généralement  $\text{div}(\beta) \in L^\infty(D)$  où  $D$  désigne le domaine de travail. Cette hypothèse est cependant systématiquement violée pour les applications considérées en optimisation de formes.

Nous montrons enfin l'efficacité de cette méthode variationnelle en revisitant l'implémentation des contraintes d'épaisseurs maximale et minimale en optimisation de formes. Nous retrouvons des résultats numériques analogues à ceux obtenus dans les travaux antérieurs [30, 234] s'appuyant sur le calcul explicite (difficile) des intégrales le long des rayons.

La majeure partie de ce chapitre est à paraître dans la publication [154]:

F. FEPPON, G. ALLAIRE, AND C. DAPOGNY, *A variational formulation for computing shape derivatives of geometric constraints along rays*, To appear in M2AN, (2019).

## Chapitre 5 : optimisation topologique d'échangeurs thermiques 2-d

Les ingrédients développés dans les parties précédentes sont appliqués dans ce chapitre au problème de l'optimisation d'échangeurs de chaleur bidimensionnels.

Dans une première partie, nous mettons en œuvre nos méthodes pour l'optimisation topologique d'échangeurs de chaleur bi-tube dont le comportement physique est décrit par le modèle faiblement couplé du chapitre 2. Nous formulons une contrainte de non mélange des deux phases liquides au moyen d'une contrainte de géométrie impliquant la fonction de distance signée à l'une des phases, que nous traitons par la méthode variationnelle du précédent chapitre 4 (Figure 12a).

Dans une seconde partie, nous rendons compte de la réalisation d'un cas test d'optimisation de formes pour des échangeurs de chaleur 2-d issu d'une collaboration avec Safran Aero Boosters : l'objectif est de déterminer la forme optimale de la section de canaux d'huiles transverses et réfrigérés par un flux d'air froid en entrée. La physique de ce cas test est modélisée par les équations d'états introduites dans le chapitre 2, à un changement de condition aux limites près pour le champ de température. Nous avons également intégré une contrainte d'épaisseur minimale pour la phase huileuse, et une contrainte sur la perte de charge maximale pour la phase constituée d'air. Notre étude démontre les capacités de notre méthodologie numérique d'optimisation de formes pour générer des designs non classiques dans un contexte simplifié (Figure 12b).

## Chapitre 6 : vers des applications 3-d et industrielles : stratégies d'implémentation pour une variété de cas tests numériques

Ce chapitre traite de l'implémentation de cas tests numériques qui approchent des applications industrielles plus avancées.

Une première section décrit succinctement certains paradigmes d'implémentation de notre code d'optimisation topologique développé en `python` et `FreeFEM` [183] pour les applications de cette thèse. Nous discutons ensuite des détails d'implémentation spécifiques à la 3-d, avec en particulier l'utilisation inévitable de techniques de décomposition de domaine associées à des préconditionneurs spécifiques aux physiques considérées.

Dans une seconde partie, nous présentons plusieurs nouveaux résultats d'optimisation de formes tridimensionnelle pour différentes physiques. Nos cas tests considèrent l'optimisation de la forme de consoles mécaniques en traction ou torsion, d'une distribution de matériaux pour la conduction thermique (Figure 13), de profils aérodynamiques permettant de générer une portance maximale tout en

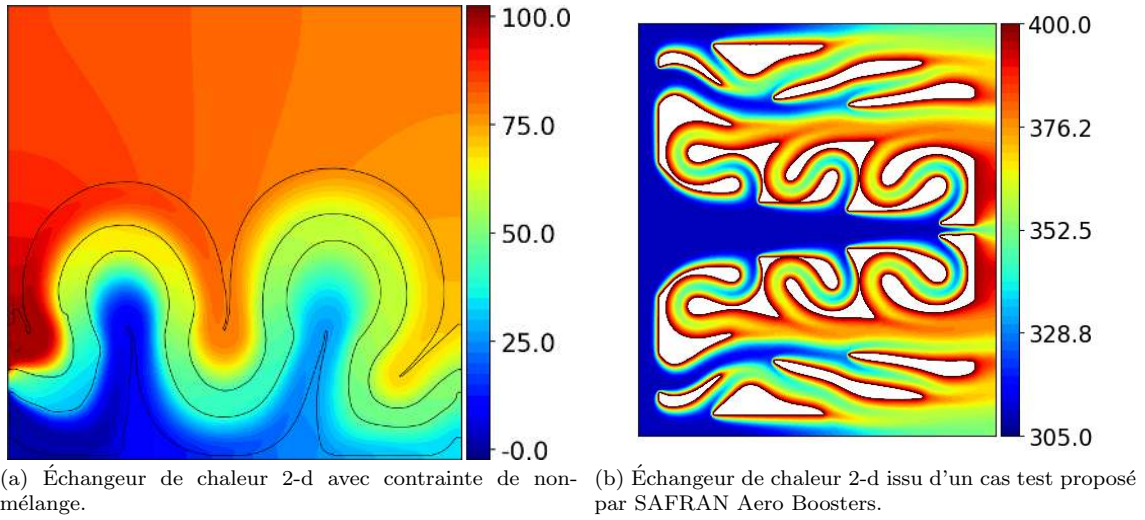


Figure 12: Deux cas tests d'optimisations considérés dans le chapitre 5.

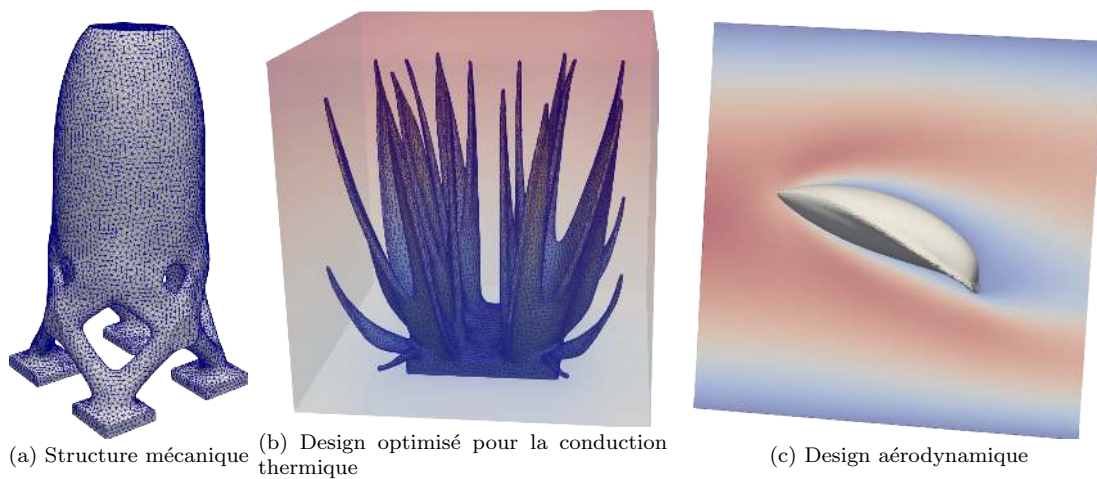


Figure 13: Quelques formes optimisées 3-d obtenues dans le chapitre 6.

limitant les efforts de frictions. Nous traitons enfin un cas test plus complexe d'un système en interaction fluide-structure. Grâce aux travaux récents de Moulin et. al. [238], nous avons été en mesure de simuler la mécanique des fluides de systèmes discrétisés par des maillages présentant jusqu'à deux millions d'éléments.

### Chapitre 7 : équations homogénéisées d'ordres élevés pour les problèmes perforés, vers l'optimisation topologique des fluides par la méthode d'homogénéisation

Ce dernier chapitre constitue une ouverture de la thèse vers l'utilisation de la méthode d'homogénéisation pour l'optimisation topologique de systèmes multiphysiques comportant des phases fluides.

Il est établi qu'en général, les solutions des problèmes d'optimisation de formes sont des matériaux composites. La théorie de l'homogénéisation permet de caractériser la physique effective des matériaux obtenus par des mélanges arbitrairement complexes de deux phases homogènes. Dans le contexte de la conception optimale des structures, celle-ci donne a donné lieu à des méthodes efficaces pour générer des formes approchant les microstructure optimales [27, 254] (telle que celle illustrée sur la Figure 6). Pour des applications en mécanique des fluides, la théorie de l'homogénéisation paraît *a priori* plus difficilement applicable : la littérature classique [103, 11, 15, 265] identifie trois régimes possibles de modèles homogénéisés, déterminés par le rapport d'échelles entre la taille et l'écartement des obstacles périodiques considérés. Aujourd'hui les méthodes d'optimisation topologique par densité utilisent quasi-exclusivement le modèle de Brinkman, qui se trouve être l'un des trois régimes homogénéisés possibles.

Dans cette partie, nous proposons des équations homogénéisées d'ordres élevés pour le système des équations de Stokes en milieu poreux qui permettent de capturer ces trois régimes simultanément. Ces modèles homogénéisés pourraient être utilisés pour l'optimisation topologique de systèmes fluides grâce à l'approche numérique proposée par [254, 27]. Notre motivation originale provient de la complexité des échangeurs thermiques industriels, dont la forme combine une silhouette macroscopique avec des motifs "microscopiques" en fines ailettes périodiques. L'optimisation de formes "classique" reposant par exemple sur la méthode de Hadamard permet de déterminer des formes optimisées *macroscopiques*, mais elle n'est pas très bien adaptée au cas où les géométries optimales sont caractérisées par de tels motifs multi-échelles.

L'un de nos résultats principaux montre que les trois régimes homogénéisés classiques peuvent être capturés par une unique équation homogénéisée d'ordre élevé : chacun des trois modèles est obtenu en passant à la limite de faible fraction volumique pour des régimes particuliers de taille d'obstacle. Plus généralement, des équations homogénéisées bien posées sont déterminées à tout ordre grâce à une méthode inspirée des travaux de Bakhvalov et Panasenko [53], Smyshlyaev et Cherednichenko [287] et Allaire, Lamacz et Rauch [33].

D'un point de vue pédagogique, nous calculons ces équations homogénéisées pour plusieurs problèmes elliptiques perforés successifs (dont les solutions s'annulent sur les petits obstacles distribués périodiquement) présentant un ordre de difficulté croissant. Nous considérons tout d'abord le cas du problème (scalaire) de Poisson perforé, puis nous étendons nos résultats au cadre vectoriel par la considération du problème analogue en élasticité. Nous étudions enfin le problème de Stokes pour lequel le terme de pression requiert un travail additionnel. Une caractéristique surprenante est l'apparition d'opérateurs *d'ordre impair* très étranges dans les équations homogénéisées d'ordres élevés : ce résultat n'est pas standard et ne semble pas avoir été observé dans la littérature classique proposant des termes correcteurs pour ces modèles.

Ce chapitre final ne contient pas de résultats numériques ; ceux-ci feront l'objet de futurs travaux.



## INTRODUCTION – SUMMARY OF THE THESIS

This work is devoted to shape and topology optimization of multiphysics systems motivated by applications of the aeronautic industry. Due to several industrial constraints that are detailed further on, we primarily consider the method of Hadamard for such purposes, although topology optimization by the homogenization method is also investigated in a final independent chapter.

Shape optimization is the mathematical art of generating shapes that best fulfill a proposed objective. Industrial applications are generally concerned with the issue of determining physical systems that achieve optimal performance. What is the shape of the most rigid structure for a given prescribed amount of material? What is the most aerodynamic airfoil generating a prescribed lift force? Which distribution of fluid pipes does achieve the best heat transfer? How to design the shape of a ship hull? These are very classical questions of industrial interest.

### 1 FROM SHAPE TO TOPOLOGY OPTIMIZATION: HADAMARD’S, DENSITY AND HOMOGENIZATION METHODS

The method of Hadamard can be traced back to the 1908 seminal memoir [180]. It consists in computing the sensitivity of the problem to small deformations of the boundary (Figure 14), which allows to optimize shapes by finding particular deformations which gradually yield better designs. With the advent of

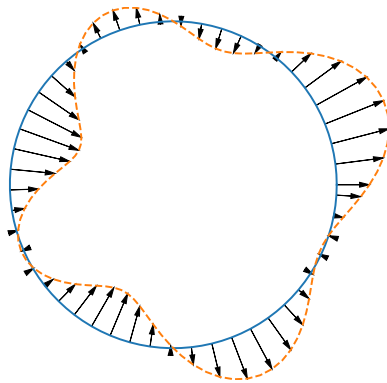


Figure 14: The boundary variation method of Hadamard: evaluating the sensitivity of the problem with respect to small shape deformations.

computers enabling to simulate numerically larger and larger physical systems, the method started to be the object of many developments in the 1970s with the seminal works of C ea [86, 87], Murat, Simon [241, 242], Dervieux and Palmerio [251, 252] and Pironneau [259, 260]. These authors were among the first to develop numerical and theoretical frameworks for solving shape optimization problems constrained by partial differential equations (PDEs). This very soon raised a growing interest in the industry with the publication of many works devoted to a variety of physical applications, dealing for instance with optimal design in mechanics [290, 182, 57], shell theory [89], heat conduction problems [121, 83, 91] and fluid mechanics [90, 172, 191, 192, 226, 123].

Since then, the field has kept improving both on theoretical [78, 79, 186] and computational aspects [179, 229]. The first numerical algorithms were based on mesh deformations [260, 291, 34]: a proposed initial guess is discretized by mean of a mesh which is deformed to a better shape (Figure 15). A major breakthrough arose in the beginning of the 2000s with the introduction of level set methods in shape optimization [32, 311, 250]: the shape to be optimized is not discretized explicitly by a mesh but rather described implicitly as the negative value set of a *level-set* function and evolved on a fixed mesh; this allows to handle complex topological changes of the optimized shape such as holes merging or boundaries collapsing (Figure 16). This ability to handle topological changes is very much desired in order to generate optimal designs without any *a priori* on the final shape, such as the number or the location of holes (of holes and arches in 3-d); for this reason the *level set method* used in conjunction with the method of Hadamard is often referred to as a *topology optimization* method, while the one based on mesh deformation is called *geometric optimization* [17]. A recent trend is to incorporate remeshing into shape and topology optimization in order to combine advantages of both methods [24, 107]: it allows to keep an explicit discretization of shapes throughout the optimization process while still allowing for topological changes (Figure 17).

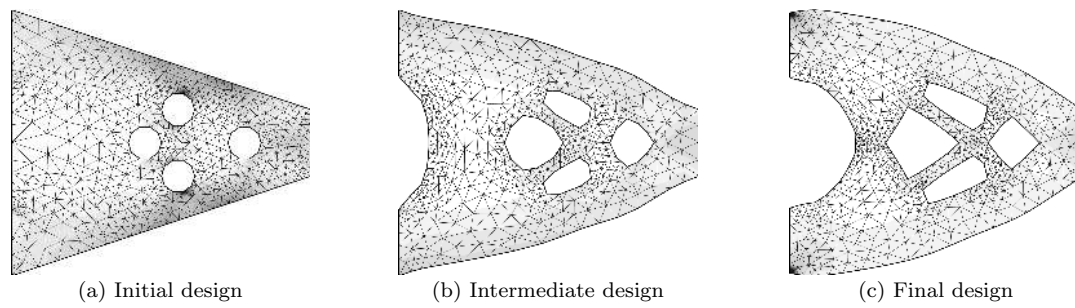


Figure 15: Hadamard's method of boundary variation implemented for the optimization of the shape of a 2-d cantilever beam with the method of mesh deformation (figures from [17]). Topological changes are difficult to handle: the number of holes of the initial and final design is unchanged.

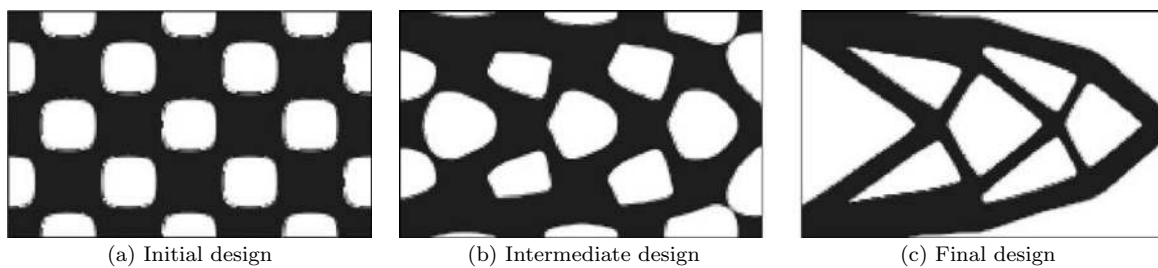


Figure 16: Hadamard's method of boundary variation implemented for the optimization of the shape of a 2-d cantilever beam with the level set method (figures from [31]). Topological changes are handled: some holes have merged from the initial to the final design.

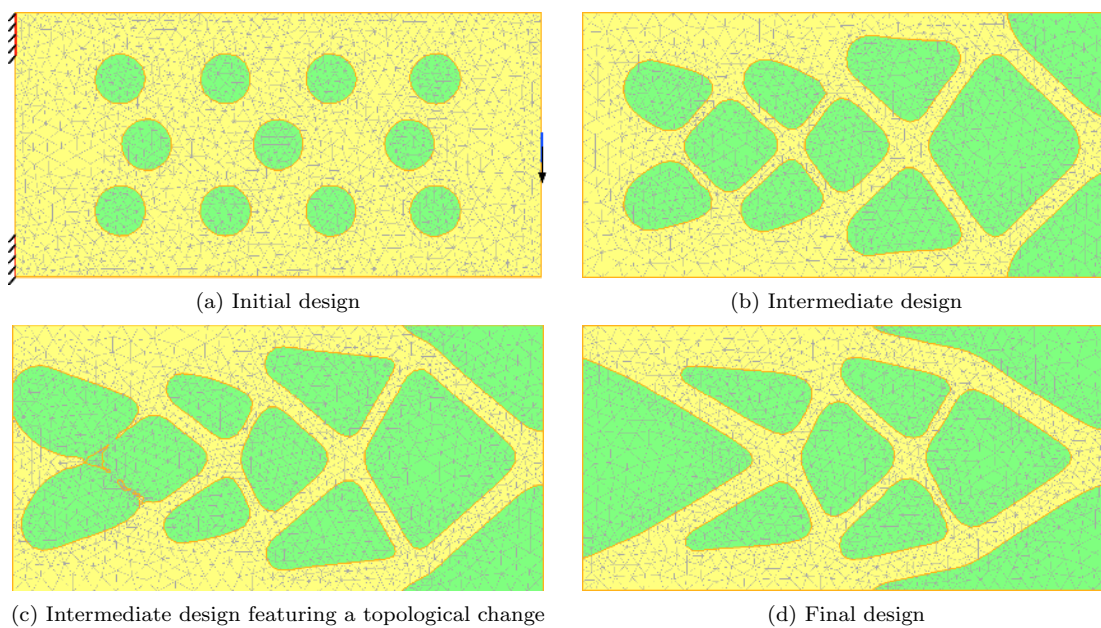


Figure 17: Hadamard's method of boundary variation implemented for the optimization of the shape of a 2-d cantilever beam with a level-set based mesh evolution method (figures from [107]). Topological changes are handled while keeping an explicit, conformal discretization of shapes.

Topology optimization is actually a rather large field which includes many other methods, such as the topological gradient method [85, 145, 165, 289, 40], level-set methods not based on the method of Hadamard [133, 316], phase field methods [73, 299], evolutionary topology optimization [193], density [66] and homogenization based methods [18, 93] (see [284] for a more complete review). Density based methods do not optimize the explicit location of the shape boundary but rather a density function  $\rho$  which represents the presence ( $\rho = 1$ ) or the absence ( $\rho = 0$ ), or intermediate states ( $0 < \rho < 1$ ) of material; some penalization scheme is used in order to obtain convergence towards true shapes, i.e. black and white designs (Figure 18). Homogenization methods replace the issue of finding an optimal shape with the one of finding an optimal composite material characterized by a varying microstructure. A first step consists in optimizing the parameters of the microstructure such as material density or cell orientation fields. In a second step (if no penalization scheme is used), these fields are interpreted in order to determine a shape approximating in some sense the optimal microstructure (Figure 19) [254, 27, 166, 177].

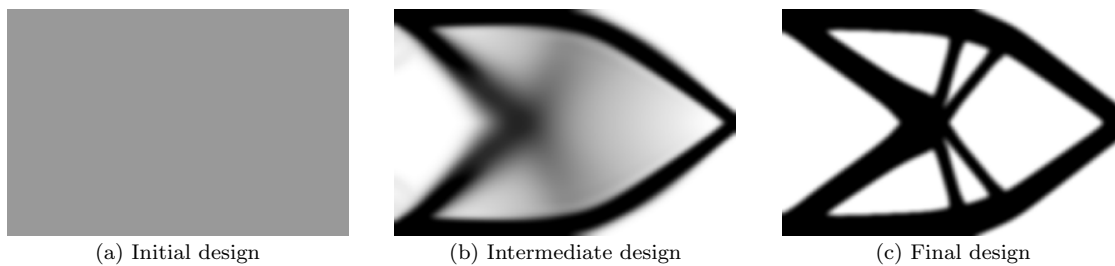


Figure 18: Topology optimization of a 2-d cantilever beam by a density method (here the so-called SIMP—Solid Isotropic Material Penalization—method). Figure obtained from the source code described in [43].

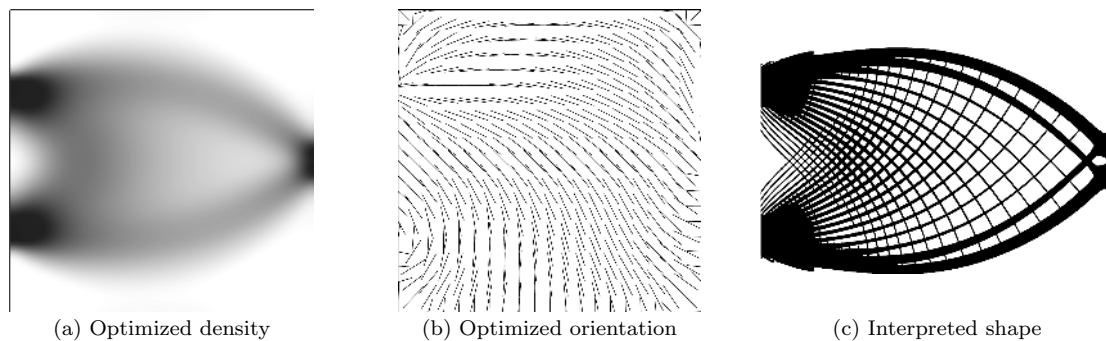


Figure 19: Topology optimization of a 2-d cantilever beam by a homogenization method. Figure from [166].

## 2 SOME CURRENT CHALLENGES IN SHAPE AND TOPOLOGY OPTIMIZATION FOR THE AERONAUTIC INDUSTRY

Today, engineers devise industrial systems most often with the assistance of Computer Aided Design (CAD) based geometry optimization software programs. Industries have been very much dependent on CAD formats because of their full compatibility with all stages of the design process, from physical numerical simulations on commercial software to actual manufacturing by automated machines. An *a priori* guess topology for the system to devise is first proposed by engineers. It is characterized for instance by a set of Bézier control points, spline surfaces parameters, or more general geometric parameters. Parametric optimization with respect to CAD parameters is then performed in order to improve the proposed design [181]. Usually, analytic shape sensitivities are not available; these are sometimes obtained by automatic differentiation (a very popular tool in aerodynamic design [94, 139, 130]) [235, 288].

However, the optimization is achieved in many cases by exploring a sufficiently large subset of the design space (e.g. with the commercial software Optimus [245]), see e.g. [142, 286]. Since these methods heavily rely on a parameterization choice of the shape geometry, they usually yield very small design updates of the initially proposed geometry (Figure 20). Such is fine for industrial applications to the extent that these small modifications may yield substantial gains of performance. However, it is easily

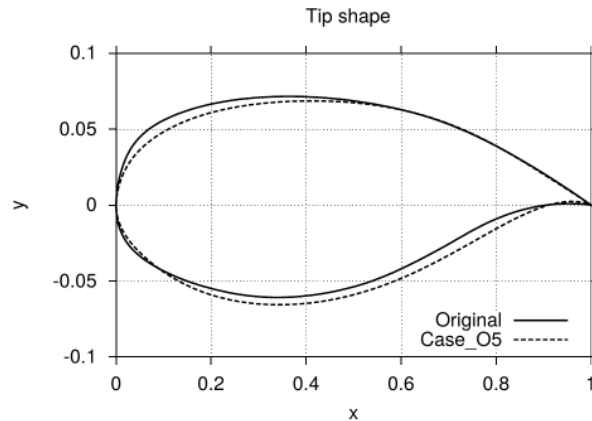


Figure 37. Tip wing section. Original vs optimized ones.

Figure 20: Lift–Drag optimization of a CAD–based airfoil design. Figure from [142].

imagined that even better performance could be obtained thanks to topology optimization since it could allow to seek very innovative new designs among much more unconstrained sets of shapes.

The rise of additive manufacturing since the 1990s has enabled the industry to be capable of fabricating more and more complex designs hardly described by CAD parameterizations, which feeds today a renewed enthusiasm for topology optimization. Many applications receiving currently a generous amount of effort are issued from the aeronautic industry and constitute some of the long term motivations at the origin of this work. One of the key challenges to overcome in order to make topology optimization applicable to aeronautic systems is the need for tackling inherent multiphysics aspects: coupled fluid, thermal, and mechanical constraints must very often be accounted for simultaneously when designing aircraft engines components.

A very representative issue drawing currently a substantial amount of attention in the topology optimization community lies in the design of heat exchangers [255, 258, 189, 271, 275]; these are devices used to cool down hot engine fluids by conveying them in the vicinity of some refrigerating gas or liquid. Industrial heat exchangers usually include many tubes and fins shaped in order to maximize the exchange surface area between hot and cold phases (Figure 21). Naturally, various additional multiphysics design constraints come into play, such as the need for controlling the loss of pressure induced by the system on the input fluid, or the mechanical resistance of the whole structure to the elevated thermal load.

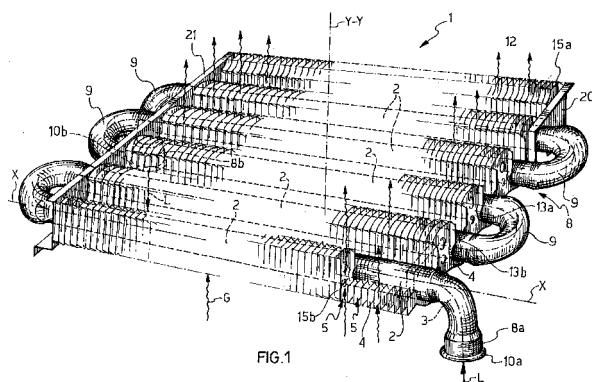


Figure 21: Industrial gas–liquid heat exchanger design. Figure from [169].

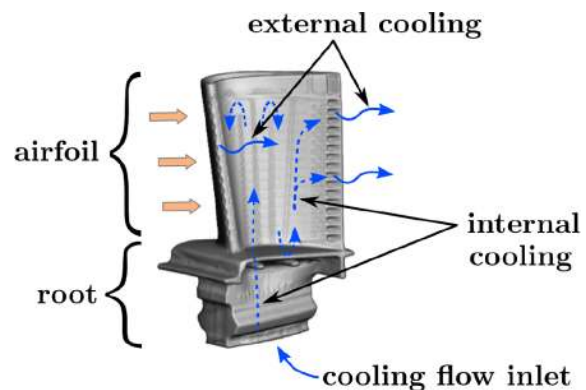


Figure 22: Turbine blade cooling system featuring internal cooling channels. Figure from [162].

It is suspected that many other components of aircraft engines could benefit as well from shape and topology optimization, such as the internal cooling channels system of turbine blades [62, 162] (Figure 22).

### 3 SCOPE OF THE THESIS

To date, the method of Hadamard used as a topology optimization method (in which the geometry of the optimized design is allowed to deform freely) is not yet completely mature for true industrial applications. Most of the test cases featured in the literature are rather exclusively concerned with linear elasticity [25, 108] or heat conduction [9, 325], with very few large scale 3-d test cases [202]. It is only relatively recently that the method is being applied on more complicated physics such as, elastoplasticity [230] or electromagnetism [160, 213]. This trend is also visible in the the density based topology optimization community which also features an increasing number of works on multidisciplinary applications [116, 129, 128, 318, 319, 240, 288].

The objective of this work is to develop theory and methodologies for shape and topology optimization of multiphysics systems, keeping in view longer term industrial requirements. The content of this thesis is mainly concerned with the boundary variation method of Hadamard, but it also includes an opening [chapter 7](#) investigating mathematical aspects of the homogenization method for future fluid applications.

By specific “industrial requirements”, it is meant several identified major industrial needs that have guided our research and which are described in the next five paragraphs. Detailed contributions of the thesis are then summarized chapter by chapter in the next section.

#### Multiphysics systems

As motivated by the industrial applications of the previous section, there is an increasing demand for topology optimization of systems involving several interacting physics. Our study shall focus on systems featuring coupled fluid, thermal and mechanical properties. Mathematically, these are characterized by a set of physical variables denoted as follows in this entire manuscript:

- $\mathbf{v}$  and  $p$  for the velocity and pressure field associated to one or more fluid phases flowing through the system;
- $T$  for the temperature field in solid and fluid phases;
- $\mathbf{u}$  for the elastic displacement of solid mechanical structures.

These variables are mathematically determined as the solutions of a set of Partial Differential Equations (PDEs), which are themselves derived from physical modelling choices. An academic setting for the characterization of  $\mathbf{v}, p, T$  and  $\mathbf{u}$ , sufficiently generic for our purposes, shall be described in [chapter 2](#).

A significant motivation of the present work lies in the design of heat exchangers. A 2-d case study for the optimization of two different models of heat exchangers is proposed in [chapter 5](#).

#### Constrained optimization

Most industrial systems feature a variety of load specifications that must be satisfied in realistic conditions of use. For instance, the overall mechanical stress or the temperature of a solid structure may be required to remain below a given bound in order to avoid premature fatigue. In other words, a design optimization problem reduces to determine the shape of a system that achieves the best performance subject to a given set of physical constraints.

Such problems can generically be modeled as mathematical programs of the form

$$\begin{aligned} \min_{\Gamma} \quad & J(\Gamma, \mathbf{v}(\Gamma), p(\Gamma), T(\Gamma), \mathbf{u}(\Gamma)) \\ \text{s.t.} \quad & \begin{cases} g_i(\Gamma, \mathbf{v}(\Gamma), p(\Gamma), T(\Gamma), \mathbf{u}(\Gamma)) = 0, & 1 \leq i \leq p, \\ h_j(\Gamma, \mathbf{v}(\Gamma), p(\Gamma), T(\Gamma), \mathbf{u}(\Gamma)) \leq 0, & 1 \leq j \leq q, \end{cases} \end{aligned} \quad (1)$$

where  $\Gamma$  denotes the boundary of the system to be optimized, typically the interface between a solid and a fluid structure (this may also include the case of a single phase).  $J$  refers to a given cost function which quantifies the system performance and which is to be minimized. Classical cost functions commonly encountered in shape optimization include the compliance of a mechanical structure, the drag force induced by an airfoil, or the heat stored into a thermal system. Equality and inequality constraints  $g_i$  and  $h_j$  model physical load specifications; they refer to target values some physical quantities need to reach (e.g. a desired location for the center of mass of the structure) or should not exceed (e.g. an upper bound limit for the overall temperature or a minimum feature size on the geometry). Both

objective function  $J$  and constraints  $g_i, h_j$  depend both on the shape  $\Gamma$  and on the physical variables  $\mathbf{v}(\Gamma), p(\Gamma), \mathbf{u}(\Gamma), T(\Gamma)$ , which depend themselves on  $\Gamma$  through physical state equations.

In density methods, the design variable  $\rho$  belongs to some euclidean space  $\mathbb{R}^n$  after discretization, which makes mathematical program of the form of (1) solvable with classical first order optimization methods [244, 298, 310]. In the context of the method of Hadamard, the optimization variable  $\Gamma$  does not belong to any euclidean space but rather a kind of infinite dimensional manifold which makes the resolution of (1) more delicate. Most available works using the method of Hadamard convert (1) into an unconstrained minimization problem by adding penalization terms to the objective function  $J$  e.g. as in the Augmented Lagrangian Method [234, 107]. More complex methods, such as Sequential Linear Programming (SLP) or the Method of Feasible Directions (MFD) have been used with some success in few published works [135, 150]. However, these methods still present some unsatisfying qualities for industrial goals because they heavily rely on the selection of case dependent meta-parameters that can be difficult to tune, especially in the presence of a large number of constraints.

In chapter 3, a novel gradient flow type algorithm is described for the resolution of constrained optimization problems adapted to the context of the method of Hadamard.

### Non-intrusiveness

For long term industrial applications, the shape optimization process should be fully compatible with the whole system conception chain. For instance, topology optimization tools should evaluate state variables by calling the same physical solvers used for industrial validation. Many topology optimization methods, including density based and level set methods, do not fulfill this condition, because they rely on artificial modifications of the state equations for pure numerical reasons.

Such need constituted one of the motivations for the systematic use of the level-set based mesh evolution algorithm of [25] (illustrated on Figure 17) in all our numerical topology optimization test cases. The main principles of the method are recalled in chapter 1. One of its main advantages is to preserve an explicit mesh discretization of the optimized shape; this allows to solve for state equations without any modification or approximation of the original physics, which could in principle be done by an external solver. For our remeshing purposes, we relied on the open source software `mmg` [108].

An additional non-intrusiveness constraint which retained our attention requires that only a strict minimal amount of information should be expected from an external user wishing to solve topology optimization test cases. In particular, the user should be allowed to specify arbitrary objective and constraints without being required the (rather technical) knowledge of the full derivation of Hadamard's shape sensitivities. In chapter 2, we provide mathematical expressions of shape sensitivities of *arbitrary* functionals which can be numerically automatically assembled from the sole knowledge of their partial derivatives.

### Geometric distance constraints

Taking into account geometric distance constraints is a classical industrial need in topology optimization due to the finite precision of the manufacturing process. These can for instance manifest into minimum thickness or overhang constraints. They also occur in heat exchanger applications; in chapter 5, we shall see that a non-mixing condition between two different fluid phases can also be formulated as a geometric distance constraint.

Many works have been devoted to the integration of such constraints in topology optimization, see [234], Chapter 3, for a review. In the context of the method of Hadamard, geometric constraints can be taken into account thanks to formulations based on the signed distance function [30]. However the practical implementation of shape sensitivities is highly delicate and depends very much on the type of shape discretization used. In chapter 4, a new variational method is proposed for the numerical evaluation of shape sensitivities of geometric constraints, which can conveniently be implemented in any finite element setting.

### Large scale problems

One last significant industrial requirement concerns the application of topology optimization to highly resolved 3-d physical systems. The physical state equations must be solved about a hundred times during a typical optimization process which is very computationally demanding for large scales problems, all the more when fluid mechanics is involved. This issue has been taken into account in various works relying on density based methods [1, 10, 71, 232, 39] and is starting to be addressed as well in evolutionary topology

optimization [227] and in level-set methods [202]. [chapter 6](#) presents and discusses the implementation in FreeFEM [183] of moderately large 3-d numerical applications (involving up to 2 million degrees of freedom) exploiting parallel computing, preconditioned iterative solvers and recent high performance domain decomposition methods [238].

Finally, industrial heat exchanger designs such as the one illustrated on [Figure 21](#) exemplify the need for generating optimized multiphysics designs featuring fine multiscale structures. In [chapter 7](#), we derive higher order homogenized models for fluid systems which could allow for the future applicability of shape optimization by the homogenization method for such applications.

## 4 SUMMARY OF CHAPTERS

### Chapter 1: Shape optimization based on Hadamard’s boundary variation method

This preliminary chapter introduces background context and material on design optimization with the method of Hadamard. We outline well-established ingredients required for an implementation workflow of topology optimization using the level-set based mesh evolution method of [25]. We expose a summary of the method of Hadamard including classical mathematical techniques that allow to calculate the shape sensitivities of PDE dependent functionals. We mention the use of gradient based optimization adapted to the infinite dimensional context. We discuss the computation of the signed distance function and detail its mathematical properties. Finally we summarize the domain evolution method of [25] based on level set advection and a remeshing step.

### Chapter 2: Hadamard’s shape derivatives for a coupled thermal fluid structure problem

This chapter introduces a simplified thermal fluid-structure model which shall be of interest throughout this whole work: the variables  $\mathbf{v}, p, T, \mathbf{u}$  are the solutions of three state equations which are weakly coupled in the sense that each of them can be solved successively and independently. This model is based on the steady state incompressible Navier-Stokes equations for the fluid velocity and pressure, convection diffusion for the temperature field, and linear thermo-elasticity for the elastic displacement. It includes previous literature studies where only one of the physics is active, while allowing for new coupled test cases such as fluid-structure interaction in a small deformation regime or convective heat transfer.

A significant novelty of this part is the derivation of Hadamard’s shape derivatives formulas for *arbitrary* objective functionals, which we detail and compare to more classical techniques in details. In the literature concerned with the method of Hadamard, hypotheses on the form of the shape functional are usually assumed: this imposes to redo the full derivation of shape derivatives for each change of objective function or constraint, which is a burden for practical implementation in industrial software. Our formulas allow to automate the computation of shape derivatives: these are assembled after the resolution of adjoint equations from the sole knowledge of the partial derivatives of the objective functional, an information which can be in principle provided easily by an external user.

Another contribution is to include in our analysis a type of fluid-structure interaction problem. A surprising property lies in the coupling of the adjoint system involved in the computation of shape sensitivities: the equality of normal stress constraints at the fluid–structure interface turns into the matching of the respective adjoint variables values on this boundary.

The validity of our formulas is finally verified numerically on several 2-d test cases activating either one, two, or three physics simultaneously.

Most of the content of this chapter has been published in [153]:

F. FEPPON, G. ALLAIRE, F. BORDEU, J. CORTIAL, AND C. DAPOGNY, *Shape optimization of a coupled thermal fluid–structure problem in a level set mesh evolution framework*, SeMA Journal, (2019), pp. 1–46.

Let us mention, however, that this chapter includes many more numerical test cases, and that some of the test cases of the published work [153] have been improved.

### Chapter 3: Null space gradient flows for constrained shape optimization

The physical model and its shape optimization context being introduced, this chapter describes the optimization algorithm developed specifically for all our numerical test cases. We propose an optimization scheme that can be interpreted as the discretization of a gradient flow able to see equality and inequality constraints (a generalization of the dynamical system approaches of [317, 300]). The gradient flow is called

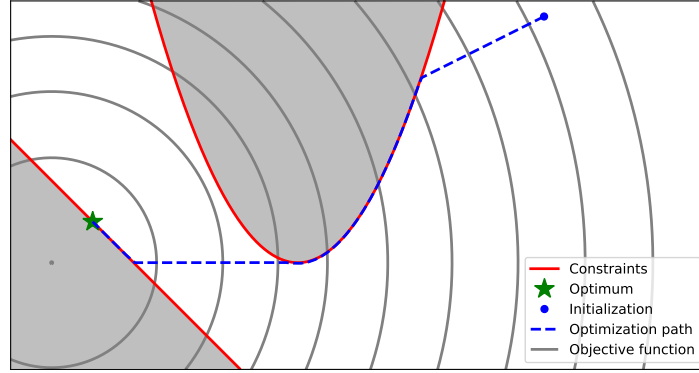


Figure 23: Illustration of our null space gradient flow for constrained optimization. Trajectories always follow the best feasible descent direction.

*null space* because the descent direction is decomposed into a *null space* step decreasing the objective function tangentially to the constraints, and a *range space* step orthogonal to the null space step and which gradually carries the optimization path towards the feasible domain. Optimization trajectories are guided by the resolution of dual quadratic subproblems which indicate when to come back into the interior of the optimization domain (Figure 23). One feature characterizing the algorithm is the ability to decrease objective function values while maintaining constraints satisfied.

Convergence properties of the continuous trajectories are established; they guarantee the improvement of optimized solutions until a local minimum is found up to the selection of a sufficiently small discretization step. The robustness and efficiency of our method for shape optimization is demonstrated numerically on multiple load bridge test cases featuring 10 constraints.

Most of the content of this chapter has been the object of the submitted preprint [155]:

F. FEPPON, G. ALLAIRE, AND C. DAPOGNY, *Null space gradient flows for constrained optimization with applications to shape optimization*, submitted, (2019).

The chapter contains in addition several numerical comparisons with classical algorithms on academic test cases.

#### Chapter 4: A variational method for computing shape derivatives of geometric constraints along rays

Heat exchangers applications featuring two distinct fluid phases motivate the need for taking into account non penetration constraints in shape optimization. A very natural approach consists in requiring each of the two phases to remain at a prescribed minimum distance from one another, which can be formulated mathematically conveniently by means of the signed distance function. More generally, many other geometric constraints can be formulated in this fashion, such as minimum or maximum thickness and minimum member's distance.

Previous works [30, 234] have proposed mathematical expressions for the derivatives of shape functionals depending on the signed distance function, and have demonstrated numerically their ability to enforce geometric constraints such as maximum or minimum thickness. However, their direct implementation is difficult because it requires numerical integration along the normal rays to the shape (Figure 24) and accurate estimates of shape skeleton or its principal curvatures. The implementation of these operations is notoriously difficult and very much depends on the dimension (2-d or 3-d) and on the type of shape discretization used (implicit or conforming, on a structured or unstructured mesh).

In this chapter, we demonstrate that shape derivatives of geometric constraints can in fact be obtained from a variational problem which can be solved conveniently by the finite element method: the value of integrals along the rays are retrieved from the variational solution values at boundary vertices. Our method amounts, in full generality, to compute integral quantities along the characteristic curves of a given velocity field  $\beta$  without requiring the explicit knowledge of these curves on the spatial discretization. The implementation is very much simplified because it requires neither the computations of rays nor of the shape curvatures but only that of a rather arbitrary weight vanishing approximately on the skeleton set.

The well-posedness of the proposed variational formulation is established thanks to a detailed analysis of the weighted graph space of the advection operator  $\beta \cdot \nabla$ . In the shape optimization context,  $\beta$  is



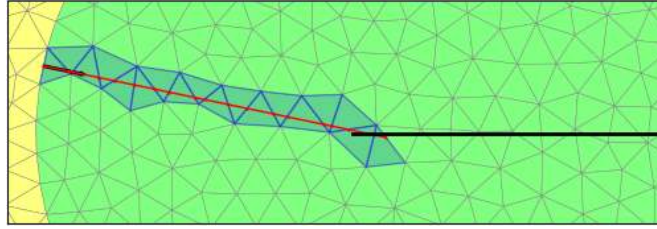


Figure 24: Classical methods for computing the shape derivatives of geometric distance constraints involve the computation of the normal rays of the shape (in red), which requires to travel along the mesh discretization.

given by the gradient of the signed distance function of the working domain, however our analysis also includes more general vector fields satisfying a suitable set of assumptions. One novelty of our approach is the ability to handle velocity fields with possibly unbounded divergence; classical works performing similar analysis [144] usually assume  $\text{div}(\beta) \in L^\infty(D)$  where  $D$  denotes the working domain, which in general does not hold for shape optimization applications.

The ease of implementation of our method is demonstrated for maximum and minimum thickness constraints in structural design: we have been able to retrieve numerical results analogous to those of previous works [30, 234] relying on explicit ray integrals.

Most of the content of this chapter is to appear in the publication [154]:

F. FEPPON, G. ALLAIRE, AND C. DAPOGNY, *A variational formulation for computing shape derivatives of geometric constraints along rays*, To appear in M2AN, (2019).

## Chapter 5: Topology optimization of 2-d heat exchangers

The material developed in the previous chapters is now applied to design optimization of 2-d heat exchangers.

In a first part, our methods are applied on the topology optimization of a liquid-liquid heat exchanger problem with the weakly coupled model introduced in chapter 2. A non-mixing constraint for two liquid phases exchanging their heat is formulated as a geometric distance constraint and treated with the variational method of the previous chapter 4 (Figure 25a).

In a second part, we present a 2-d heat exchanger case study issued from a collaboration with Safran Aero Boosters: the objective is to determine the optimal shape of transverse oil pipe cross sections refrigerated by an input cold air flow. The test case fits well in the physical setting introduced in chapter 2 up to a small change of boundary conditions for the thermal profile. The optimization problem featured a minimum thickness constraint for the oil phase and a maximum pressure loss constraint for the air phase. Our study demonstrates the ability of our shape optimization method to generate a variety of non-classical designs in a simplified setting (Figure 25b).

## Chapter 6: Towards 3-d and industrial applications: implementation recipes for a variety of numerical test cases

This chapter discusses the implementation of numerical test cases approaching more advanced industrial applications.

A first section describes succinctly programming paradigms of our topology optimization code developed in python and FreeFEM [183] in the context of this thesis. We then discuss delicate implementation details specific to the 3-d context, for instance the unavoidable use of domain decomposition techniques in conjunction with physics-dependent preconditioners.

The next section presents then a variety of new 3-d shape optimization results for different physics. Our case studies include the classical benchmark example of the cantilever beam subject to traction or torsion, optimal designs for heat conduction (Figure 13), optimal lift-drag profiles, up to more complex fluid-structure interaction test cases. Thanks to the work of Moulin et. al. [238], we have been able to solve problems featuring Navier-Stokes flows with mesh discretizations involving up to two millions of mesh elements.

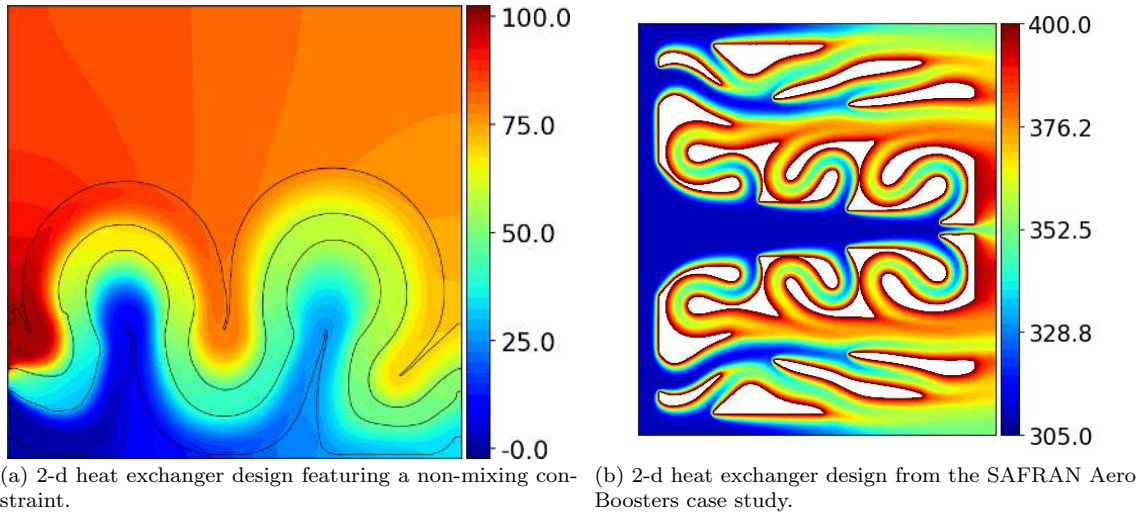


Figure 25: Two of the optimization test cases presented in [chapter 5](#).

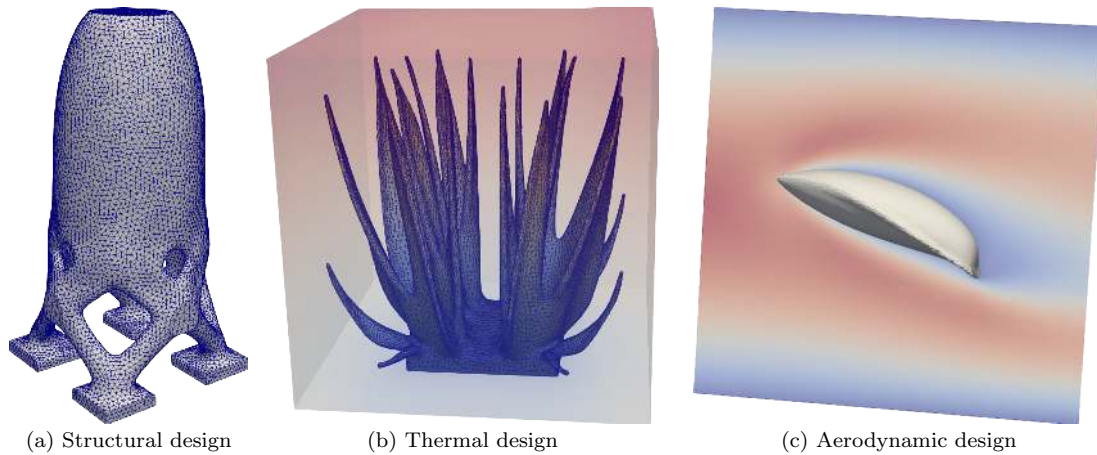


Figure 26: A few 3-d test cases obtained in [chapter 6](#).

## Chapter 7: High order homogenized equations for perforated problems, towards fluid topology optimization by the homogenization method

This final chapter constitutes an opening of the thesis towards the use of the homogenization method for topology optimization of multiphysics applications involving fluids.

It is well known that in general, optimal designs in shape optimization are composite structures. The theory of homogenization allows to characterize the effective physics of the materials obtained by arbitrarily complex phase mixtures: this enables to develop practical shape and topology optimization algorithms [18]. In the context of optimal design for the mechanics of structures, the theory of homogenization is well known and efficient methods have been developed for generating shapes approaching optimal composite micro-structures [27, 254] (such as those depicted on [Figure 19](#)). For fluid applications, the theory of homogenization is less clear, because classical literature [103, 11, 15, 265] identifies three possible homogenized regimes depending on how periodic obstacles scale within their periodic cell. To date, density based methods rely exclusively on the Brinkman model for such purposes, which turns to be only one of these possible homogenized the three.

This part of the thesis investigates high order homogenized equations for the Stokes system in a porous medium which capture all three regimes. These homogenized models which could serve for topology optimization of fluid systems by applying the approach of [254, 27]. Our motivation originates from the observation of the complexity of industrial heat exchanger designs, which combine a macroscopic shape design together with highly resolved periodic blade structures. Classical shape optimization methods such as the one developed above are able to determine *macroscopic* geometries, but are not very well suited in the case where optimal geometries exhibit such multi-scale patterns.

One of our main results shows that the classical three homogenized regimes can be captured by a single higher order homogenized equation: each of them is retrieved, in low-volume fraction limits, for particular obstacle's shape scalings. More generally, well-posed homogenized equations are derived at any order thanks to a method inspired from the work of Bakhvalov and Panasenko [53], Smyshlyaev and Cherednichenko [287] and Allaire et. al. [33].

From a pedagogical perspective, we provide derivations and a mathematical analysis of high order homogenized equations for several elliptic perforated problems (whose solutions vanish on small periodic holes) featuring an increasing order of complexity. We first consider the case of the perforated scalar Poisson problem, the extension to the vectorial case is then treated by studying the analogous perforated problem in elasticity, before turning eventually to the Stokes problem where the pressure term requires some additional work. A very striking feature is the occurrence of strange differential operators of *odd order* in these homogenized equations: this fact is not standard and does not seem to have been noticed in classical literature seeking higher order model corrections.

This final chapter does not include numerical results, which are the object of future works.



# CHAPTER 1

## SHAPE AND TOPOLOGY OPTIMIZATION BASED ON HADAMARD'S BOUNDARY VARIATION METHOD

### Contents

---

<b>1.1 Notation</b> . . . . .	<b>37</b>
<b>1.2 The boundary variation method of Hadamard for shape sensitivity of PDE constrained problems</b> . . . . .	<b>38</b>
1.2.1 Shape derivatives in the sense of Hadamard . . . . .	38
1.2.2 Shape derivatives of volume and surface functionals . . . . .	40
1.2.3 Shape derivatives of PDE constrained functionals . . . . .	43
<b>1.3 On the signed distance function and its main properties</b> . . . . .	<b>48</b>
1.3.1 Definition of the signed distance function and its first two derivatives . . . . .	48
1.3.2 Differentiability of the signed distance function with respect to the domain and shape derivatives of distance functionals . . . . .	51
1.3.3 Numerical methods for the computation of the signed distance function . . . . .	52
<b>1.4 A classical shape optimization numerical workflow using a level-set based mesh evolution method</b> . . . . .	<b>53</b>
1.4.1 Gradient based optimization in the context of the method of Hadamard: identification, regularization, and extensions of shape derivatives . . . . .	54
1.4.2 Numerical representations and updates of shapes: mesh deformations, level set methods, and level set based mesh evolution method . . . . .	55

---

This chapter is a review of classical background material on shape and topology optimization based on the method of Hadamard. The first [section 1.1](#) is a reference section for the notation that is used throughout this whole thesis. Then [section 1.2](#) summarizes common definitions and properties of shape derivatives. For later comparison with the method proposed in [chapter 2](#), classical techniques are recalled for the derivation of shape derivatives in the context of PDE constrained problem. [Section 1.3](#) is devoted to the signed distance function, a central object for taking into account geometric constraints in shape optimization. Important properties that shall be used intensively in [chapter 4](#) are referenced, including the sensitivity of the signed distance function with respect to the domain. Finally, [section 1.4](#) focuses on the practical integration of the previous ingredients into numerical algorithms. Classical implementation steps of the Hadamard's method for topology optimization are summarized. The gradient method and its adaptation to the (infinite dimensional) shape optimization context is reviewed as a preliminary to the dedicated [chapter 3](#). Finally, we discuss the important issue of numerically representing and evolving shapes. We mention mesh deformation and level set methods, and we summarize the level-set based mesh evolution algorithm of [\[24\]](#) which we used for all our numerical shape optimization test cases.

### 1.1 NOTATION

In what follows, for  $d$  an integer and a given open set  $D \subset \mathbb{R}^d$ , we denote by

$$L^p(D) := \left\{ v \text{ measurable} \mid \int_D |v|^p dx < +\infty \right\}$$

the space of  $p$ -integrable functions and by  $L^\infty(D)$  the space of bounded functions almost everywhere on  $D$ . The Sobolev space  $W^{m,p}(D)$  for  $m \in \mathbb{N}$  and  $p \in (1, +\infty)$  (see [\[76\]](#)) is the space of functions with derivatives up to the order  $k$  in  $L^p(D)$ .

$$W^{m,p}(D) := \{v \text{ measurable} \mid \forall 0 \leq k \leq m, \forall 1 \leq i_1 \leq \dots \leq i_k \leq d, \partial_{i_1 \dots i_k}^k v \in L^p(D)\}.$$

The Hilbert space  $W^{m,2}(D)$  shall also be denoted by  $H^m(D)$ .

The norm naturally associated to a Banach space  $V$  is denoted by  $\|\cdot\|_V$ , for instance

$$\forall v \in W^{m,p}(D), \|v\|_{W^{m,p}(D)} = \left( \sum_{0 \leq k \leq m} \sum_{1 \leq i_1 \dots i_k \leq d} |\partial_{i_1 \dots i_k}^k v|^p \right)^{1/p}.$$

When  $V$  is a Hilbert space, its scalar product is denoted by  $\langle \cdot, \cdot \rangle_V$ . Its natural norm is then  $\|\cdot\|_V := \langle \cdot, \cdot \rangle_V^{1/2}$ . Space of vector valued functions whose components belong to  $L^p(D)$ ,  $W^{m,p}(D)$  or  $H^m(D)$  are respectively denoted by  $L^p(D, \mathbb{R}^d)$ ,  $W^{m,p}(D, \mathbb{R}^d)$  and  $H^m(D, \mathbb{R}^d)$ . The dual space of a Banach space  $V$  is denoted by  $V^*$ , and the identity mapping by  $I : V \rightarrow V$  (the dependence with respect to  $V$  being implicit if clear from the context).

The usual euclidean norm on  $\mathbb{R}^d$  is denoted by  $\|\cdot\|$  (without index notation) and the scalar product of two vectors  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$  by  $\mathbf{a} \cdot \mathbf{b}$ . We shall also use the notation

$$\mathbf{a} \cdot M \cdot \mathbf{b} := \mathbf{a}^T M \mathbf{b}$$

for scalar products in  $\mathbb{R}^d$  featuring a product by a  $d \times d$  matrix  $M \in \mathbb{R}^{d \times d}$ . In general, elements  $\mathbf{f} \in \mathbb{R}^d$  that must be thought of as vectors are written in boldface, while a non bold notation  $x \in \mathbb{R}^d$  is preferred those that must be rather considered as points.

Finally, the gradient of a differentiable function  $f$  is the *column* vector  $\nabla f := (\partial_i f)_{1 \leq i \leq d}$ . The Jacobian matrix of a differentiable vector field  $\mathbf{f}$  is the  $d \times d$  matrix  $\nabla \mathbf{f} := (\partial_j f_i)_{1 \leq i, j \leq d}$ .

## 1.2 THE BOUNDARY VARIATION METHOD OF HADAMARD FOR SHAPE SENSITIVITY OF PDE CONSTRAINED PROBLEMS

This section reviews the classical theory of shape differentiation by the method of Hadamard for PDE constrained functionals. Definitions of shape derivatives and the Hadamard's structure theorem are recalled in [section 1.2.1](#). Reference formulas for the shape differentiation of volume and boundary integrals are summarized in [section 1.2.2](#). Finally, the classical steps involved in the calculation of shape derivatives for PDE constrained functionals are reviewed in [section 1.2.3](#). Sketch of proofs are often included for the reader's convenience, who will find more complete material in the classical textbooks [[17](#), [184](#), [291](#)].

### 1.2.1 Shape derivatives in the sense of Hadamard

Shape optimization is concerned with the problem of finding the shape of a bounded Lipschitz domain  $\Omega \subset \mathbb{R}^d$  ( $d = 2$  or  $3$  in usual applications) which minimizes a given cost function  $J$ :

$$\min_{\Omega} J(\Omega). \tag{1.2.1}$$

For the applications considered in this thesis, we shall also be interested in the slightly more general case where the shape to optimize is a Lipschitz codimension one manifold  $\Gamma$ , possibly with non empty boundary. For instance, in [chapters 2](#) and [6](#),  $\Gamma$  represents the interface between two subdomains occupied by fluid and solid phases.

The purpose of the method of Hadamard [[184](#), [17](#), [291](#)] is to introduce a notion of differentiation with respect to the position of the shape  $\Omega$  (or  $\Gamma$ ) which ultimately allows to solve [\(1.2.1\)](#) (or more complex problems such as [\(1\)](#) in the introduction) with gradient based optimization algorithms. The principle of the method is to consider shape deformations of the form

$$\Omega_{\boldsymbol{\theta}} = (I + \boldsymbol{\theta})\Omega, \text{ where } \boldsymbol{\theta} \in W^{1,\infty}(\mathbb{R}^d, \mathbb{R}^d), \|\boldsymbol{\theta}\|_{W^{1,\infty}(\mathbb{R}^d, \mathbb{R}^d)} < 1, \tag{1.2.2}$$

$$\Gamma_{\boldsymbol{\theta}} = (I + \boldsymbol{\theta})\Gamma, \text{ where } \boldsymbol{\theta} \in W^{1,\infty}(\mathbb{R}^d, \mathbb{R}^d), \|\boldsymbol{\theta}\|_{W^{1,\infty}(\mathbb{R}^d, \mathbb{R}^d)} < 1, \tag{1.2.3}$$

in either case where the shape to optimize is a bounded domain  $\Omega$  or an interface  $\Gamma$ .

The variable  $\boldsymbol{\theta}$  denotes a small vector field which moves all points of  $\Omega$  from  $x$  to the deformed location  $x + \boldsymbol{\theta}(x)$ , as is illustrated on [Figure 1.1](#) below. In order to estimate the sensitivity of the shape with respect to such deformations, it is sufficient to consider "small" displacements  $\boldsymbol{\theta}$ ; this smallness is generally measured with the norm  $\|\cdot\|_{W^{1,\infty}(\mathbb{R}^d, \mathbb{R}^d)}$  of the set  $W^{1,\infty}(\mathbb{R}^d, \mathbb{R}^d)$ , which is the space of Lipschitz bounded vector fields [[76](#)]:

$$W^{1,\infty}(\mathbb{R}^d, \mathbb{R}^d) := \{\boldsymbol{\theta} \in L^\infty(\mathbb{R}^d, \mathbb{R}^d) \mid \nabla \boldsymbol{\theta} \in L^\infty(\mathbb{R}^d, \mathbb{R}^{d \times d})\}.$$

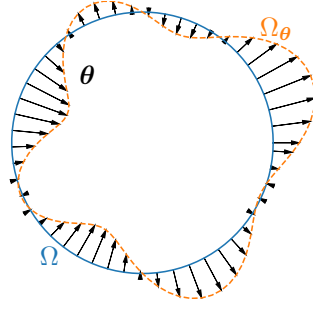


Figure 1.1: Deformation of a domain  $\Omega$  with the method of Hadamard. A small vector field  $\boldsymbol{\theta}$  is used to deform  $\Omega$  into  $\Omega_{\boldsymbol{\theta}} = (I + \boldsymbol{\theta})\Omega$ .

A classical result states that  $I + \boldsymbol{\theta}$  is a Lipschitz diffeomorphism whenever  $\boldsymbol{\theta}$  has a norm smaller than one:

**Lemma 1.1** (see [17], Lemma 6.13 p.129). *For any  $\boldsymbol{\theta} \in W^{1,\infty}(\mathbb{R}^d, \mathbb{R}^d)$  such that  $\|\boldsymbol{\theta}\|_{W^{1,\infty}(\mathbb{R}^d, \mathbb{R}^d)} < 1$ , the map  $I + \boldsymbol{\theta}$  is a bijection satisfying  $(I + \boldsymbol{\theta})^{-1} - I \in W^{1,\infty}(\mathbb{R}^d, \mathbb{R}^d)$ .*

*Sketch of proof.* Formally, the inverse map is given by

$$(I + \boldsymbol{\theta})^{-1} = \sum_{k=0}^{+\infty} (-1)^k \overbrace{\boldsymbol{\theta} \circ \dots \circ \boldsymbol{\theta}}^{k \text{ times}},$$

where the above series is convergent in the norm of  $W^{1,\infty}(\mathbb{R}^d, \mathbb{R}^d)$ . For a complete proof, see [17].  $\square$

The shape derivative of a given functional  $J(\Omega)$  is then defined by differentiation of  $J(\Omega_{\boldsymbol{\theta}})$  with respect to  $\boldsymbol{\theta}$ :

**Definition 1.1.** A shape functional  $J(\Omega)$  is said shape differentiable if the mapping

$$\begin{aligned} W^{1,\infty}(\mathbb{R}^d, \mathbb{R}^d) &\longrightarrow \mathbb{R} \\ \boldsymbol{\theta} &\longmapsto J(\Omega_{\boldsymbol{\theta}}) \end{aligned} \tag{1.2.4}$$

is Fréchet differentiable at  $\boldsymbol{\theta} = 0$ , *i.e.* if there exists a continuous linear form

$$DJ(\Omega) \in W^{1,\infty}(\mathbb{R}^d, \mathbb{R}^d)^* \tag{1.2.5}$$

such that the following asymptotics holds true:

$$J(\Omega_{\boldsymbol{\theta}}) = J(\Omega) + DJ(\Omega)(\boldsymbol{\theta}) + o(\boldsymbol{\theta}), \quad \text{where } \frac{|o(\boldsymbol{\theta})|}{\|\boldsymbol{\theta}\|_{W^{1,\infty}(\mathbb{R}^d, \mathbb{R}^d)}} \xrightarrow{\boldsymbol{\theta} \rightarrow 0} 0. \tag{1.2.6}$$

**Remark 1.1.** In case where the shape to optimize is an interface  $\Gamma$ , a functional  $J(\Gamma)$  is said shape differentiable if  $\boldsymbol{\theta} \mapsto J(\Gamma_{\boldsymbol{\theta}})$  is differentiable and the shape derivative  $DJ(\Gamma)(\boldsymbol{\theta})$  is defined analogously to (1.2.6).

**Remark 1.2.** It will be convenient to write (1.2.6) with a  $d/d\boldsymbol{\theta}$  differential notation:

$$\left. \frac{d}{d\boldsymbol{\theta}} \right|_{\boldsymbol{\theta}=0} [J(\Omega_{\boldsymbol{\theta}})](\boldsymbol{\theta}) := DJ(\Omega)(\boldsymbol{\theta}),$$

where with a little abuse of notations, we have also denoted by  $\boldsymbol{\theta}$  the direction in which  $\boldsymbol{\theta} \mapsto J(\Omega_{\boldsymbol{\theta}})$  is differentiated.

**Remark 1.3.**  $W^{1,\infty}(\mathbb{R}^d, \mathbb{R}^d)^*$  is the dual space of  $W^{1,\infty}(\mathbb{R}^d, \mathbb{R}^d)$ . The notation (1.2.5) requires consequently the existence of some constant  $C(\Omega)$  independent of  $\boldsymbol{\theta}$  such that

$$\forall \boldsymbol{\theta} \in W^{1,\infty}(\mathbb{R}^d, \mathbb{R}^d), |DJ(\Omega)(\boldsymbol{\theta})| \leq C(\Omega) \|\boldsymbol{\theta}\|_{W^{1,\infty}(\mathbb{R}^d, \mathbb{R}^d)}.$$

An important theoretical result known as Hadamard's structure theorem states that under suitable smoothness assumptions, the shape derivative  $DJ(\Omega)(\boldsymbol{\theta})$  of a functional  $J(\Omega)$  depends only on the normal trace component  $\boldsymbol{\theta} \cdot \mathbf{n}$  of the vector field  $\boldsymbol{\theta}$  deforming the bounded open set  $\Omega$ .

**Proposition 1.1** (Hadamard's structure theorem [17, 184]). *Let  $\Omega$  a smooth bounded open set of  $\mathbb{R}^d$  and  $J(\Omega)$  a shape differentiable functional. If  $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in W^{1,\infty}(\mathbb{R}^d, \mathbb{R}^d)$  are such that  $\boldsymbol{\theta}_2 - \boldsymbol{\theta}_1 \in \mathcal{C}^1(\mathbb{R}^d, \mathbb{R}^d)$  and  $\boldsymbol{\theta}_1 \cdot \mathbf{n} = \boldsymbol{\theta}_2 \cdot \mathbf{n}$  on  $\partial\Omega$ , then it holds*

$$DJ(\Omega)(\boldsymbol{\theta}_1) = DJ(\Omega)(\boldsymbol{\theta}_2).$$

### 1.2.2 Shape derivatives of volume and surface functionals

The purpose of this part is to recall classical theorems involved in shape derivative calculus of volume and surface integrals. The main tool underlying all subsequent properties is the classical change of variable formula [184, 147, 212].

**Proposition 1.2.** *If  $\Phi$  is a Lipschitz diffeomorphism of  $\mathbb{R}^d$  and  $\Omega \subset \mathbb{R}^d$  an open set, then for any  $f \in L^1(\Phi(\Omega))$ ,  $f \circ \Phi$  belongs to  $L^1(\Omega)$  and it holds*

$$\int_{\Phi(\Omega)} f dx = \int_{\Omega} f \circ \Phi \|\nabla \Phi\| dx. \quad (1.2.7)$$

We shall use as well its variant for the change of variables on co-dimension one surfaces [17, 184]:

**Proposition 1.3.** *Let  $\Gamma$  a  $\mathcal{C}^1$  codimension one surface and  $\Phi$  a  $\mathcal{C}^1$  diffeomorphism of  $\mathbb{R}^d$ . Then for any function  $f \in L^1(\Phi(\Gamma))$ , it holds  $f \circ \Phi \in L^1(\Gamma)$  and*

$$\int_{\Phi(\Gamma)} f ds = \int_{\Gamma} f \circ \Phi |\det(\nabla \Phi)| |(\nabla \Phi)^{-T} \mathbf{n}| ds, \quad (1.2.8)$$

where  $\mathbf{n}$  is any normal vector field to  $\Gamma$ .

*Sketch of proof.* Formula (1.2.8) can be deduced from the generalization of (1.2.7) to arbitrary manifolds [212]. In that case, the change of variable from  $\Gamma$  to  $\Phi(\Gamma)$  reads:

$$\int_{\Phi(\Gamma)} f ds = \int_{\Gamma} f \circ \Phi |\det(\nabla \Phi|_{\mathcal{T}(s)})| ds$$

where  $|\det(\nabla \Phi|_{\mathcal{T}(s)})|$  is the determinant of the restriction of  $\nabla \Phi$  to the tangent space  $\mathcal{T}(s)$  of  $\Gamma$  at  $s$  and onto the tangent space  $\mathcal{T}(\Phi(s))$  of  $\Phi(\Gamma)$  at  $\Phi(s)$ . It is clear that  $\mathcal{T}(\Phi(s))$  is given by

$$\mathcal{T}(\Phi(s)) = \{\nabla \Phi(s) \boldsymbol{\xi} \mid \boldsymbol{\xi} \in \mathcal{T}(s)\}.$$

Denote  $\mathbf{n}'(s)$  a normal vector to  $\mathcal{T}(\Phi(s))$ . The Jacobian matrix  $\nabla \Phi(s)$  can be written in the form of a block matrix with respect to the decomposition  $\mathbb{R}^d = \mathbb{R}\mathbf{n}(s) \oplus \mathcal{T}(s)$  onto  $\mathbb{R}^d = \mathbb{R}\mathbf{n}'(s) \oplus \mathcal{T}(\Phi(s))$ :

$$\nabla \Phi(s) = \begin{array}{c} \left[ \begin{array}{c|c} \mathbb{R}\mathbf{n}(s) & \mathcal{T}(s) \\ \hline \nabla \Phi(s) \mathbf{n}(s) & \nabla \Phi(s)|_{\mathcal{T}(s)} \end{array} \right] \begin{array}{l} \left. \vphantom{\begin{array}{c|c} \mathbb{R}\mathbf{n}(s) & \mathcal{T}(s) \\ \hline \nabla \Phi(s) \mathbf{n}(s) & \nabla \Phi(s)|_{\mathcal{T}(s)} \end{array}} \right\} \mathbb{R}\mathbf{n}'(s) \\ \left. \vphantom{\begin{array}{c|c} \mathbb{R}\mathbf{n}(s) & \mathcal{T}(s) \\ \hline \nabla \Phi(s) \mathbf{n}(s) & \nabla \Phi(s)|_{\mathcal{T}(s)} \end{array}} \right\} \mathcal{T}(\Phi(s)) \end{array}$$

This implies that the determinant of  $\nabla \Phi(s)$  is given by  $\det(\nabla \Phi(s)) = \mathbf{n}'(s) \cdot \nabla \Phi(s) \cdot \mathbf{n}(s) \det(\nabla \Phi|_{\mathcal{T}(s)})$ . Finally, the relation  $\mathbf{n}'(s) \cdot \nabla \Phi(s) \cdot \boldsymbol{\xi} = \boldsymbol{\xi} \cdot \nabla \Phi(s)^T \cdot \mathbf{n}'(s) = 0$  which holds true for any  $\boldsymbol{\xi} \in \mathcal{T}(\Phi(s))$  implies that  $\mathbf{n}(s)$  is proportional to  $\nabla \Phi(s)^T \mathbf{n}'(s)$ . This yields the following expression for  $\mathbf{n}'(s)$  (up to a sign change):

$$\mathbf{n}'(s) = \frac{\nabla \Phi^{-T} \mathbf{n}(s)}{\|\nabla \Phi^{-T}(s) \mathbf{n}(s)\|}, \quad (1.2.9)$$

from where we infer  $\langle \nabla \Phi(s) \mathbf{n}(s), \mathbf{n}'(s) \rangle = 1/\|\nabla \Phi^{-T}(s) \mathbf{n}(s)\|$  and the result.  $\square$



**Proposition 1.4.** *Let  $\Omega$  be a bounded open set of  $\mathbb{R}^d$ . For any  $f \in W^{1,1}(\mathbb{R}^d)$ , the functional  $J(\Omega)$  defined by*

$$J(\Omega) := \int_{\Omega} f(x) dx$$

*is shape differentiable and it holds*

$$DJ(\Omega)(\boldsymbol{\theta}) = \int_{\Omega} \operatorname{div}(f\boldsymbol{\theta}) dx = \int_{\Omega} (\nabla f \cdot \boldsymbol{\theta} + f \operatorname{div}(\boldsymbol{\theta})) dx, \quad \boldsymbol{\theta} \in W^{1,\infty}(\mathbb{R}^d, \mathbb{R}^d).$$

*If in addition  $\Omega$  is smooth then the above formula can be rewritten as*

$$DJ(\Omega)(\boldsymbol{\theta}) = \int_{\partial\Omega} f \boldsymbol{\theta} \cdot \mathbf{n} ds, \quad \boldsymbol{\theta} \in W^{1,\infty}(\mathbb{R}^d, \mathbb{R}^d).$$

*where  $\mathbf{n}$  denotes the outward normal to  $\Omega$ .*

*Sketch of proof.* The application of the change of variable formula yields

$$J(\Omega_{\boldsymbol{\theta}}) = \int_{(I+\boldsymbol{\theta})\Omega} f dx = \int_{\Omega} f \circ (I + \boldsymbol{\theta}) \det(I + \nabla \boldsymbol{\theta}) dx = \int_{\Omega} (\nabla f \cdot \boldsymbol{\theta} + f \operatorname{div}(\boldsymbol{\theta})) dx + o(\boldsymbol{\theta}).$$

The boundary integral formula follows by integration by part.  $\square$

The derivation of the shape derivatives of surface integrals is slightly more involved and requires some notions of differential geometry and tangential calculus [293]. For a given smooth codimension one surface  $\Gamma$  with a prescribed smooth normal vector field  $\mathbf{n}$ , we denote by  $\nabla_{\Gamma}$  the tangential gradient:

$$\forall g \in H^1(\mathbb{R}^d), \nabla_{\Gamma} g := \nabla g - (\mathbf{n} \cdot \nabla g) \mathbf{n}.$$

The tangential divergence of a vector field  $\mathbf{f} \in H^1(\mathbb{R}^d, \mathbb{R}^d)$  is defined by the trace of the linear operator  $\nabla \mathbf{f}$  restricted to the tangent spaces of  $\Gamma$ ; it reads

$$\operatorname{div}_{\Gamma}(\mathbf{f}) := \operatorname{div}(\mathbf{f}) - \mathbf{n} \cdot \nabla \mathbf{f} \cdot \mathbf{n}.$$

Let us now recall the classical definition of principal curvatures for the codimension one manifold  $\partial\Omega$  (see [293] for proofs and further material on differential geometry).

**Proposition 1.5** (Principal curvatures). *Let  $\Gamma$  be a  $C^2$  manifold and let  $\mathbf{n}$  be any differentiable unit vector field normal to  $\Gamma$ . The gradient of the normal  $\nabla \mathbf{n}$  satisfies:*

1.  $\forall y \in \Gamma, \nabla \mathbf{n}(y) \cdot \mathbf{n}(y) = 0,$
2.  $\forall y \in \Gamma, \nabla \mathbf{n}^T = \nabla \mathbf{n}.$

*In other words, for any  $y \in \Gamma$ ,  $\nabla \mathbf{n}(y)$  is a symmetric matrix which leaves the tangent space of  $\Gamma$  at  $y$  invariant. Consequently, it can be diagonalized with  $d - 1$  eigenvectors  $(\boldsymbol{\tau}_i(y))_{1 \leq i \leq d-1} \in \mathbb{R}^{d \times d}$  belonging to the tangent space at  $y$  and associated to  $d - 1$  eigenvalues  $(\kappa_i(y))_{1 \leq i \leq d-1} \in \mathbb{R}^d$ :*

$$\forall y \in \Gamma, \nabla \mathbf{n}(y) = \sum_{i=1}^{d-1} \kappa_i(y) \boldsymbol{\tau}_i(y) \boldsymbol{\tau}_i(y)^T.$$

*The real numbers  $(\kappa_i(y))_{1 \leq i \leq d-1}$  and their related basis of eigenvectors  $(\boldsymbol{\tau}_i(y))_{1 \leq i \leq d-1}$  are respectively called principal curvatures and principal directions of  $\Gamma$  at  $y$ .*

*The mean curvature of  $\Gamma$  is the real number  $\kappa(y)$  defined (up to a sign change) by*

$$\kappa(y) := \sum_{i=1}^{d-1} \kappa_i(y) = \operatorname{div}(\mathbf{n}(y)).$$

For the derivation of shape derivatives of boundary integrals, we shall also need the following identity which is a variant of the Stokes formula on manifolds:

**Lemma 1.2.** *Let  $\Gamma$  be a smooth codimension one surface of  $\mathbb{R}^d$  with boundary  $\partial\Gamma$ ,  $\mathbf{f} \in H^1(\Gamma, \mathbb{R}^d)$  and  $g \in H^1(\Gamma)$ . Then it holds*

$$\int_{\Gamma} g \operatorname{div}_{\Gamma}(\mathbf{f}) ds = \int_{\Gamma} (-\nabla_{\Gamma} g \cdot \mathbf{f} + \kappa g(\mathbf{f} \cdot \mathbf{n})) ds + \int_{\partial\Gamma} g \mathbf{f} \cdot \boldsymbol{\tau} ds,$$

where  $\boldsymbol{\tau}$  is the outward normal to  $\partial\Gamma$  tangent to  $\Gamma$ .

*Sketch of proof.* The projection of  $\mathbf{f}$  tangent to  $\Gamma$  is  $\mathbf{f}_{\Gamma} := \mathbf{f} - (\mathbf{f} \cdot \mathbf{n})\mathbf{n}$ . Stokes formula on manifolds then reads [212]

$$\int_{\Gamma} g \operatorname{div}_{\Gamma}(\mathbf{f}_{\Gamma}) ds = - \int_{\Gamma} \nabla_{\Gamma} g \cdot \mathbf{f}_{\Gamma} ds + \int_{\partial\Gamma} g \mathbf{f}_{\Gamma} \cdot \boldsymbol{\tau} dl.$$

The results follows from the following identities

$$\operatorname{div}_{\Gamma}(\mathbf{f}) = \operatorname{div}_{\Gamma}(\mathbf{f}_{\Gamma}) + \kappa \mathbf{f} \cdot \mathbf{n}, \quad (1.2.10)$$

$$\nabla_{\Gamma} g \cdot \mathbf{f}_{\Gamma} = \nabla_{\Gamma} g \cdot \mathbf{f}, \quad (1.2.11)$$

$$\mathbf{f}_{\Gamma} \cdot \boldsymbol{\tau} = \mathbf{f} \cdot \boldsymbol{\tau}. \quad (1.2.12)$$

□

The previous lemma 1.2 allows to compute shape derivatives of surface integrals.

**Proposition 1.6.** *Let  $\Gamma$  a smooth codimension one surface of  $\mathbb{R}^d$  with boundary  $\partial\Gamma$ . For any  $f \in W^{2,1}(\mathbb{R}^d)$ , the functional  $J(\Gamma)$  defined by*

$$J(\Gamma) := \int_{\Gamma} f ds$$

is shape differentiable and the shape derivative reads

$$\begin{aligned} DJ(\Gamma)(\boldsymbol{\theta}) &= \int_{\Gamma} (\operatorname{div}(f\boldsymbol{\theta}) - \mathbf{n} \cdot \nabla \boldsymbol{\theta} \cdot \mathbf{n} f) ds \\ &= \int_{\Gamma} \left( \frac{\partial f}{\partial \mathbf{n}} + \kappa f \right) (\boldsymbol{\theta} \cdot \mathbf{n}) ds + \int_{\partial\Gamma} f \boldsymbol{\theta} \cdot \boldsymbol{\tau} dl, \end{aligned}$$

where  $\boldsymbol{\tau}$  denotes the outward normal to  $\partial\Gamma$  tangent to  $\Gamma$ .

*Proof.* The first equality follows easily from the application of (1.2.8) and differentiation with respect to  $\boldsymbol{\theta}$  of

$$\int_{\Gamma_{\boldsymbol{\theta}}} f ds = \int_{\Gamma} f \circ (I + \boldsymbol{\theta}) | \det(I + \nabla \boldsymbol{\theta}) | | (I + \nabla \boldsymbol{\theta})^{-T} \mathbf{n} | ds. \quad (1.2.13)$$

The use of the previous lemma 1.2 yields the second equality because

$$\begin{aligned} \int_{\Gamma} \operatorname{div}(f\boldsymbol{\theta}) ds &= \int_{\Gamma} (\operatorname{div}_{\Gamma}(f\boldsymbol{\theta}) + \mathbf{n} \cdot \nabla(f\boldsymbol{\theta}) \cdot \mathbf{n}) ds \\ &= \int_{\Gamma} \left( \kappa f(\boldsymbol{\theta} \cdot \mathbf{n}) + \frac{\partial f}{\partial \mathbf{n}}(\boldsymbol{\theta} \cdot \mathbf{n}) + \mathbf{n} \cdot \nabla \boldsymbol{\theta} \cdot \mathbf{n} f \right) ds + \int_{\partial\Gamma} f(\boldsymbol{\theta} \cdot \boldsymbol{\tau}) dl. \end{aligned} \quad (1.2.14)$$

□

The last result of this part provides formulas for the shape derivatives of surface integrals involving the normal (which also depends on the shape). It is slightly less classical but is obtained thanks to very similar techniques.

**Proposition 1.7.** *Let  $\Gamma$  a smooth codimension one surface of  $\mathbb{R}^d$  with boundary  $\partial\Gamma$  and differentiable normal vector field  $\mathbf{n}$ . For any  $\mathbf{f} \in W^{2,1}(\mathbb{R}^d, \mathbb{R}^d)$ , the functional  $J(\Gamma)$  defined by*

$$J(\Gamma) := \int_{\Gamma} \mathbf{f} \cdot \mathbf{n} ds$$

is shape differentiable and the shape derivative reads

$$\begin{aligned} DJ(\Gamma)(\boldsymbol{\theta}) &= \int_{\Gamma} (\mathbf{n} \cdot \nabla \mathbf{f} \cdot \boldsymbol{\theta} - \mathbf{n} \cdot \nabla \boldsymbol{\theta} \cdot \mathbf{f} + (\mathbf{f} \cdot \mathbf{n}) \operatorname{div}(\boldsymbol{\theta})) ds \\ &= \int_{\Gamma} \operatorname{div}(\mathbf{f})(\boldsymbol{\theta} \cdot \mathbf{n}) ds + \int_{\partial\Gamma} [(\mathbf{f} \cdot \mathbf{n})(\boldsymbol{\theta} \cdot \boldsymbol{\tau}) - (\mathbf{f} \cdot \boldsymbol{\tau})(\boldsymbol{\theta} \cdot \mathbf{n})] ds, \end{aligned} \quad (1.2.15)$$

where  $\boldsymbol{\tau}$  denotes the outward normal to  $\partial\Gamma$  tangent to  $\Gamma$ .

*Proof.* The proposition implicitly understands that the transported unit normal  $\mathbf{n}_\theta$  to  $\Gamma_\theta = (I + \theta)\Gamma$  is given by (see formula (1.2.9))

$$\mathbf{n}_\theta \circ (I + \theta) := \frac{(I + \nabla\theta)^{-T}\mathbf{n}}{\|(I + \nabla\theta)^{-T}\mathbf{n}\|}.$$

Under this circumstance,  $J(\Gamma_\theta)$  reads

$$\begin{aligned} J(\Gamma_\theta) &= \int_{\Gamma_\theta} \mathbf{f} \cdot \mathbf{n}_\theta ds = \int_\Gamma \mathbf{f} \circ (I + \theta) \cdot \mathbf{n}_\theta \circ (I + \theta) |\det(I + \nabla\theta)| \|(I + \nabla\theta)^{-T}\mathbf{n}\| ds \\ &= \int_\Gamma \mathbf{f} \circ (I + \theta) \cdot (I + \nabla\theta)^{-T} \cdot \mathbf{n} |\det(I + \nabla\theta)| ds. \end{aligned}$$

Differentiating this expression with respect to  $\theta$  yields then the first equality. We then use the fact that  $\theta \cdot \nabla \mathbf{n} \cdot \mathbf{f} = \mathbf{f} \cdot \nabla \mathbf{n} \cdot \theta$  and  $\nabla \mathbf{n} \cdot \mathbf{n} = 0$  in order to write

$$\int_\Gamma (\mathbf{n} \cdot \nabla \mathbf{f} \cdot \theta - \mathbf{n} \cdot \nabla \theta \cdot \mathbf{f}) ds = \int_\Gamma (\mathbf{n} \cdot \nabla \mathbf{f} \cdot \mathbf{n} (\theta \cdot \mathbf{n}) - \mathbf{n} \cdot \nabla \theta \cdot \mathbf{n} (\mathbf{f} \cdot \mathbf{n}) + \theta \cdot \nabla_\Gamma (\mathbf{f} \cdot \mathbf{n}) - \mathbf{f} \cdot \nabla_\Gamma (\theta \cdot \mathbf{n})) ds. \quad (1.2.16)$$

Applying now lemma 1.2 to treat the term depending on  $\text{div}(\theta)$  in the first line of (1.2.15), we obtain:

$$\begin{aligned} \int_\Gamma (\mathbf{f} \cdot \mathbf{n}) \text{div}(\theta) ds &= \int_\Gamma [(\mathbf{f} \cdot \mathbf{n}) \text{div}_\Gamma(\theta) + (\mathbf{f} \cdot \mathbf{n}) \mathbf{n} \cdot \nabla \theta \cdot \mathbf{n}] ds \\ &= \int_\Gamma (-\nabla_\Gamma (\mathbf{f} \cdot \mathbf{n}) \cdot \theta + \kappa (\mathbf{f} \cdot \mathbf{n}) (\theta \cdot \mathbf{n}) + (\mathbf{f} \cdot \mathbf{n}) \mathbf{n} \cdot \nabla \theta \cdot \mathbf{n}) ds + \int_{\partial\Gamma} (\mathbf{f} \cdot \mathbf{n}) (\theta \cdot \boldsymbol{\tau}) ds. \end{aligned} \quad (1.2.17)$$

It remains to sum (1.2.16) and (1.2.17) to obtain:

$$DJ(\Gamma)(\theta) = \int_\Gamma (-\nabla_\Gamma (\theta \cdot \mathbf{n}) \cdot \mathbf{f} + \kappa (\mathbf{f} \cdot \mathbf{n}) (\theta \cdot \mathbf{n}) + \mathbf{n} \cdot \nabla \mathbf{f} \cdot \mathbf{n} (\theta \cdot \mathbf{n})) ds + \int_{\partial\Gamma} (\mathbf{f} \cdot \mathbf{n}) (\theta \cdot \boldsymbol{\tau}) ds \quad (1.2.18)$$

which is equivalent to the formula (1.2.15) after a second application of lemma 1.2.  $\square$

**Remark 1.4.** In the case where  $\Gamma = \partial\Omega$  is the boundary of a bounded smooth open set  $\Omega$ , the result is trivial because Stokes formula allows to write

$$\int_\Gamma \mathbf{f} \cdot \mathbf{n} ds = \int_{\partial\Omega} \mathbf{f} \cdot \mathbf{n} ds = \int_\Omega \text{div}(\mathbf{f}) dx$$

from where the shape derivative can be obtained directly by applying proposition 1.4.

### 1.2.3 Shape derivatives of PDE constrained functionals

In this last subsection we review classical methods for the calculation of shape derivatives of functionals depending on the solutions to Partial Differential Equations (PDEs). For simplicity, we consider the model problem of the Laplace equation

$$\begin{cases} -\Delta u = f & \text{in } \Omega \\ u = 0 & \text{on } \Gamma_D \\ \frac{\partial u}{\partial \mathbf{n}} = 0 & \text{on } \Gamma_N \end{cases} \quad (1.2.19)$$

set on a Lipschitz domain  $\Omega$ , where the boundary  $\partial\Omega = \Gamma_D \cup \Gamma_N$  is divided into respective Dirichlet and Neumann parts  $\Gamma_D$  and  $\Gamma_N$  (the setting is represented on Figure 1.2). Let  $V = \{v \in H^1(\Omega) \mid v = 0 \text{ on } \Gamma_D\}$ . Recall that  $u(\Omega) \in V$  is the unique function solving the variational problem

$$\text{Find } u \in V \text{ such that } \forall v \in V, \int_\Omega \nabla u \cdot \nabla v dx = \int_\Omega f v dx. \quad (1.2.20)$$

In order to make the integral term of the right hand side shape differentiable, the source term  $f$  is assumed to belong to  $H^1(\mathbb{R}^d)$ , although the problem makes sense in general for  $f \in H^{-1}(\Omega)$ .

We are concerned with the derivation of the shape derivative of a functional  $J(\Omega, u(\Omega))$  depending both on the shape  $\Omega$  and on the solution  $u(\Omega)$  to (1.2.19) (the dependence of  $u$  with respect to  $\Omega$  shall be omitted when the context is clear). Both Neumann and Dirichlet boundaries  $\Gamma_N$  and  $\Gamma_D$  are allowed

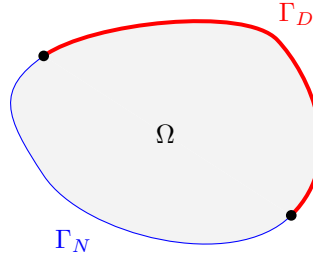


Figure 1.2: Setting for the Poisson problem of (1.2.19).

to deform. As is customary in the classical shape optimization literature, we assume (in this chapter only)  $J(\Omega, u(\Omega))$  can be written in the form of a volume integral:

$$J(\Omega, u(\Omega)) := \int_{\Omega} j(x, u(x)) dx \quad (1.2.21)$$

where  $j : \mathbb{R}^d \times \mathbb{R} \rightarrow \mathbb{R}$  is a  $C^1$  function satisfying  $|\partial_u j(x, u)| \leq c|u|$  for some constant  $c > 0$ .

As a preliminary to the derivations of chapter 2 involving more complex PDEs, the full, rigorous, calculation of shape derivative based on Lagrangian derivatives is reviewed and compared to the so-called “fast” but only formal derivation method of C ea in the next paragraph. We shall see later in chapter 2 how to adapt (and even shorten, see the proof of chapter 2, proposition 2.2) these derivations to obtain formulas for arbitrary functionals  $J(\Omega, u(\Omega))$ .

### Rigorous derivation based on Lagrangian derivatives

In this part, we follow the derivation of [184] which allows to obtain the shape sensitivity of  $J(\Omega, u(\Omega))$  in a “safe” mathematically rigorous manner. The first step is to perform the usual change of variable  $x = (I + \theta)(y)$ :

$$J(\Omega_{\theta}, u(\Omega_{\theta})) = \int_{\Omega_{\theta}} j(x, u(\Omega_{\theta})(x)) dx = \int_{\Omega} j((I + \theta)(y), u(\Omega_{\theta}) \circ (I + \theta)(y)) \det(I + \nabla \theta) dy. \quad (1.2.22)$$

This change of variable brings into play the transported solution  $u_{\theta} \in V$  defined by

$$u_{\theta} := u(\Omega_{\theta}) \circ (I + \theta).$$

Importantly, the function  $u_{\theta}$  has the nice property to be defined on the fixed space  $V$ , while  $u(\Omega_{\theta})$  is defined on the space  $H^1(\Omega_{\theta})$  which depends on  $\theta$ . The same change of variable in (1.2.20) yields a variational formulation satisfied by  $u_{\theta}$ :

Find  $u_{\theta} \in V$  such that  $\forall v \in V$ ,

$$\int_{\Omega} (I + \nabla \theta)^{-T} \nabla u_{\theta} \cdot (I + \nabla \theta)^{-T} \nabla v \det(I + \nabla \theta) dx = \int_{\Omega} f \circ (I + \theta) \det(I + \nabla \theta) v dx, \quad (1.2.23)$$

We have used the elementary but very useful lemma 1.3 below when performing change of variables involving gradients:

**Lemma 1.3.** *Let  $f \in H^1(\mathbb{R}^d)$  and  $\mathbf{f} \in H^1(\mathbb{R}^d, \mathbb{R}^d)$  be respectively scalar and vectorial functions, and  $\theta \in W^{1,\infty}(\mathbb{R}^d, \mathbb{R}^d)$  with  $\|\theta\|_{W^{1,\infty}(\mathbb{R}^d, \mathbb{R}^d)} < 1$ . It holds*

$$(\nabla f) \circ (I + \theta) = (I + \nabla \theta)^{-T} \nabla (f \circ (I + \theta)) \quad (1.2.24)$$

$$(\nabla \mathbf{f}) \circ (I + \theta) = \nabla (\mathbf{f} \circ (I + \theta)) (I + \nabla \theta). \quad (1.2.25)$$

**Remark 1.5.** The apparent asymmetry between (1.2.24) and (1.2.25) is only related to an unfortunate convention for the gradient notation  $\nabla$ . For a vector field  $\mathbf{f}$ , the notation  $\nabla \mathbf{f} := (\partial_j f_i)_{1 \leq i, j \leq d}$  for the Jacobian matrix is not consistent with the one for the gradient  $\nabla f := (\partial_i f)_{1 \leq i \leq d}$  considered as a column vector (the Jacobian matrix of  $f$  is a linear form, hence it should be a row vector). We shall however maintain this convention commonly assumed, introducing a different notation in chapter 3 in contexts where a clear distinction between gradients and differential is needed.

**Remark 1.6.** From a practical point of view, (1.2.24) and (1.2.25) imply that when differentiating a variational formulation with respect to the shape:

1. the differentiation of a scalar gradient term  $\nabla f$  yields a term  $-\nabla\theta^T\nabla f$ ;
2. the differentiation of a vectorial gradient term  $\nabla \mathbf{f}$  yields a term  $-\nabla \mathbf{f}\nabla\theta$ ;
3. the differentiation of  $dx$  yields a term  $\operatorname{div}(\theta)dx$ ;
4. test functions  $v$  are not differentiated (because  $v \circ (I + \theta)^{-1}$  can be chosen in (1.2.20)) (see [17])

These properties could be used formally e.g. for obtaining (1.2.28) below directly from the differentiation of (1.2.20) without writing actually (1.2.23).

The variational formulation (1.2.23) can be rewritten into an equation of the form  $F(\theta, u_\theta) = 0$  where  $F$  is the functional

$$\begin{aligned} F : W^{1,\infty}(\mathbb{R}^d, \mathbb{R}^d) \times V &\rightarrow V^* \\ (\theta, u) &\mapsto v \mapsto \int_{\Omega} (I + \nabla\theta)^{-T} \nabla u \cdot (I + \nabla\theta)^{-T} \nabla v \det(I + \nabla\theta) dx \\ &\quad - \int_{\Omega} f \circ (I + \theta) \det(I + \nabla\theta) v dx. \end{aligned} \quad (1.2.26)$$

It can be verified that the partial derivative  $\partial F/\partial u$  at  $(0, u(\Omega))$  is the linear operator associated to the bilinear form  $(u, v) \mapsto \int_{\Omega} \nabla u \cdot \nabla v dx$ , which is invertible from standard Lax-Milgram theory (see [184, 185]). This allows to apply the implicit function theorem ([212], Chap. 1, Theorem 5.9), which states that  $\theta \mapsto u_\theta$  is differentiable at  $\theta = 0$  (see also chapter 2, lemma 2.1 for more details). Its Fréchet derivative in a direction  $\theta$  is denoted

$$\dot{u}(\theta) := \left. \frac{du_\theta}{d\theta} \right|_{\theta=0}(\theta) \quad (1.2.27)$$

(the dependence with respect to  $\Omega$  is omitted) and it is called the *Lagrangian derivative* of  $u(\Omega)$ . A variational formulation for  $\dot{u}(\theta)$  is obtained by differentiating (1.2.23) with respect to  $\theta$ :

$$\forall v \in V, \int_{\Omega} \nabla \dot{u}(\theta) \cdot \nabla v dx = \int_{\Omega} (\nabla\theta + \nabla\theta^T - \operatorname{div}(\theta)I) \nabla u \cdot \nabla v dx + \int_{\Omega} \operatorname{div}(f\theta) v dx. \quad (1.2.28)$$

**Remark 1.7.** The Lagrangian derivative  $\dot{u}(\theta)$  is the Fréchet derivative of  $\theta \mapsto u(\Omega_\theta) \circ (I + \theta)$ . Another possible way to differentiate the domain dependent function  $u(\Omega)$  is to consider the (Gâteaux) differential of the map  $\theta \mapsto u(\Omega_\theta)(x)$  for a fixed  $x$ , which is denoted

$$\forall x \in \Omega, u'(\theta)(x) := \frac{d}{d\theta} u(\Omega_\theta)(x)$$

and is called the *Eulerian derivative* of  $u(\Omega)$ . Although this definition could seem more natural at first glance, the Eulerian derivative does not systematically exist in the same solution space than  $u(\Omega)$ . Indeed, a formal computation shows that Eulerian and Lagrangian derivatives are related by the formula

$$u'(\theta)(x) = \dot{u}(\theta)(x) - \nabla u(x) \cdot \theta(x).$$

In our case, the Lagrangian derivative  $\dot{u}(\theta)$  is an element of  $H^1(\Omega)$ , however  $\nabla u(x)$  only belongs to  $L^2(\Omega)$ , which implies that  $u'(\theta)$  has less regularity than  $\dot{u}(\theta)$ .

Coming back to  $J(\Omega_\theta, u(\Omega_\theta))$  in (1.2.22), the use of the chain rules yields then

$$J(\Omega_\theta, u(\Omega_\theta)) = J(\Omega) + \int_{\Omega} \left( \nabla_x j(x, u(x)) \cdot \theta + j(x, u(x)) \operatorname{div}(\theta) + \frac{\partial j}{\partial u}(x, u(x)) \dot{u}(\theta) \right) dx + o(\theta). \quad (1.2.29)$$

The formula is not yet satisfactory because (i) it involves the computation of  $\dot{u}(\theta)$  which requires the resolution of (1.2.28) for any vector field  $\theta$ , and (ii) it does not satisfy the Hadamard's structure theorem of proposition 1.1. The classical trick is to introduce an adjoint variable  $p \in V$  solution to the variational problem

$$\text{Find } p \in V \text{ such that } \forall v \in V, \int_{\Omega} \nabla p \cdot \nabla v dx = \int_{\Omega} \frac{\partial j}{\partial u}(x, u(x)) v dx. \quad (1.2.30)$$

The above equation means that  $p$  is the solution of the Poisson problem

$$\begin{cases} -\Delta p = \frac{\partial j}{\partial u}(x, u(x)) \text{ in } \Omega \\ p = 0 \text{ on } \Gamma_D \\ \frac{\partial p}{\partial \mathbf{n}} = 0 \text{ on } \Gamma_N. \end{cases} \quad (1.2.31)$$

Inserting  $v = \dot{u}(\boldsymbol{\theta})$  in (1.2.30) and using (1.2.28) with  $v = p$  yields indeed

$$\int_{\Omega} \frac{\partial j}{\partial u}(x, u(x)) \dot{u}(\boldsymbol{\theta}) dx = \int_{\Omega} \nabla p \cdot \nabla \dot{u}(\boldsymbol{\theta}) dx = \int_{\Omega} (\nabla \boldsymbol{\theta} + \nabla \boldsymbol{\theta}^T - \operatorname{div}(\boldsymbol{\theta})I) \nabla u \cdot \nabla p dx + \int_{\Omega} \operatorname{div}(f\boldsymbol{\theta}) p dx.$$

Replacing the above formula into (1.2.29) yields the well known volumetric expression of the shape derivative of  $J(\Omega, u(\Omega))$ .

**Proposition 1.8.** *Assume  $\Omega \subset D$  is a Lipschitz bounded open set and  $f \in H^1(\mathbb{R}^d)$ . The functional  $J(\Omega, u(\Omega))$  defined by (1.2.21) is shape differentiable and the shape derivative reads*

$$\begin{aligned} \frac{d}{d\boldsymbol{\theta}} [J(\Omega_{\boldsymbol{\theta}}, u(\Omega_{\boldsymbol{\theta}}))](\boldsymbol{\theta}) &= \int_{\Omega} \left( \frac{\partial j}{\partial x}(x, u(x)) \cdot \boldsymbol{\theta} + j(x, u(x)) \operatorname{div}(\boldsymbol{\theta}) \right) dx \\ &\quad + \int_{\Omega} [(\nabla \boldsymbol{\theta} + \nabla \boldsymbol{\theta}^T - \operatorname{div}(\boldsymbol{\theta})I) \nabla u \cdot \nabla p + p \operatorname{div}(f\boldsymbol{\theta})] dx. \end{aligned} \quad (1.2.32)$$

**Remark 1.8.** Classically and importantly, no further regularity of the state solutions  $u$  and  $p$  is required for the volume expression (1.2.32) to make sense. However, further regularity is needed for either rewriting it in the form of a boundary integral (see below).

The expression (1.2.32) is better than (1.2.29) in the sense that it does not depend on  $\dot{u}(\boldsymbol{\theta})$ , however it can still be simplified in order to obtain a boundary expression satisfying Hadamard's structure theorem. Such is obtained through an integration by part eliminating the derivatives of  $\boldsymbol{\theta}$ . Higher order derivatives of  $u$  and  $p$  are going to appear in the process: specifically, we need  $H^2$  regularity of  $u$  and  $p$ , which generally holds if  $\Omega$  is smooth except on a neighborhood of the interface  $\Gamma_D \cap \Gamma_N$  between the Neumann and Dirichlet boundary (see [176] and [114] for a rigorous treatment of such issue in shape optimization).

Therefore, we assume for now that  $\Omega$  is a smooth domain and that a neighborhood  $\omega$  of  $\Gamma_D \cap \Gamma_N \subset \omega$  is fixed:  $\boldsymbol{\theta} = 0$  on  $\omega$ . Therefore, it holds  $u, p \in H^2(\Omega \setminus \omega)$  which allows us to perform the following integration by part:

$$\begin{aligned} \int_{\Omega} [(\nabla \boldsymbol{\theta} + \nabla \boldsymbol{\theta}^T - \operatorname{div}(\boldsymbol{\theta})I) \nabla u \cdot \nabla p] dx &= - \int_{\Omega} \boldsymbol{\theta} \cdot (\operatorname{div}(\nabla p \otimes \nabla u + \nabla u \otimes \nabla p) - \nabla(\nabla u \cdot \nabla p)) dx \\ &\quad + \int_{\partial\Omega} [(\boldsymbol{\theta} \cdot \nabla p)(\nabla u \cdot \mathbf{n}) + (\boldsymbol{\theta} \cdot \nabla u)(\nabla p \cdot \mathbf{n}) - (\nabla u \cdot \nabla p)(\boldsymbol{\theta} \cdot \mathbf{n})] ds. \end{aligned} \quad (1.2.33)$$

A small calculation allows to simplify the volume integrand:

$$\begin{aligned} -[\operatorname{div}(\nabla p \otimes \nabla u + \nabla u \otimes \nabla p) - \nabla(\nabla u \cdot \nabla p)] &= -\Delta u \nabla p - \Delta p \nabla u \\ &= f \nabla p + \frac{\partial j}{\partial u}(\cdot, u) \nabla u. \end{aligned} \quad (1.2.34)$$

Then the volume integrals involving  $j$  in (1.2.32) and (1.2.34) can be rewritten as follows:

$$\begin{aligned} \int_{\Omega} \left( \frac{\partial j}{\partial x}(x, u(x)) \cdot \boldsymbol{\theta} + j(x, u(x)) \operatorname{div}(\boldsymbol{\theta}) + \nabla u \cdot \boldsymbol{\theta} \frac{\partial j}{\partial u}(x, u(x)) \right) dx \\ = \int_{\Omega} \operatorname{div}(j(x, u(x)) \boldsymbol{\theta}) dx = \int_{\partial\Omega} j(x, u(x)) \boldsymbol{\theta} \cdot \mathbf{n} ds. \end{aligned} \quad (1.2.35)$$

The volume terms involving  $p$  can also be treated in the same manner:

$$\int_{\Omega} (p \operatorname{div}(f\boldsymbol{\theta}) + \nabla p \cdot \boldsymbol{\theta} f) dx = \int_{\Omega} \operatorname{div}(fp\boldsymbol{\theta}) dx = \int_{\partial\Omega} fp\boldsymbol{\theta} \cdot \mathbf{n} ds.$$

Finally, using the boundary conditions satisfied by  $u$  and  $p$ , the surface integral term of (1.2.33) can be rewritten

$$\begin{aligned} \int_{\partial\Omega} [(\boldsymbol{\theta} \cdot \nabla p)(\nabla u \cdot \mathbf{n}) + (\boldsymbol{\theta} \cdot \nabla u)(\nabla p \cdot \mathbf{n}) - (\nabla u \cdot \nabla p)(\boldsymbol{\theta} \cdot \mathbf{n})] ds \\ = \int_{\Gamma_D} \frac{\partial u}{\partial \mathbf{n}} \frac{\partial p}{\partial \mathbf{n}} \boldsymbol{\theta} \cdot \mathbf{n} ds - \int_{\Gamma_N} \nabla u \cdot \nabla p (\boldsymbol{\theta} \cdot \mathbf{n}) ds. \end{aligned}$$

All in all, we have obtained

**Proposition 1.9.** *Assume  $\Omega$  is smooth and  $f \in H^1(\mathbb{R}^d)$ . If  $\boldsymbol{\theta} = 0$  on a neighborhood of  $\Gamma_D \cap \Gamma_N$ , then the shape derivative of  $J(\Omega, u(\Omega))$  given by (1.2.32) rewrites as a boundary integral involving only the normal trace component  $\boldsymbol{\theta} \cdot \mathbf{n}$  of  $\boldsymbol{\theta}$ :*

$$\frac{d}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=0} [J(\Omega_{\boldsymbol{\theta}}, u(\Omega_{\boldsymbol{\theta}}))](\boldsymbol{\theta}) = \int_{\partial\Omega} (j(x, u(x)) + fp) \boldsymbol{\theta} \cdot \mathbf{n} ds + \int_{\Gamma_D} \frac{\partial u}{\partial \mathbf{n}} \frac{\partial p}{\partial \mathbf{n}} \boldsymbol{\theta} \cdot \mathbf{n} ds - \int_{\Gamma_N} \nabla u \cdot \nabla p (\boldsymbol{\theta} \cdot \mathbf{n}) ds. \quad (1.2.36)$$

**Remark 1.9.** As emphasized in e.g. [168, 312], the volume expression (1.2.32) requires less regularity on the data (only the one obtained from the variational setting) than (1.2.36). For example, it is well-known that  $u_f, u_s$  may fail to be  $H^2$  functions when the domains  $\Omega_s$  or  $\Omega_f$  involve corners [76, 176]; in such case (1.2.32) remains valid while the surface expression (1.2.36) may become invalid. Furthermore, some authors have found evidence that the Fréchet derivatives of shape functional are better approximated when discretizing the volume form [188]. Let us note, however, that regularity assumptions are still needed (i.e.  $DJ(\Omega)$  continuous over  $H^1(D, \mathbb{R}^d)$  and not only  $W^{1,\infty}(D, \mathbb{R}^d)$ ) if one wants to identify  $DJ(\Omega)$  with a  $H^1$  scalar product.

### Céa's “fast” derivation method

The previous result (1.2.36) shall now be retrieved with Céa's classical method, which is faster in the sense that it circumvents the need for computing the Lagrangian derivative  $\dot{u}$  (eqn. (1.2.27)). However, the computation is only formal and is prone to errors. This part follows [17], section 6.4.3, in a slightly more complicated context. The principle of the method is to introduce a Lagrangian function

$$\mathcal{L}(\Omega, \hat{u}, \hat{p}, \hat{\lambda}) := \int_{\Omega} j(x, \hat{u}(x)) dx + \int_{\Omega} (f + \Delta \hat{u}) \hat{p} dx + \int_{\Gamma_D} \hat{\lambda} \hat{u} ds + \int_{\Gamma_N} \hat{\lambda} \frac{\partial \hat{u}}{\partial \mathbf{n}} ds \quad (1.2.37)$$

defined for  $\hat{u}, \hat{p}, \hat{\lambda} \in H^1(\mathbb{R}^d)$  (they are defined on  $\mathbb{R}^d$  and not on  $\Omega$  to make them independent from the shape  $\Omega$ ). The variable  $\hat{\lambda}$  is a Lagrange multiplier for boundary constraints, and  $\hat{p}$  is a Lagrange multiplier for the PDE constraint (1.2.19) satisfied by the solution  $u(\Omega)$ . In other words, the partial derivatives of the Lagrangian  $\mathcal{L}$  with respect to  $\hat{p}$  and  $\hat{\lambda}$  cancel at  $\hat{u} = u$  where  $u$  is the solution to the Poisson problem (1.2.19):

$$\forall \hat{p}, \hat{\lambda} \in H^1(\mathbb{R}^d), \quad \frac{\partial \mathcal{L}}{\partial \hat{p}}(\Omega, u, \hat{p}, \hat{\lambda}) = 0 \quad \text{and} \quad \frac{\partial \mathcal{L}}{\partial \hat{\lambda}}(\Omega, u, \hat{p}, \hat{\lambda}) = 0.$$

The values of  $\hat{p}$  and  $\hat{\lambda}$  are sought in order to cancel the partial derivative  $\partial \mathcal{L} / \partial \hat{u}$  for  $\hat{u} = u$ . Performing a double integration by part, it holds:

$$\begin{aligned} \forall v \in H^1(\mathbb{R}^d), \quad \frac{\partial \mathcal{L}}{\partial \hat{u}}(\Omega, u, \hat{p}, \hat{\lambda})(v) &= \int_{\Omega} \frac{\partial j}{\partial u}(x, \hat{u}(x)) v dx + \int_{\Omega} \hat{p} \Delta v dx + \int_{\Gamma_D} \hat{\lambda} v ds + \int_{\Gamma_N} \hat{\lambda} \frac{\partial v}{\partial \mathbf{n}} ds \\ &= \int_{\Omega} \frac{\partial j}{\partial u}(x, \hat{u}(x)) v dx + \int_{\Omega} v \Delta \hat{p} dx + \int_{\partial\Omega} \left( \hat{p} \frac{\partial v}{\partial \mathbf{n}} - v \frac{\partial \hat{p}}{\partial \mathbf{n}} \right) ds + \int_{\Gamma_D} \hat{\lambda} v ds + \int_{\Gamma_N} \hat{\lambda} \frac{\partial v}{\partial \mathbf{n}} ds. \end{aligned} \quad (1.2.38)$$

Therefore this partial derivative vanishes by setting  $\hat{p} = p$  where  $p$  is the adjoint variable defined by (1.2.30), and

$$\hat{\lambda} = \lambda := \begin{cases} \frac{\partial p}{\partial \mathbf{n}} & \text{on } \Gamma_D \\ -p & \text{on } \Gamma_N. \end{cases}$$

The boundary expression for the shape derivative of  $J(\Omega, u(\Omega))$  follows then by remarking that

$$J(\Omega_{\boldsymbol{\theta}}, u(\Omega_{\boldsymbol{\theta}})) = \mathcal{L}(\Omega_{\boldsymbol{\theta}}, u(\Omega_{\boldsymbol{\theta}}), \widehat{p}, \widehat{\lambda})$$

for any  $\widehat{p}, \widehat{\lambda} \in H^1(\mathbb{R}^d)$ , which allows to write

$$\begin{aligned} \frac{d}{d\boldsymbol{\theta}}[J(\Omega_{\boldsymbol{\theta}}, u(\Omega_{\boldsymbol{\theta}}))](\boldsymbol{\theta}) &= \frac{d}{d\boldsymbol{\theta}}[\mathcal{L}(\Omega_{\boldsymbol{\theta}}, u(\Omega_{\boldsymbol{\theta}}), \widehat{p}, \widehat{\lambda})] \\ &= \frac{\partial}{\partial \boldsymbol{\theta}}(\mathcal{L}(\Omega_{\boldsymbol{\theta}}, u(\Omega), \widehat{p}, \widehat{\lambda}))(\boldsymbol{\theta}) + \frac{\partial \mathcal{L}}{\partial \widehat{u}}(\Omega, u(\Omega), \widehat{p}, \widehat{\lambda})(u'(\boldsymbol{\theta})). \end{aligned} \quad (1.2.39)$$

where

$$u'(\boldsymbol{\theta}) = \frac{d}{d\boldsymbol{\theta}}u(\Omega_{\boldsymbol{\theta}}) \quad (1.2.40)$$

is the *Eulerian* derivative of  $u(\Omega)$  ([remark 1.7](#)). Setting  $\widehat{p} = p$  and  $\widehat{\lambda} = \lambda$  in (1.2.39) allows to cancel the term involving  $u'(\boldsymbol{\theta})$ , which yields then

$$\begin{aligned} \frac{d}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=0} [J(\Omega_{\boldsymbol{\theta}}, u(\Omega_{\boldsymbol{\theta}}))](\boldsymbol{\theta}) &= \frac{\partial}{\partial \boldsymbol{\theta}}[\mathcal{L}(\Omega_{\boldsymbol{\theta}}, u(\Omega), p, \lambda)] \\ &= \int_{\partial\Omega} j(x, u(x)) \boldsymbol{\theta} \cdot \mathbf{n} ds + \int_{\Gamma_D} \left( \kappa + \frac{\partial}{\partial \mathbf{n}} \right) \left( \frac{\partial p}{\partial \mathbf{n}} u \right) \boldsymbol{\theta} \cdot \mathbf{n} ds - \int_{\Gamma_N} \operatorname{div}(p \nabla u) \boldsymbol{\theta} \cdot \mathbf{n} ds \\ &= \int_{\partial\Omega} j(x, u(x)) \boldsymbol{\theta} \cdot \mathbf{n} ds + \int_{\Gamma_D} \frac{\partial u}{\partial \mathbf{n}} \frac{\partial p}{\partial \mathbf{n}} \boldsymbol{\theta} \cdot \mathbf{n} ds + \int_{\Gamma_N} (fp - \nabla u \cdot \nabla p) \boldsymbol{\theta} \cdot \mathbf{n} ds. \end{aligned}$$

Note the use of [proposition 1.7](#) with the assumption that  $\boldsymbol{\theta} = 0$  on  $\partial\Gamma_N$  in the shape differentiation of the term

$$\int_{\Gamma_N} \lambda \frac{\partial u}{\partial \mathbf{n}} ds = \int_{\Gamma_N} (\lambda \nabla u) \cdot \mathbf{n} ds.$$

This expression coincides with (1.2.36) remembering the boundary conditions satisfied by  $u$  and  $p$ . Albeit this method allows to avoid to explicit Eulerian or Lagrangian derivatives  $u'(\boldsymbol{\theta})$  and  $\dot{u}(\boldsymbol{\theta})$ , it remains rather technical and does not work without hypothesis on the objective functional. Furthermore, this method remains formal (and may even yield wrong formulas, see e.g. [253]) because (i) the Eulerian derivative  $u'(\boldsymbol{\theta})$  (eqn. (1.2.40)) has less regularity than  $u(\Omega)$ , which could imply

$$\frac{\partial \mathcal{L}}{\partial \widehat{u}}(\Omega, u(\Omega), p, \lambda)(u'(\boldsymbol{\theta})) \neq 0,$$

and (ii) the Lagrangian  $\mathcal{L}$  defined in (1.2.37) implicitly assume very high regularity for  $u$  and  $p$ .

### 1.3 ON THE SIGNED DISTANCE FUNCTION AND ITS MAIN PROPERTIES

The signed distance function to a domain (as illustrated on [Figure 1.3](#)) is a mathematical object commonly used in a variety of applications of the level-set method [96, 248, 280, 247] (reviewed in [section 1.4.2](#) below). In the field of shape optimization based on the method of Hadamard, it also allows to formulate geometric distance constraints such as minimum or maximum thickness [234, 30]. Its definition and regularity properties are summarized in [section 1.3.1](#). [Section 1.3.2](#) recalls then differentiability results with respect to the domain and their use for obtaining the shape derivatives of distance constraints. Finally, a few words are given in [section 1.3.3](#) regarding the numerical computation of the signed distance function on discretization meshes. The reader is referred to [64, 119, 122, 23] for proofs and much more exhaustive material.

#### 1.3.1 Definition of the signed distance function and its first two derivatives

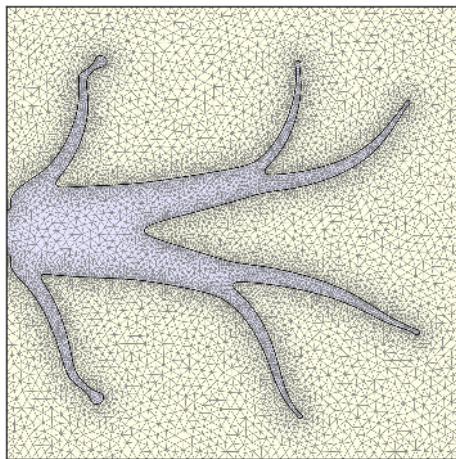
**Definition 1.2** (Signed distance function). The signed distance function to a bounded open domain  $\Omega \subset \mathbb{R}^d$  is the function

$$d_{\Omega} : \mathbb{R}^d \rightarrow \mathbb{R}$$

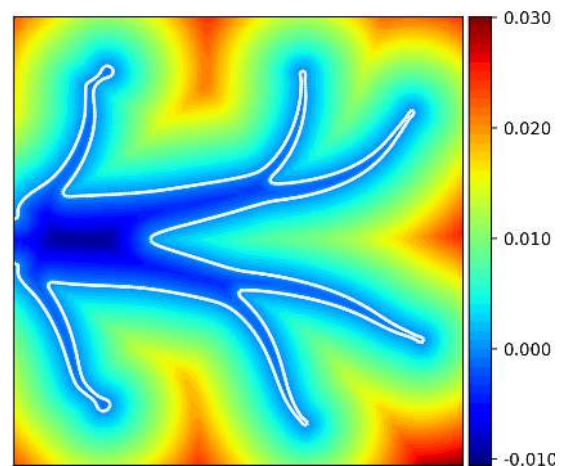
defined for any  $x \in \mathbb{R}^d$  by

$$d_{\Omega}(x) = \begin{cases} - \inf_{y \in \partial\Omega} \|x - y\| & \text{if } x \in \Omega, \\ \inf_{y \in \partial\Omega} \|x - y\| & \text{if } x \notin \Omega. \end{cases}$$

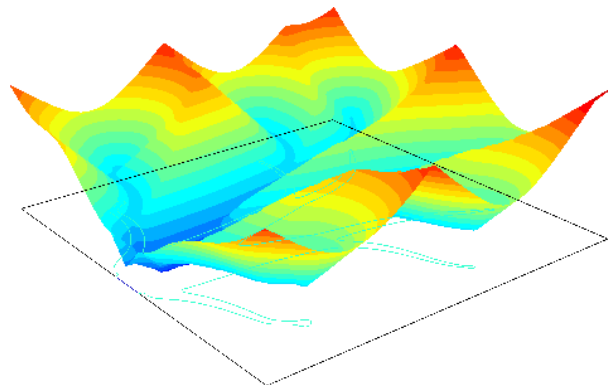




(a) Meshed subdomain  $\Omega \subset D$  (in blue) of a computational domain  $D$ .



(b) Isocontours of the signed distance function  $d_\Omega$ .



(c) 3-d plot of  $d_\Omega$ .

Figure 1.3: Example of signed distance function  $d_\Omega$  numerically computed on a meshed domain.

It is obvious that (i)  $d_\Omega(x) = 0$  for any  $x \in \partial\Omega$ , and (ii) that the infimum involved in the definition of  $d_\Omega$  is attained in possibly several points of  $\partial\Omega$  that are called *projections* onto  $\partial\Omega$ . It is useful to distinguish points that admit a unique projection from the others.

**Definition 1.3** (Skeleton set and projection).

1. The set of points  $x \in \mathbb{R}^d$  for which the minimization problem

$$\min_{y \in \partial\Omega} \|x - y\| \quad (1.3.1)$$

admits several minimizers is called the *skeleton* of  $\Omega$  and is denoted by  $\Sigma$ .

2. For any  $x \in \mathbb{R}^d \setminus \Sigma$ , the unique minimizer of (1.3.1) is denoted  $p_{\partial\Omega}(x)$  and is called the (orthogonal) *projection* of  $x$  onto  $\partial\Omega$ , in that case it holds

$$\forall x \in \mathbb{R}^d \setminus \Sigma, d_\Omega(x) = \begin{cases} -\|x - p_{\partial\Omega}(x)\| & \text{if } x \in \Omega, \\ \|x - p_{\partial\Omega}(x)\| & \text{if } x \notin \Omega. \end{cases}$$

The skeleton  $\Sigma$  of a domain  $\Omega$  and the projection  $p_{\partial\Omega}$  are illustrated on [Figure 1.4](#).

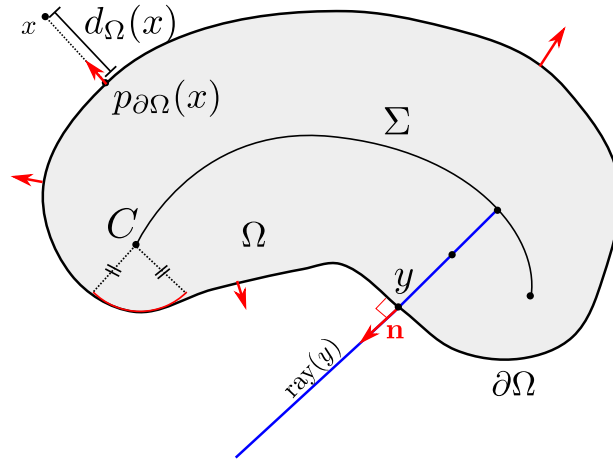


Figure 1.4: Shape  $\Omega \subset \mathbb{R}^d$ , skeleton  $\Sigma$  with normal rays  $(\text{ray}(y))_{y \in \partial\Omega}$ , outward normal vector  $\mathbf{n}$ , center of curvature  $C$  and orthogonal projection of a point  $x \in D$  onto  $\partial\Omega$ .

It is elementary to show that without any further assumption,  $d_\Omega$  is a Lipschitz map with Lipschitz constant smaller than one. Therefore, from Rademacher's theorem [146, 147],  $d_\Omega$  is differentiable almost everywhere, and it can be proved that  $d_\Omega$  is actually differentiable on  $\mathbb{R}^d \setminus (\Sigma \cup \partial\Omega)$  which implies that  $\Sigma$  has zero Lebesgue measure (see [122, 146]). It is possible to be more explicit in the case where  $\Omega$  is a  $\mathcal{C}^1$  domain:

**Proposition 1.10** (Differentiability of  $d_\Omega$ ). *Assume  $\Omega$  is a  $\mathcal{C}^1$  domain with outward normal  $\mathbf{n}$ . The signed distance function  $d_\Omega$  is differentiable at any point  $x \in \mathbb{R}^d \setminus \Sigma$ , and it is not differentiable on  $\Sigma$ . The gradient  $\nabla d_\Omega$  is an extension of the unit normal vector  $\mathbf{n}$  to  $\partial\Omega$  pointing outward  $\Omega$ :*

$$\forall x \in \mathbb{R}^d \setminus \Sigma, \nabla d_\Omega(x) = \mathbf{n}(p_{\partial\Omega}(x)). \quad (1.3.2)$$

In particular,  $d_\Omega$  solves the so-called “Eikonal” equation:

$$\begin{cases} \|\nabla d_\Omega\| = 1 & \text{in } \mathbb{R}^d \setminus \Sigma, \\ d_\Omega = 0 & \text{on } \partial\Omega. \end{cases} \quad (1.3.3)$$

Formulas for the differential of the projection map  $p_{\partial\Omega}$  and the Hessian  $\nabla^2 d_\Omega$  are available under the additional assumption that  $\Omega$  is a domain of class  $\mathcal{C}^2$ , referring to [2, 122] for less regular contexts. In that case, the use of the basis  $(\boldsymbol{\tau}_i)_{1 \leq i \leq d-1}$  associated with the principal curvatures  $(\kappa_i)_{1 \leq i \leq d-1}$  of  $\partial\Omega$  (proposition 1.5) yields (see [36, 64, 82, 157] for proofs):

**Proposition 1.11** (Differentiability of  $p_{\partial\Omega}$  and Hessian of  $d_\Omega$ ). *If  $\Omega$  is a  $C^2$  bounded domain, then the following properties hold:*

1. for any point  $x \in \mathbb{R}^d \setminus \bar{\Sigma}$ , one has:

$$1 + \kappa_i(p_{\partial\Omega}(x))d_\Omega(x) > 0, \quad i = 1, \dots, d-1;$$

2. the projection  $p_{\partial\Omega}$  is differentiable on  $\mathbb{R}^d \setminus \bar{\Sigma}$  and

$$\forall x \in \mathbb{R}^d \setminus \bar{\Sigma}, \quad \nabla p_{\partial\Omega}(x) = \sum_{i=1}^{d-1} \frac{1}{1 + \kappa_i(p_{\partial\Omega}(x))d_\Omega(x)} \boldsymbol{\tau}_i(p_{\partial\Omega}(x)) \boldsymbol{\tau}_i(p_{\partial\Omega}(x))^T; \quad (1.3.4)$$

3. the signed distance function  $d_\Omega$  is twice differentiable on  $\mathbb{R}^d \setminus \bar{\Sigma}$  and the Hessian  $\nabla^2 d_\Omega$  is given by

$$\forall x \in \mathbb{R}^d \setminus \bar{\Sigma}, \quad \nabla^2 d_\Omega(x) = \sum_{i=1}^{d-1} \frac{\kappa_i(p_{\partial\Omega}(x))}{1 + \kappa_i(p_{\partial\Omega}(x))d_\Omega(x)} \boldsymbol{\tau}_i(p_{\partial\Omega}(x)) \boldsymbol{\tau}_i(p_{\partial\Omega}(x))^T.$$

If in addition,  $\Omega$  is a  $C^3$  domain, then  $\mathbb{R}^d \setminus \bar{\Sigma}$  has zero Lebesgue measure (see [223]).

**Remark 1.10.** The projection  $p_{\partial\Omega}$  is not differentiable on the boundary  $\partial\Sigma$  of the skeleton  $\Sigma$ . Indeed (see Figure 1.4 and [82, 157]), a point  $C \in \mathbb{R}^d$  belongs to the closure  $\bar{\Sigma}$  either because it lies on  $\Sigma$ , or because it is a center of curvature of  $\partial\Omega$ , i.e. there exists  $y \in \Pi_{\partial\Omega}(C)$  and  $i = 1, \dots, d-1$  such that

$$1 + d_\Omega(C)\kappa_i(y) = 0.$$

**Remark 1.11.** Formula (1.3.4) can be generalized to arbitrary embedded smooth manifolds, see e.g. [157, 156] where it is used in the context of matrix manifolds for the differentiation of matrix decompositions.

### 1.3.2 Differentiability of the signed distance function with respect to the domain and shape derivatives of distance functionals

Delfour and Zolésio [120] followed by [23, 107] demonstrated that the function  $d_\Omega$  is differentiable with respect to the shape  $\Omega$  in the following sense:

**Proposition 1.12** ([23], Proposition 3.5). *For any  $x \notin \Sigma$ , the map  $\boldsymbol{\theta} \mapsto d_{(I+\boldsymbol{\theta})\Omega}(x)$  is Gâteaux-differentiable at  $\boldsymbol{\theta}$  as an application from  $W^{1,\infty}(\mathbb{R}^d, \mathbb{R}^d)$  into  $\mathbb{R}$  and its derivative reads*

$$d'_\Omega(\boldsymbol{\theta})(x) = -\boldsymbol{\theta}(p_{\partial\Omega}(x)) \cdot \boldsymbol{n}(p_{\partial\Omega}(x)).$$

*Sketch of proof.* A formal proof of this result is obtained as follows. Recall the signed distance function  $d_\Omega$  satisfies

$$\begin{cases} \|\nabla d_\Omega(x)\| = 1 & \text{in } \mathbb{R}^d \setminus \Sigma \\ d_\Omega(x) = 0 & \text{on } \partial\Omega. \end{cases}$$

Differentiating the first equality with respect to  $\Omega$  yields formally  $\nabla d_\Omega \cdot \nabla(d'_\Omega(\boldsymbol{\theta})) = 0$ , which means that  $d'_\Omega(\boldsymbol{\theta})$  is constant along the rays normal to the shape. One then rewrites the second equality in a weak sense on the deformed domain  $\Omega_\boldsymbol{\theta} = (I + \boldsymbol{\theta})\Omega$ :

$$\forall \phi \in H^1(\mathbb{R}^d), \quad \int_{\partial\Omega_\boldsymbol{\theta}} \phi \circ (I + \boldsymbol{\theta})^{-1} d_{\Omega_\boldsymbol{\theta}}(y) dy = 0.$$

Differentiating this equation with respect to  $\boldsymbol{\theta}$  then yields

$$\forall \phi \in H^1(\mathbb{R}^d), \quad \int_{\partial\Omega} \phi (\kappa d_\Omega + \nabla d_\Omega \cdot \boldsymbol{n}) \boldsymbol{\theta} \cdot \boldsymbol{n} dy + \int_{\partial\Omega} \phi d'_\Omega(\boldsymbol{\theta}) dy = 0,$$

from where  $d'_\Omega(\boldsymbol{\theta}) = -\boldsymbol{\theta} \cdot \boldsymbol{n}$  on  $\partial\Omega$  follows. □

In practice, this result is used for taking into account geometric constraints expressed as an integral criterion involving the signed distance function. Let  $D \subset \mathbb{R}^d$  be a “hold-all” domain containing  $\Omega \subset D$  as a subdomain. A geometric constraint functional  $P(\Omega)$  is typically written in the form

$$P(\Omega) := \int_D j(d_\Omega(x)) dx \quad (1.3.5)$$

where  $j : \mathbb{R} \rightarrow \mathbb{R}$  is a given  $\mathcal{C}^1$  function. Maximum thickness, minimum thickness, minimum member's distance [30, 234], or the stress energy of a multimaterial medium involving a smooth interface [23] can be formulated in this form. The application of proposition 1.12 yields that the shape derivative of  $P(\Omega)$  is (in a Gâteaux differentiability sense):

$$P'(\Omega)(\boldsymbol{\theta}) = \int_D j'(d_\Omega(x)) d'_\Omega(\boldsymbol{\theta})(x) dx. \quad (1.3.6)$$

This formula is, at first glance, not satisfying because it does not satisfy Hadamard's structure theorem (proposition 1.1). A more explicit expression is obtained thanks to a suitable change of variable which involves integrals along the normal rays to the shape. More precisely, for  $y \in \partial\Omega$ , the ray emerging from  $y$  is defined to be the one-dimensional segment

$$\text{ray}(y) := \{x \in D \setminus \bar{\Omega}, p_{\partial\Omega}(x) = y\}.$$

The application of the coarea formula (decomposing  $D$  onto the level sets of  $p_{\partial\Omega}$ , see [23]), or a suitable change of variable (detailed later in chapter 4) yields:

$$P'(\Omega)(\boldsymbol{\theta}) = \int_{\partial\Omega} \left( \int_{z \in \text{ray}(y)} j'(d_\Omega(z)) \prod_{i=1}^{d-1} (1 + \kappa_i(y) d_\Omega(z)) dz \right) \boldsymbol{\theta}(y) \cdot \mathbf{n}(y) dy$$

which verifies Hadamard's structure theorem. This formula remains however very delicate to implement in practice; we shall introduce a new characterization for  $P'(\Omega)(\boldsymbol{\theta})$  in chapter 4 which allows to evaluate the above integral without the need for explicitly calculating  $\text{ray}(y)$  nor the shape curvatures  $\kappa_i(y)$  from the discretization of  $\partial\Omega$ .

### 1.3.3 Numerical methods for the computation of the signed distance function

In practical PDE constrained shape optimization algorithms, the domain  $\Omega \subset D$  to optimize is known only in a discrete form. Most common types of discretization can be broadly divided into two categories:

1. explicit discretization: the shape  $\Omega$  is explicitly discretized into a finite element mesh, such as in the example of Figure 1.3a.
2. implicit discretization: the shape  $\Omega$  is implicitly represented by nodal values of a *level set* function  $\phi$  discretized on a *fixed* meshed domain  $D$ :

$$\Omega = \{x \in D \mid \phi(x) < 0\}.$$

In that case, the shape  $\Omega$  is retrieved from the 0 isocontour of the level set  $\phi$ , as visible on Figure 1.3b.

Various numerical algorithms exist for computing the signed distance function  $d_\Omega$  in these various discrete settings. These generally rely on the fact that  $d_\Omega$  is a viscosity solution of the Eikonal equation (1.3.3), which can be computed as a steady-state solution of the so called “redistanciation” equation:

$$\partial_t \phi + \text{sign}(\phi)(\|\nabla \phi\| - 1) = 0. \quad (1.3.7)$$

It can be shown that the viscosity solution  $\phi(t, \cdot)$  of (1.3.7) (see [61]) converges indeed to  $d_\Omega$  where  $\Omega = \{x \in D \mid \phi(0, x) < 0\}$ . This can be used:

1. to devise numerical schemes on *fixed grid* transforming an input level set function  $\phi(0, \cdot)$  into a signed distance function  $d_\Omega$ , see e.g. [247], chapter 7. This is of particular interest in implicit methods because the signed distance function is a particular level set function characterized by good numerical properties [247];

2. to devise numerical algorithms on triangular or tetrahedral mesh to build the signed distance function of a subdomain discretized as a submesh, provided the distance has been computed first on vertex adjacent to the subdomain boundary. This approach is e.g. followed by [111] and implemented in the software `mshdist`, which we used in our own implementations of shape optimization test cases.

Several other methods exist for the computation of the signed distance function such as the Fast Marching Method (see [280], or [204] for a version on simplicial meshes), or the Fast Sweeping method [323]; the reader is referred to [107], section 1.3.1. for a review.

#### 1.4 A CLASSICAL SHAPE OPTIMIZATION NUMERICAL WORKFLOW USING A LEVEL-SET BASED MESH EVOLUTION METHOD

This last section describes the classical main steps involved in the numerical implementation of the method of Hadamard in order to solve shape optimization problems of the form

$$\min_{\Omega} J(\Omega), \quad (1.4.1)$$

where  $J$  may depend on the solutions to some partial differential equations as in (1.2.21).

In practice, additional equality or inequality constraints come into play in most applications. For now, it is sufficient to assume that these can be treated by penalization techniques which reduce constrained problems to unconstrained one of the same form of (1.4.1). This approach has been actually commonly used in literature dealing with problems featuring only a few constraints [190, 203, 29, 234]. For instance, if there is only one constraint  $P(\Omega) = 0$ , the Augmented Lagrangian Method [244] considers

$$\min_{\Omega} J(\Omega) + \lambda P(\Omega) + \frac{1}{2} \mu P(\Omega)^2, \quad (1.4.2)$$

where  $\lambda$  and  $\mu$  are tuning parameters which may be updated in the course of the optimization process. This method is unsatisfying for a number of reasons; for constrained optimization with an arbitrary number of constraints, all our numerical examples rely in fact on a different algorithm which is described in details in the dedicated [chapter 3](#).

The main steps of a typical topology optimization algorithm for the minimization of PDE constrained functionals are summarized in [algorithm 1.1](#) below. These algorithms involve two classical cornerstones: the computation of a descent direction (a vector field  $\boldsymbol{\theta} \in W^{1,\infty}(\mathbb{R}^d, \mathbb{R}^d)$  yielding improved design shapes), and the representation and practical updates of shapes in the numerical setting. These two steps are reviewed in details in the next two sections. [section 1.4.1](#) reviews the well known necessity to regularize and extend the shape derivative and its consistency with respect to the gradient method. Then, the problem of numerically evolving shapes is discussed in [sections 1.4.1](#) and [1.4.2](#), which includes a summary of the level set mesh evolution method which we used in our numerical test cases in the latter section.

---

**Algorithm 1.1** Shape optimization with the method of Hadamard and a level set based mesh evolution algorithm.

---

Generate an initial mesh  $\mathcal{T}_0$  for a computational domain  $D$  adapted to a first guess subdomain  $\Omega_0 \subset D$ .

**for**  $n = 0, 1, 2 \dots$  **do**

1. Solve the physical equations posed on  $\Omega_n$  (e.g. (1.2.19)).
2. Assemble the shape derivative  $DJ(\Omega_n)$  (e.g. (1.2.36) and [section 1.2.2](#)) after solving adjoint equations (such as (1.2.30)).
3. Identify the shape derivative  $DJ(\Omega_n)$  to a gradient  $\nabla J(\Omega_n)$  by solving an extension regularization problem of the form of (1.4.4) (see [section 1.4.1](#)).
4. Update the shape according to the descent direction  $\boldsymbol{\theta}_n = -\Delta t \nabla J(\Omega_n)$  where  $\Delta t$  is a small discretization step. This can be done by using mesh deformations (with (1.4.3)) or with the level set based mesh evolution algorithm of [algorithm 1.2](#) (with (1.4.12)), see [section 1.4.2](#).

**end for**

---

### 1.4.1 Gradient based optimization in the context of the method of Hadamard: identification, regularization, and extensions of shape derivatives

Numerical algorithms for solving shape optimization problems of the form

$$\min_{\Omega} J(\Omega).$$

compute a minimizing sequence of domains  $(\Omega_n)_{n \in \mathbb{N}}$  which gradually decrease the objective function  $J(\Omega)$ . In the context of the method of Hadamard, the domain  $\Omega_{n+1}$  is obtained by deformation of  $\Omega_n$  along a vector field  $\boldsymbol{\theta}_n \in W^{1,\infty}(\mathbb{R}^d, \mathbb{R}^d)$ :

$$\Omega_{n+1} = (I + \boldsymbol{\theta}_n)\Omega_n, \quad \forall n \geq 0. \quad (1.4.3)$$

The principle of the gradient method is to use the knowledge of the shape derivative in order to build a descent direction  $\boldsymbol{\theta}_n$  that decreases the objective function  $J(\Omega)$  in some optimal sense.

In this infinite dimensional context, the shape derivative  $DJ(\Omega)$  of the objective function  $J(\Omega)$  at a given domain  $\Omega$  is not strictly speaking a descent direction: indeed,  $DJ(\Omega)$  belongs to the dual space  $W^{1,\infty}(\mathbb{R}^d, \mathbb{R}^d)^*$  and not to  $W^{1,\infty}(\mathbb{R}^d, \mathbb{R}^d)$  itself. This issue is numerically related to the extension and regularization step of the shape derivative, which is very classical issue acknowledged in a number of works since the early developments of shape optimization [87, 86, 49, 48, 81, 34, 117, 236]. Practically, a descent direction can be obtained by identifying the linear form  $DJ(\Omega) \in W^{1,\infty}(D, \mathbb{R}^d)^*$  to an actual gradient  $\nabla J(\Omega) \in W^{1,\infty}(D, \mathbb{R}^d)$ . To do so, it is sufficient to consider a Hilbert space  $V$  sufficiently “small” for the inclusion  $V \subset W^{1,\infty}(\Omega, \mathbb{R}^d)$  to hold. The converse inclusion for the dual spaces,  $W^{1,\infty}(\Omega, \mathbb{R}^d)^* \subset V^*$ , implies the continuity of  $DJ(\Omega)$  on  $V$ . Therefore, denoting by  $\langle \cdot, \cdot \rangle_V$  the scalar product of  $V$ , the standard Riesz representation theorem [76, 49, 81, 117] yields the existence of a unique element  $\nabla J(\Omega) \in V$  such that

$$\forall \boldsymbol{\theta} \in V, \langle \nabla J(\Omega), \boldsymbol{\theta} \rangle_V = DJ(\Omega)(\boldsymbol{\theta}). \quad (1.4.4)$$

The opposite of the vector field  $\nabla J(\Omega) \in V \subset W^{1,\infty}(D, \mathbb{R}^d)$  can then serve as a descent direction to obtain a better shape. Indeed, it holds for any sufficiently small step size  $\Delta t > 0$

$$\begin{aligned} J((I - \Delta t \nabla J(\Omega))\Omega) &= J(\Omega) - \Delta t DJ(\Omega)(\nabla J(\Omega)) + o(\Delta t) \\ &= J(\Omega) - \Delta t \langle \nabla J(\Omega), \nabla J(\Omega) \rangle_V + o(\Delta t), \end{aligned} \quad (1.4.5)$$

with  $-\langle \nabla J(\Omega), \nabla J(\Omega) \rangle_V \leq 0$ , which shows that updating  $\Omega$  by deformation with the vector field  $\boldsymbol{\theta} = -\Delta t \nabla J(\Omega)$  leads to a non-increase of the objective function. Furthermore,  $-\nabla J(\Omega)$  is a “best” descent direction in the sense that it is better than any other element of  $V$  having the same norm: in [chapter 3](#), we recall that if  $DJ(\Omega) \neq 0$ , then

$$-\frac{\nabla J(\Omega)}{\|\nabla J(\Omega)\|_V} = \min_{\substack{\boldsymbol{\theta} \in V, \\ \|\boldsymbol{\theta}\|_V \leq 1}} DJ(\Omega)(\boldsymbol{\theta}).$$

As for the choice of the Hilbert space  $V \subset W^{1,\infty}(D, \mathbb{R}^d)$  used in the identification (1.4.4), one can set  $V = H^m(D, \mathbb{R}^d)$  with  $m > 1 + d/2$ , equipped with its standard inner product: indeed, the inclusion  $H^m(D, \mathbb{R}^d) \subset W^{1,\infty}(D, \mathbb{R}^d)$  holds as a consequence of the Sobolev embedding theorem [76]. In this case, the identification problem (1.4.4) reduces to a linear elliptic problem of order  $2m$ . What is more, when  $DJ(\Omega)(\boldsymbol{\theta})$  can be written in the form of a boundary integral, namely when there exists  $v_J(\Omega) \in L^1(\partial\Omega)$  such that

$$\forall \boldsymbol{\theta} \in W^{1,\infty}(D, \mathbb{R}^d), \quad DJ(\Omega)\boldsymbol{\theta} = \int_{\partial\Omega} v_J(\Omega) \boldsymbol{\theta} \cdot \boldsymbol{n} ds, \quad (1.4.6)$$

then  $\nabla J(\Omega) \in H^m(D, \mathbb{R}^d)$  is a regularized extension of the field  $v_J(\Omega)\boldsymbol{n}$  to the whole domain  $D$ , which turns to be very convenient for practical numerical algorithms.

A very common strategy in the literature (see for instance [34, 49, 81, 153, 117, 236]), though, consists in taking simply  $V = H^1(D, \mathbb{R}^d)$  with the inner product

$$\forall \boldsymbol{\theta}, \boldsymbol{\theta}' \in V, \langle \boldsymbol{\theta}, \boldsymbol{\theta}' \rangle_V := \int_D (\gamma^2 \nabla \boldsymbol{\theta} : \nabla \boldsymbol{\theta}' + \boldsymbol{\theta} \cdot \boldsymbol{\theta}') dx, \quad (1.4.7)$$

where  $\gamma > 0$  is a user-defined parameter which can physically be interpreted as a length-scale for the regularity of the deformations  $\boldsymbol{\theta}$ ; typically,  $\gamma$  is of the order of the minimum mesh element size of the

discretization. Variants can be considered for tuning more finely the smoothness of such extensions, or to prescribe non optimizable boundaries by setting zero Dirichlet boundary conditions for  $\nabla J(\Omega)$  in (1.4.6). This choice of  $V$  is *a priori* only formal because  $V$  is not rigorously a subspace of  $W^{1,\infty}(D, \mathbb{R}^d)$ : it works only if  $DJ(\Omega)$  is continuous over  $H^1(D, \mathbb{R}^d)$ . This condition is actually not a limitation because it turns to be often satisfied for PDE constrained problems under suitable regularity assumptions on  $\Omega$  and the data. Furthermore, this choice is very convenient numerically because this space is easily discretized with  $\mathbb{P}_1$  finite elements. Since this leads to very good results in practice, we rely on this strategy in our own numerical implementation.

Last, it is possible to even reduce the cost of the computation of  $\nabla J(\Omega)$  to the one of the resolution of a simple scalar elliptic problem if one considers  $V = \{v \nabla d_\Omega \mid v \in H^1(\Omega)\}$  equipped with the scalar product  $\langle \cdot, \cdot \rangle_V$  defined by

$$\forall v, w \in H^1(\Omega), \langle v \nabla d_\Omega, w \nabla d_\Omega \rangle_V := \int_\Omega (\gamma \nabla v \cdot \nabla w + vw) dx. \quad (1.4.8)$$

Remembering that  $\nabla d_\Omega(x) = \mathbf{n}(p_{\partial\Omega}(x))$  is a continuous extension of the outward vector field normal to  $\partial\Omega$  (proposition 1.10), (1.4.4) reduces to the identification of a single normal component  $v$  for the updating vector field  $\boldsymbol{\theta} = -v \nabla d_\Omega$ . In practice, one resorts to (1.4.8) when using boundary integral expressions of shape derivatives such as (1.2.36), and to (1.4.7) when relying rather on volumetric expressions of the form of (1.2.32).

### 1.4.2 Numerical representations and updates of shapes: mesh deformations, level set methods, and level set based mesh evolution method

One of the key issues of shape optimization algorithms lies in the numerical implementation of the shape update step (1.4.3). This step involves (i) a numerical representation of the shape  $\Omega_n$  and (ii) a process that allows to obtain a new shape  $\Omega_{n+1}$  in this representation from a displacement vector field  $\boldsymbol{\theta}_n$ . Furthermore, in a context involving partial differential equations, the numerical representation of  $\Omega$  should be compatible with physical solvers relying e.g. on the Finite Element method.

In what follows, two complementary numerical representations of shapes together with their related method for updating them are briefly reviewed: explicit methods which evolve a meshed representation, and *level set* methods. Then, we briefly describe the main steps of the level set mesh evolution algorithm of Allaire, Dapogny and Frey [25], which combines advantages of both representations.

#### Explicit mesh representations and nodal displacements

The first implementations of the method of Hadamard for shape optimization were considering explicit discretizations of shapes [259, 86, 34]. An open domain  $\Omega$  is typically represented by a simplicial mesh (a set of triangles or tetrahedra) which is globally *conforming* in the sense that the intersection between any two triangles or tetrahedra is either empty, or reduced to a vertex, or an edge, or a triangle of the mesh (Figure 1.5).

This representation is convenient for computing solutions to partial differential equations on  $\Omega$  with the finite element method domain such as (1.2.19), as well as volume or surface integrals involved in the evaluation of shape derivatives such as (1.2.32) and (1.2.36). Furthermore, it is easily amenable to *small* deformation updates by translations of mesh vertices. If  $\Omega$  is discretized into a mesh involving  $N$  vertices  $x_1, \dots, x_N$ , then a deformed shape  $(I + \boldsymbol{\theta})\Omega$  with  $\boldsymbol{\theta} \in W^{1,\infty}(\mathbb{R}^d, \mathbb{R}^d)$  is obtained by applying the transformation

$$x_i \leftarrow x_i + \boldsymbol{\theta}(x_i), \quad 1 \leq i \leq N.$$

The process is not computationally expensive and can be iterated in order to implement (1.4.3) for obtaining a minimizing sequence  $\Omega_n$  from descent updates  $\boldsymbol{\theta}_n$  until a final design. However, it is well-known [259] that in very few iterations, the quality of the deformed mesh may quickly decrease due to the occurrence of nearly flat or inverted mesh elements, calling for the use of remeshing, or for the use of other numerical shape representations.

#### Level set methods: advection and Hamilton-Jacobi equations

The level set method was initially introduced by S. Osher and J. Sethian for solving free boundary problems in fluid mechanics [248], before being used for shape optimization in [31, 32, 250, 281, 311].

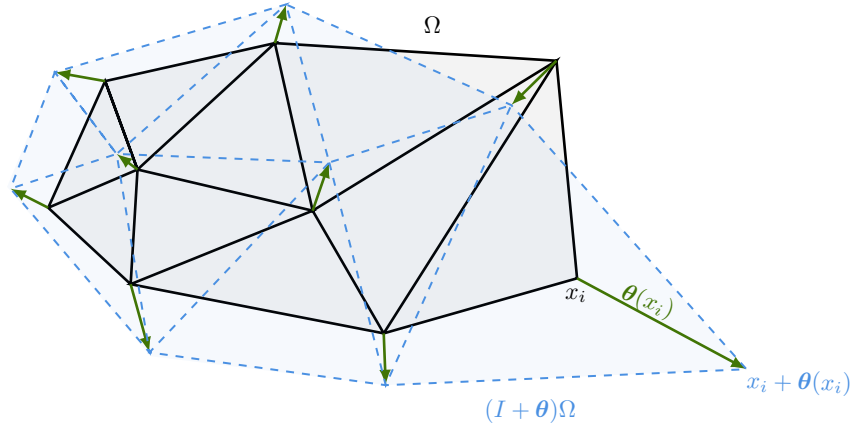


Figure 1.5: Discretization of a 2-d domain  $\Omega$  in a simplicial mesh and its deformation by application of a displacement vector field  $\theta$ .

Given a fixed computational domain  $D$ , it amounts to representing any subdomain  $\Omega \subset D$  as the negative subdomain of a scalar “level set” function  $\phi : D \rightarrow \mathbb{R}$ :

$$\begin{cases} \phi(x) < 0 & \text{if } x \in \Omega, \\ \phi(x) = 0 & \text{if } x \in \partial\Omega, \\ \phi(x) > 0 & \text{if } x \in D \setminus \Omega. \end{cases} \quad (1.4.9)$$

The motion of a domain  $\Omega(t)$  in  $D$  according to a vector velocity field  $\theta(x)$  translates then in terms of an associated level set function  $\phi(t, x)$  as the following advection equation:

$$\begin{cases} \frac{\partial \phi}{\partial t}(t, x) + \theta(t, x) \cdot \nabla \phi(t, x) = 0, & x \in D, \\ \phi(0, x) = \phi_0(x), & x \in D, \end{cases} \quad (1.4.10)$$

where  $\phi_0$  is any level set function for  $\Omega$ . This representation is very convenient for describing the evolution of domains at the discrete level on *fixed* meshes, because topological changes are handled naturally and automatically (Figure 1.6).

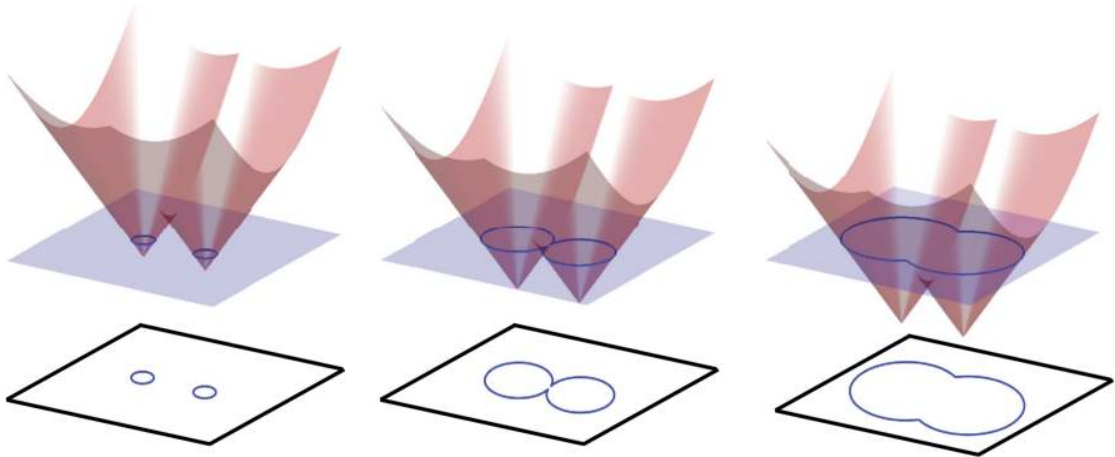


Figure 1.6: Motion of an implicit domain featuring a topological change and its representation by level set functions. Figure from [170].

When the time dependent domain  $\Omega(t)$  evolves according to a motion  $\theta(t, x)$  normal to the interface  $\partial\Omega(t)$ , i.e.  $\theta(t, x) = v(t, x)\mathbf{n}$  on  $\partial\Omega$  for some scalar function  $v(t, x)$ , (1.4.10) rewrites as the following Hamilton-Jacobi equation:

$$\frac{\partial \phi}{\partial t}(t, x) + v(t, x) \|\nabla \phi(t, x)\| = 0, \quad t > 0, \quad x \in \Omega. \quad (1.4.11)$$



Generally, the level set function is initialized to the signed distance function of the initial domain and reinitialized from times to times, which allows to reduce numerical diffusion and round off errors [96, 247]. In applications where the computational domain  $D$  is discretized into a finite difference or finite volume cartesian grid, one usually relies on the Hamilton-Jacobi equation (1.4.11) for the evolution of implicit domains  $\Omega(t)$  represented by level set functions  $\phi(t, \cdot)$ . Indeed, efficient and accurate numerical schemes are available for the resolution of (1.4.11) on structured grids [249, 247]. In contexts involving finite element unstructured meshes with level set functions  $\phi$  discretized as a set of nodal values, it can be more convenient to solve the evolution equation (1.4.10), e.g. with the method of characteristics. This is the method we favored for all our numerical examples; we relied on the implementation available in the software `advect` issued from [80] and which shares common features with the semi-Lagrangian method of J. Strain [295].

In shape optimization algorithms (such as in [algorithm 1.1](#)) relying on the level set method, either of the evolution equations (1.4.10) or (1.4.11) is solved with respectively  $\theta(t, x) \equiv \theta_n(x)$  or  $v(t, x) \equiv v_n(x)\mathbf{n}(t, x)$ , where the (time-independent) descent directions  $\theta_n(x)$  and  $v_n(x)$  are computed to evolve the current shape  $\Omega_n$  into the next iterate  $\Omega_{n+1}$ .

Let us note that the level set method for moving domain boundaries (either with (1.4.10) or (1.4.11)) actually involves domain updates different of those (1.4.3) considered by the method of Hadamard : if the current domain  $\Omega_n$  is to be deformed according to a descent direction  $\theta_n \in W^{1,\infty}(D, \mathbb{R}^d)$ , then the updated domain  $\Omega_{n+1}$  is obtained by

$$\Omega_{n+1} = \tilde{\rho}_{\Omega_n}(\theta_n),$$

where for a given subdomain  $\Omega \subset D$ ,  $\tilde{\rho}_{\Omega}$  is the mapping defined by

$$\tilde{\rho}_{\Omega}(\theta) := \{x \in D \mid \phi(1, x) < 0\}, \quad \theta \in W^{1,\infty}(D, \mathbb{R}^d) \quad (1.4.12)$$

with  $\phi(t, x)$  the solution to either (1.4.10) of (1.4.11). The mapping  $\tilde{\rho}_{\Omega}$  is different from the one  $\rho_{\Omega}(\theta) := (I + \theta)\Omega$  considered in the Hadamard's method. However, it can be shown that there is no loss of consistency when using either  $\tilde{\rho}_{\Omega}(\theta)$  or  $\rho_{\Omega}(\theta)$  in the implementation of first order optimization methods. In fact, when  $\phi(t, \cdot)$  is obtained with the advection equation (1.4.10), then  $\tilde{\rho}_{\Omega}$  is the update map considered in the *speed* method of Zolésio for boundary variations [291], which is equivalent at first order to the method of Hadamard. When  $\phi(t, \cdot)$  is obtained by (1.4.11),  $\tilde{\rho}_{\Omega}$  corresponds to moving the points of  $\partial\Omega$  along *bicharacteristics* associated to the Hamilton-Jacobi equation, which is also a second order perturbation of the method of Hadamard (see [21] where this is discussed in details). An equivalent way to formulate this statement is that the domain update maps  $\rho_{\Omega}$  and  $\tilde{\rho}_{\Omega}$  can be interpreted as two admissible *retractions* converting descent directions into new optimization points, this point is further explained in [chapter 3](#).

A common difficulty of the implicit representation on a fixed computational mesh is that it does not make readily possible to solve for original physical equations such as (1.2.19) posed on the implicit domain  $\Omega$ : indeed, the location of  $\Omega$  is known only through the values of the level set function  $\phi$ . Classical level set based shape optimization methods use the so-called “*ersatz*” or fictitious material approach [32] (other interpolation methods are possible in such level set descriptions, such as XFEM or cutFEM methods [140, 309]). For instance, in the context of the Laplace problem considered in [section 1.2.3](#), if only the Neumann boundary  $\Gamma_N$  is optimizable and  $\Gamma_D \subset \partial D$ , this amounts to replace (1.2.19) by an approximate problem set on the whole domain  $D$ :

$$\begin{cases} -\operatorname{div}(A_{\varepsilon}(\Omega)\nabla u_{\varepsilon}) = f & \text{in } D \\ u_{\varepsilon} = 0 & \text{on } \Gamma_D. \end{cases} \quad (1.4.13)$$

The parameter  $A_{\varepsilon}$  is a fictitious conductivity set equal to 1 in the domain  $\Omega$  and to a (small) “*ersatz*” value  $\varepsilon$  in the complementary  $D \setminus \Omega$ . The solution  $u_{\varepsilon}$  can be computed on the mesh adapted to  $D$ , and it satisfies  $u_{\varepsilon} \rightarrow u$  in the  $L^2$  norm on  $\Omega$  as  $\varepsilon \rightarrow 0$ . Using regularized Heaviside and Dirac functions [247], it becomes possible to assemble finite element matrices associated with (1.4.13) and to estimate volume integrals on  $\Omega$  or surface integrals on  $\partial\Omega$  for the implementation of shape optimization algorithms [31, 311].

Finally, let us remark that if the level set method allows *large* domain deformations up to topological changes to occur, these are not natively accounted for by the method of Hadamard. Indeed, shape derivatives ([definition 1.1](#)) quantify the sensitivity of objective functionals with respect to *small* deformations only, without topology changes. Therefore, topological changes occurring in the course of

optimization updates with the level set method are the consequence of the finite resolution of both the spatial and time discretization. For instance, very thin parts of an evolved domain  $\Omega(t)$  typically disappear because either of some numerical diffusion smearing out small details, or because the advection step chosen for the numerical resolution of (1.4.10) is too large with respect to the smallest mesh element size. There is no guarantee that an objective function  $J(\Omega)$  keeps decreasing when such an event suddenly occurs, however it can be expected that the optimization carries on smoothly afterwards until another topological change happens or until convergence.

### The level set based mesh evolution algorithm of Allaire, Dapogny and Frey [24]

In all numerical examples of the present work, we use the level set based mesh evolution method introduced by Allaire, Dapogny and Frey [24]. The main idea is to combine both numerical representations of the domain  $\Omega \subset D$  seen as a subset of a computational domain  $D$  (see Figure 1.3):

- on the one hand, a computational mesh  $\mathcal{T}$  of  $D$  is available, in which  $\Omega$  is explicitly discretized as a meshed subdomain;
- on the other hand,  $\Omega$  is *implicitly* described, using the level set method: e.g. it can be seen as the negative subdomain of a function  $\phi : D \rightarrow \mathbb{R}$ .

It becomes then possible to consistently alternate between both descriptions of the domain  $\Omega$  depending on the nature of the ongoing operation: finite element resolutions are carried out using the meshed description, while the motion of the phases is tracked using the level set method.

The idea of remeshing shapes in between shape optimization iterations was first proposed by Persson and implemented on a test case in his thesis [256]. For our applications, we rely on the subsequent work of [108] for implicit domain meshing. The main features of this method are now outlined in the case of two space dimensions; referring to [25] for further details. Our numerical shape optimization test cases relied on the library `mng` in which the remeshing step of a subdomain described by a level set is implemented both in 2-d and in 3-d.

In what follows, we assume a vector field  $\boldsymbol{\theta}$  is given at the nodes of the computational mesh  $\mathcal{T}$  discretizing  $D$  and  $\Omega \subset D$  (Figure 1.7a). The objective of the algorithm is to compute a new mesh  $\mathcal{T}'$  adapted to the discretization of the evolved domain  $\Omega' := \tilde{\rho}_\Omega(\boldsymbol{\theta})$ , where  $\tilde{\rho}_\Omega$  (eqn. (1.4.12)) encloses the advection step as dictated by the transport equation (1.4.10). The main steps of the algorithm are summarized in algorithm 1.2 and illustrated on Figure 1.8.

---

**Algorithm 1.2** Domain update by mean of a level set based mesh evolution method.

---

1. A level set function  $\phi$  associated with the domain  $\Omega \subset D$  is computed at the nodes of the mesh  $\mathcal{T}$  (Figure 1.7b). As is classical in level set methods, we choose  $\phi = d_\Omega$  to be given by the signed distance function of  $\Omega$ , although any other level set function smooth in a band around  $\partial\Omega$  would be sufficient.
  2.  $\phi$  is advected on the fixed mesh  $D$  by solving the evolution equation (1.4.10) (Figure 1.8a).
  3. The new domain  $\Omega'$  (Figure 1.8b) is discretized by splitting mesh elements of  $\mathcal{T}$  according to the zero level set of the updated level set. This part is purely combinatorial and yields a conforming but poor quality mesh.
  4. Remeshing operations (explained in details [108, 25, 107, 153]) are performed iteratively to improve the mesh quality and the approximation of the discrete surface  $\partial\Omega'$ . Such is possible thanks to the computation of a metric based on the Hausdorff distance to the surface  $\partial\Omega'$ . This yields a new high quality mesh  $\mathcal{T}'$  adapted to the new subdomain  $\Omega'$  (Figure 1.8c).
- 

Finally, let us mention the existence of the *Deformable Simplicial Complex* (DSC) method for shape and topology optimization [99, 98], which is an alternative method for evolving explicit meshed representations of shapes. The essential principle of the method is to translate vertices according to the deformation field, and to perform elementary mesh operations to repair invalid or poor quality elements. The method is substantially different, in that no adaptive refinement of the meshed interface  $\Gamma$  is performed according to a surface metric. Furthermore, topological changes are not dictated by the level set advection but are the consequences of the remeshing operation. Let us notice that as an essential user constraint, the remeshing software `mng` implementing the method of [25] forbids itself to alter the topology of the zero isosurface of an input level set function to mesh.

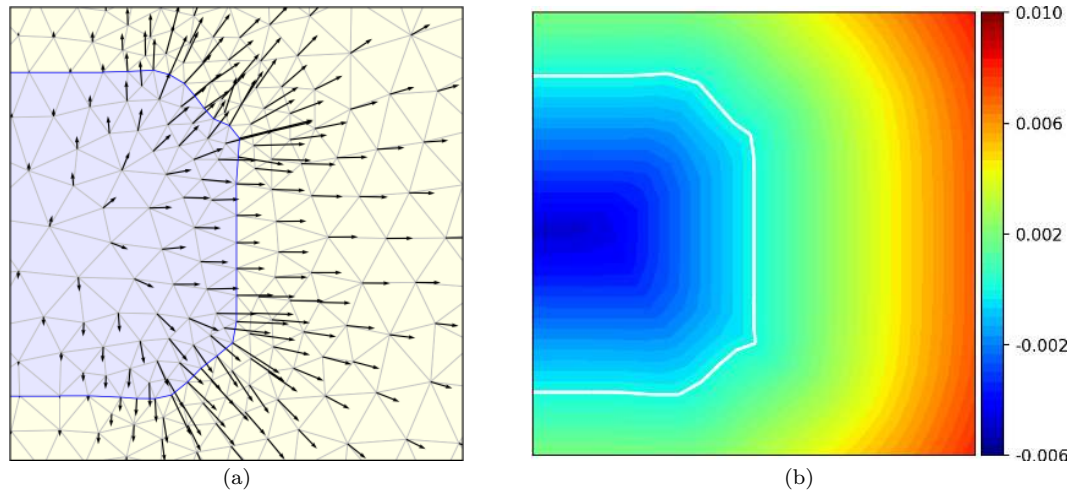


Figure 1.7: Initial setting of the level set based mesh evolution algorithm of [24]. (a) A vector field  $\theta$  is defined at the vertices of the computational mesh  $\mathcal{T}$  for the background domain  $D$  in which  $\Omega \subset D$  is explicitly; (b) the signed distance function  $d_\Omega$  is computed on the mesh  $\mathcal{T}$  as a level set  $\phi = d_\Omega$  representing  $\Omega$ .

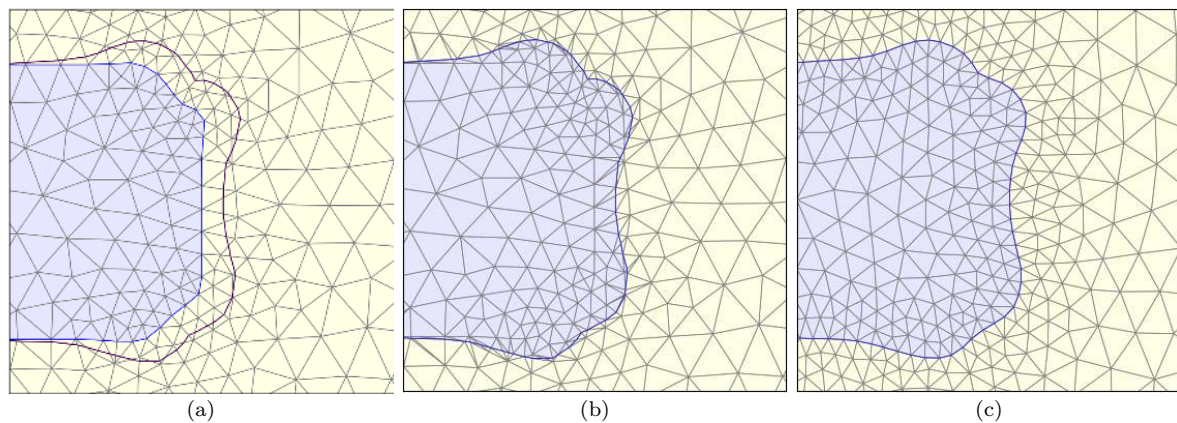


Figure 1.8: Mesh evolution algorithm: (b) the level set function  $\phi = d_\Omega$  associated with the subdomain  $\Omega$  is advected on  $\mathcal{T}$ ; (c) A poor quality mesh  $\tilde{\mathcal{T}}$  of the evolved domain  $\Omega'$  is obtained by splitting  $\mathcal{T}$  according to the updated level set function ; (d)  $\tilde{\mathcal{T}}$  is iteratively remeshed into a new mesh  $\mathcal{T}'$  of sufficient quality for finite element analysis.



## CHAPTER 2

# HADAMARD'S SHAPE DERIVATIVES FOR A COUPLED THERMAL FLUID STRUCTURE PROBLEM

### Contents

---

<b>2.1</b>	<b>Introduction</b>	<b>61</b>
<b>2.2</b>	<b>Setting of the three-physic problem</b>	<b>64</b>
2.2.1	Steady state incompressible Navier-Stokes equation for the fluid variable	65
2.2.2	Convection diffusion for the temperature variable	66
2.2.3	Thermoelasticity with fluid structure interaction for the elastic variable	67
2.2.4	Shape optimization setting	67
<b>2.3</b>	<b>Shape derivatives for a simplified scalar fluid structure interaction problem</b>	<b>68</b>
2.3.1	Presentation of the simplified problem and of its variational formulation	68
2.3.2	A fully Lagrangian setting for computing shape derivatives of arbitrary objective functionals	70
<b>2.4</b>	<b>Shape derivatives for the three-physic problem</b>	<b>76</b>
<b>2.5</b>	<b>Numerical test cases</b>	<b>80</b>
2.5.1	A few details about the numerical implementation	80
2.5.2	Cantilever beam in linearized elasticity	81
2.5.3	Optimal shapes for pure heat conduction	82
2.5.4	Optimal drag profiles for Stokes and Navier-Stokes flows	84
2.5.5	Minimum compliance problem in thermoelasticity	87
2.5.6	A steady-state fluid-structure interaction problem	89
2.5.7	Convective heat transfer	92
2.5.8	Optimization of a compliant thermoelastic solid with fluid-structure interaction	96
<b>2.6</b>	<b>Appendix</b>	<b>100</b>
2.6.1	Proof of propositions 2.3 and 2.4	100
2.6.2	Calculating the shape derivative of a particular objective functional and its adjoint system with C�ea's method	104

---

*Note* : most of the content of this chapter has been published in [153]:

F. FEPPON, G. ALLAIRE, F. BORDEU, J. CORTIAL, AND C. DAPOGNY, *Shape optimization of a coupled thermal fluid-structure problem in a level set mesh evolution framework*, SeMA Journal, (2019), pp. 1–46.

However, the introduction has been rewritten in order to fit the outline of the thesis, and the numerical section 2.5 has been substantially improved and enlarged with additional test cases.

### 2.1 INTRODUCTION

One of the ultimate motivations guiding this thesis is the optimization of mechanical structures subjected to thermal loads and cooled down by fluids. In this chapter, we investigate shape and topology optimization based on the method of Hadamard for a weakly coupled model of heat propagation, fluid flow and structure deformation. This model is detailed in section 2.2 and serves as a basis for all the numerical shape optimization test cases considered in this thesis; it is based on the incompressible Stokes or Navier Stokes equations in the fluid domain, on the convection-diffusion equation for heat propagation in both fluid and solid domains, and on the linearized thermoelasticity system for the mechanical displacement of the solid domain. Our main result is the calculation of shape derivatives of *arbitrary* objective functionals in volume and surface form; these are given in propositions 2.3 and 2.4 respectively. These formulas are implemented and verified numerically on 2-d test cases in section 2.5 with the level set based mesh evolution method of [25] (summarized in chapter 1, section 1.4.2). Our 2-d numerical test cases involve either one, two or three of the aforementioned physics simultaneously; more complex numerical test cases including 3-d problems being specifically detailed in chapter 4.

The main originality of the present work is the consideration of coupled heat propagation, fluid flow and structure deformations for topology optimization with the method of Hadamard. There are few previous contributions [318] using such complete models for shape and topology optimization, and they all rely on very different methods for parametrizing shapes and topologies, namely SIMP [65, 66] or variable density methods [72]. Some simpler variants of this full three-physics model have been more extensively studied. For example, there are quite many works dealing with convective heat transfer problems (involving coupled fluid and thermal equations, without coupled elasticity) with density based methods [224, 116, 134, 324, 275, 129, 118, 258, 263] or variants of level set methods [9, 104, 315] (not explicitly based on the method of Hadamard). Fluid-structure interacting systems (without taking into account thermal effects) have also been considered and in slightly fewer works [319, 240, 42, 194, 221, 208] featuring only 2-d results, as well as thermoelastic structures [197, 285, 151, 314, 101, 28].

Likewise, there is a relatively small amount of work considering level set methods for topology optimization of fluid systems [315, 194, 133] (which do not rely explicitly on the method of Hadamard), although its use in the context of geometry optimization (in which the shape is explicitly meshed but its topology is fixed) is a very classical problem, see the monographs [191, 260, 236, 261] and the more recent work [112].

From a numerical point of view, one specificity of our contribution is the implementation of the Hadamard method with the level set based remeshing algorithm of Allaire, Dapogny and Frey [24] (reviewed in chapter 1, section 1.4.2). This allows—in sharp contrast with density and more classical level set methods—to optimize such coupled systems without introducing a relaxed formulation for the description of a mixture of the fluid and solid domains. In the context of density methods (using namely SIMP for linear elasticity or the Brinkman penalization approach for the Navier-Stokes equations [72]), the entire domain  $D$  is assumed to be filled with a porous material containing a volume fraction  $\eta_f(x)$  of fluid and  $1 - \eta_f(x)$  of solid at every point  $x$ . Then the state equations posed in the fluid domain  $\Omega_f$  or in the solid domain  $\Omega_s$  are replaced with an approximate version (inspired from the homogenization theory) to be solved on the whole domain  $D = \Omega_s \cup \Omega_f$  with coefficients depending on the volume fraction  $\eta_f$  (if not on a local microstructure, see [18]). The nature of the approximation may depend on the physical model and the type of boundary condition considered at the interface  $\Gamma = \partial\Omega_s \cap \partial\Omega_f$  (we discuss this point further in chapter 6, section 6.1.3); for instance the classical approximation of the standard linear elasticity system with constant Hooke's tensor  $A$ ,

$$\begin{cases} -\operatorname{div}(A\nabla\mathbf{u}) = \mathbf{f}_s \text{ in } \Omega_s \\ A\nabla\mathbf{u} \cdot \mathbf{n} = 0 \text{ on } \partial\Omega_s \end{cases}$$

is an elasticity system with variable Hooke's tensor  $\tilde{A}(\eta_f)$ ,

$$-\operatorname{div}(\tilde{A}(\eta_f)\nabla\mathbf{u}) = \mathbf{f}_s \text{ in } D.$$

Likewise, the Stokes problem

$$\begin{cases} -\Delta\mathbf{v} + \nabla p = \mathbf{f}_f \text{ in } \Omega_f \\ \operatorname{div}(\mathbf{v}) = 0 \text{ in } D \\ \mathbf{v} = 0 \text{ on } \partial\Omega_f \end{cases}$$

is commonly relaxed as a Brinkman model with variable inverse permeability  $\tilde{\alpha}(\eta_f)$  [72]

$$\begin{cases} -\Delta\mathbf{v} + \nabla p + \tilde{\alpha}(\eta_f)\mathbf{v} = \mathbf{f} \text{ in } D, \\ \operatorname{div}(\mathbf{v}) = 0 \text{ in } D. \end{cases}$$

The relaxed Hooke's tensor  $\tilde{A}$  and the inverse permeability  $\tilde{\alpha}$  must vary between respectively  $\tilde{A}(0) = A$  and  $\tilde{\alpha}(0) \simeq +\infty$  in the pure solid part corresponding to  $\eta_f = 0$ , and  $\tilde{A}(1) \simeq 0$  and  $\tilde{\alpha}(1) = 0$  in the pure fluid part corresponding to  $\eta_f = 1$ . A common difficulty characterizing these density approaches is the need for determining adequate interpolation laws in order to penalize intermediate densities and to obtain in practice convergence towards “black and white” designs (featuring only  $\eta_f = 0$  or  $\eta_f = 1$  in the computational domain  $D$ ). One law must be proposed per homogenized coefficient, which makes it delicate to set when considering complex multiphysics problems involving several or many of these coefficients.

On the contrary, our level set approach keeps distinct fluid and solid domains throughout the optimization process,  $\Omega_f$  and  $\Omega_s$ , and makes the fluid-solid interface  $\Gamma$  the main optimization variable. The

method of Hadamard allows to estimate the sensitivity of the original physics with respect to local variations of the sharp interface  $\Gamma$ . Keeping track of a clear and neat interface with the use of remeshing avoids various problems such as mass conservation, load transmission or boundary layers description [281, 276]. In our context, the level set based mesh evolution algorithm of [24, 25] alleviates these difficulties while enabling shapes to change their topology in the course of optimization iterations. In numerical practice, the solid-fluid interface  $\Gamma$  is evolved with the level set method and remeshed at every iteration for the finite element analysis. This is particularly well suited for fluid-solid applications, as no alteration of the physical equations is required in contrast with density methods or fixed mesh level set methods, requiring ersatz materials. Another advantage is that this allows, in principle, to solve for each physics non intrusively, for instance with industrial solvers external to the shape optimization implementation.

The chapter is organized as follows. Section 2.2 details the considered weakly coupled model. The coupling is weak in the sense that it allows to solve for each physical equations successively and independently from the knowledge of already computed state variables. The main simplifying hypothesis is that mechanical deformations and displacements are small. In particular, the fluid domain is fixed at first order; it is independent of the deformation of the structure. This allows to solve, first, the incompressible Stokes or Navier-Stokes equations in the fluid domain. Second, we solve a convection-diffusion equation for the temperature in the fluid and solid domains with the previously computed velocity. Third, the linearized thermoelasticity system is solved for the mechanical displacement of the solid domain: the forces at play are a combination of applied external loads, of the fluid stress prescribed on the fluid-solid interface, and of the thermal dilation induced by the temperature. This weak coupling is of course a major simplification and it dramatically reduces the computational cost since no monolithic coupled system has to be solved. Let us note that despite these simplifications, this model is sufficiently general to include previous classical studies when considering only a subset of these physics separately.

The next two sections 2.3 and 2.4 are concerned with the derivations of Hadamard shape sensitivities formulas which are at the basis of our gradient-based topology optimization algorithm used in our numerical simulations. The first section 2.3 introduces a simplified model made of two scalar Laplace equations, mimicking the weak coupling between the solid and the fluid subdomains. This section is purely pedagogical: it deals with the same mathematical difficulties which could be somewhat hidden behind the relative complexity of the full multiphysics system treated in section 2.4. Two major contributions are especially highlighted.

First, we explain how to adapt the classical shape differentiation methods in order to consider *arbitrary* shape functionals. Indeed, as reviewed in chapter 1, section 1.2.3 the most commonly used method is the so-called Lagrangian method of C ea [84, 17]. It brings simplifications with respect to the ‘‘rigorous calculation’’ (see chapter 1, section 1.2.1) but it relies on the specific knowledge of a formula for the objective function. In other words, if the objective function changes (say from a volume integral to a surface integral), the calculation of the shape derivative has to be performed again from the beginning. Here, following the lead of Murat and Simon [242], we rather rely on a Lagrangian approach which allows us to treat very general objective functions, without precise formulas, under a mild assumption on the existence of some partial derivatives. Obviously, it is in the more complex framework of the full three-physic problem that our proposed Lagrangian approach is really more efficient than the classical C ea’s method.

Second, we detail the treatment of the coupling induced by the fluid on the structure, which features a quite surprising mathematical phenomenon. Classically, the order of the coupling is reversed for the adjoint system: for the direct problem, one first solves the fluid equation, and in a second step the solid equation; for the adjoint problem, it is the opposite, namely elasticity is solved first, followed by fluid mechanics. While the fluid and solid equations are coupled by a one-sided interface transmission condition of Neumann type in the direct problem (which accounts for the force exerted by the fluid on the solid), the adjoint equations are, on the contrary, coupled by a Dirichlet interface condition. This ‘transformation’ of the Neumann interface condition for the direct problem into a Dirichlet interface condition for the adjoint seems new to the best of our knowledge.

These derivations are explained in details in section 2.3 on the simplified model, before being extended to the full multiphysics setting in section 2.4 for our practical applications.

Eventually, 2-d applicative examples are considered in section 2.5 for the numerical validation of our derived analytical shape derivatives. Several optimizations are performed on test cases featuring either one, two, or three physics enabled simultaneously. Constrained optimization problems are solved in these situations with the help of a constrained optimization algorithm detailed in the dedicated chapter 3. Our preliminary test cases are drawn from classical literature with only one physics involved:

we consider the classical cantilever beam in structural mechanics, an optimal heat conduction test case, and classical optimal drag profile designs in fluid mechanics. We then address problems featuring two physics simultaneously. Our first example is a fluid-structure interaction problem, without taking into accounts thermal effects. It is inspired from a similar example in [319, 221]. Our second test case is a convective heat transfer problem where the elastic deformation of the solid phase is neglected. It was already studied in [225] with a different approach (based on a variable density relaxation of the problem). Our third test case is concerned with optimal design in thermoelasticity and is based on the previous work [314]. Finally, we present a last test case which comprises the three physics and is new, to the best of our knowledge. All three examples work nicely in the sense that the objective function is indeed minimized while respecting optimization constraints constraints, and that the obtained optimal shapes are significantly different from the initial ones sharing improved performances. These results are preliminary to the more complex heat transfer and 3-d applications considered in [chapter 6](#).

## 2.2 SETTING OF THE THREE-PHYSIC PROBLEM

Let  $D = \Omega_s \cup \Omega_f$  be a fixed, open bounded domain in  $\mathbb{R}^d$  ( $d = 2$  or  $3$  in applications), arising as the disjoint reunion of a ‘fluid’ phase  $\Omega_f$  and a ‘solid’ phase  $\Omega_s = D \setminus \overline{\Omega_f}$  (see [Figure 2.1](#)), separated by an interface  $\Gamma := \partial\Omega_f \cap \partial\Omega_s$  which is to be optimized. Throughout this chapter, the normal vector  $\mathbf{n}$  to  $\Gamma$  is pointing outward the fluid domain  $\Omega_f$ .

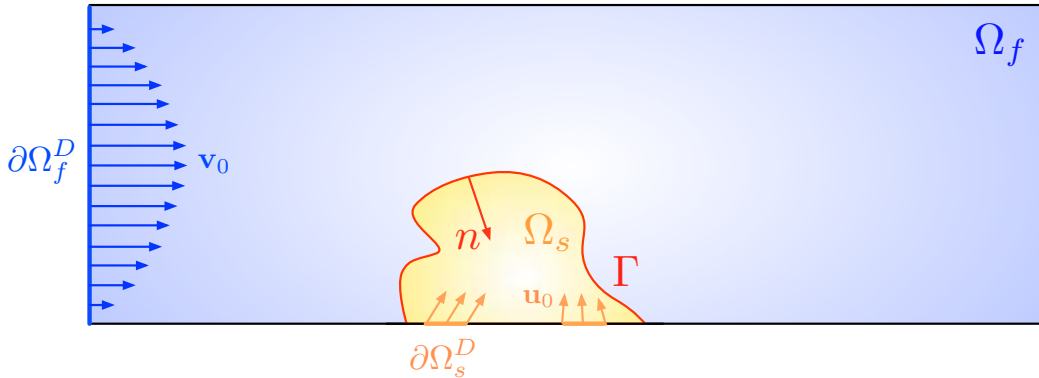


Figure 2.1: Setting of the considered three-physic problem (not all the regions of  $\partial D$  featured in (2.2.1) (2.2.2) and (2.2.3) are represented, see the numerical [section 2.5](#) for more complete settings).

The domain  $D$  is described by three physical variables which are governed by three coupled models:

- the motion of the fluid inside  $\Omega_f$  described by the velocity and pressure fields  $\mathbf{v}$  and  $p$ ;
- the diffusion of heat inside the whole domain  $D$ , and its transport by convection in the fluid domain, resulting in a temperature field  $T$ ;
- the deformation of the solid region  $\Omega_s$  as a result of the stress exerted by the fluid part and of the dilation induced by thermoelastic effects, characterized by a mechanical displacement  $\mathbf{u}$ .

The physical equations chosen for the modeling of  $(\mathbf{v}, p)$ ,  $T$  and  $\mathbf{u}$  with their relevant set of boundary conditions are described in strong form in [sections 2.2.1](#) to [2.2.3](#). The shape optimization setting considered with the method of Hadamard is then introduced in [section 2.2.4](#). Notation for the setting of boundary conditions as well as for all the physical parameters involved in the state equations are summarized respectively in [Tables 2.1](#) and [2.2](#) below.



$D = \Omega_s \cup \Omega_f$	Whole domain featuring fluid and solid phases
$\Omega_f$	Fluid phase
$\Omega_s$	Solid phase
$\Gamma = \partial\Omega_s \cup \partial\Omega_f$	Fluid-structure interface
$\partial D = \partial\Omega_T^N \cup \partial\Omega_T^D$	Boundary of the global domain $D$
$\partial\Omega_T^D$	Dirichlet (isothermal) boundary condition for the temperature variable ( $T = T_0$ on $\partial\Omega_T^D$ )
$\partial\Omega_T^N$	Neumann (adiabatic or isoflux) boundary condition for the temperature variable ( $-k_f \partial T_f / \partial \mathbf{n} = h$ on $\partial\Omega_T^N \cap \partial\Omega_f$ and $-k_s \partial T_s / \partial \mathbf{n} = h$ on $\partial\Omega_T^N \cap \partial\Omega_s$ )
$\partial\Omega_s = \partial\Omega_s^D \cup \partial\Omega_s^N \cup \Gamma$	Boundary of the solid domain $\Omega_s$
$\partial\Omega_s^D$	Dirichlet boundary for the solid variable $\mathbf{u}$
$\partial\Omega_s^N$	Neumann boundary for the solid variable $\mathbf{u}$
$\partial\Omega_f = \partial\Omega_f^D \cup \partial\Omega_f^N \cup \Gamma$	Boundary of the fluid domain $\Omega_f$
$\partial\Omega_f^D$	Dirichlet (inlet) boundary for the fluid variable ( $\mathbf{v} = \mathbf{v}_0$ on $\partial\Omega_f^D$ )
$\partial\Omega_f^N$	Neumann (outlet) boundary for the fluid variable ( $\sigma_f(\mathbf{v}, p) \cdot \mathbf{n} = 0$ on $\partial\Omega_f^N$ )

Table 2.1: Domains  $D$ ,  $\Omega_s$  and  $\Omega_f$  with their associated boundary conditions

### 2.2.1 Steady state incompressible Navier-Stokes equation for the fluid variable

The fluid domain  $\Omega_f$  is filled with a fluid characterized by its viscosity  $\nu$  and density  $\rho$ ; its velocity  $\mathbf{v}$  and pressure  $p$  satisfy the incompressible steady-state Navier-Stokes equations:

$$\left\{ \begin{array}{ll} -\operatorname{div}(\sigma_f(\mathbf{v}, p)) + \rho \nabla \mathbf{v} \mathbf{v} = \mathbf{f}_f & \text{in } \Omega_f \\ \operatorname{div}(\mathbf{v}) = 0 & \text{in } \Omega_f \\ \mathbf{v} = \mathbf{v}_0 & \text{on } \partial\Omega_f^D \\ \sigma_f(\mathbf{v}, p) \mathbf{n} = 0 & \text{on } \partial\Omega_f^N \\ \mathbf{v} = 0 & \text{on } \Gamma, \end{array} \right. \quad (2.2.1)$$

where  $\nabla \mathbf{v}$  is the Jacobian matrix  $(\nabla \mathbf{v})_{ij} = \partial_j v_i$ ,  $i, j = 1, \dots, d$ . The fluid stress tensor  $\sigma_f(\mathbf{v}, p)$  is given by the Newton law

$$\sigma_f(\mathbf{v}, p) = 2\nu e(\mathbf{v}) - pI \quad \text{with } e(\mathbf{v}) = (\nabla \mathbf{v} + \nabla \mathbf{v}^T)/2,$$

where  $I$  is the  $d \times d$  identity matrix.

**Remark 2.1.** The nonlinear term  $\nabla \mathbf{v} \mathbf{v}$  is often written  $(\mathbf{v} \cdot \nabla) \mathbf{v}$  in the literature, however this convention is not consistent with our notation  $\mathbf{v} \cdot \nabla \mathbf{v} := \mathbf{v}^T \nabla \mathbf{v}$  (section 1.1). We shall therefore keep the notation  $\nabla \mathbf{v} \mathbf{v}$  considered to be the matrix product of the Jacobian matrix  $\nabla \mathbf{v}$  with the vector  $\mathbf{v}$ .

In (2.2.1),  $\mathbf{f}_f$  is an applied body force (e.g. gravity); the boundary of the fluid phase is the disjoint reunion

$$\partial\Omega_f = \partial\Omega_f^D \cup \partial\Omega_f^N \cup \Gamma$$

of a Dirichlet (or inlet) part  $\partial\Omega_f^D$  where the flow enters with a given velocity  $\mathbf{v} = \mathbf{v}_0$ , a Neumann (or outlet) part  $\partial\Omega_f^N$  where zero normal stress is observed, and the interface  $\Gamma$  with the solid domain  $\Omega_s$ . At this stage it is assumed that the deformation of the solid domain is sufficiently small so that no slip boundary conditions hold on  $\Gamma$ :  $\mathbf{v} = 0$ ; see remark 2.2 below about this point. Therefore, the variables  $(\mathbf{v}, p)$  depend solely on the geometry of the fluid domain  $\Omega_f$ .

$\mathbf{v} \in H^1(\Omega_f, \mathbb{R}^d)$	Fluid velocity
$p \in L^2(\Omega_f)$	Fluid pressure
$\nu \in \mathbb{R}_+^*$	Fluid viscosity
$\rho \in \mathbb{R}_+^*$	Fluid density
$\mathbf{f}_f \in H^1(\Omega_f, \mathbb{R}^d)$	Volume force in the fluid phase
$\mathbf{v}_0 \in H^{1/2}(\partial\Omega_f, \mathbb{R}^d)$	Input inlet or outlet velocity
$T \in H^1(D)$	Temperature field
$T_s \in H^1(\Omega_s)$	Restriction of $T$ to the solid domain
$T_f \in H^1(\Omega_f)$	Restriction of $T$ to the fluid domain
$c_p \in \mathbb{R}_+^*$	Heat capacity of the fluid
$k_f \in \mathbb{R}_+^*$	Conductivity in fluid
$k_s \in \mathbb{R}_+^*$	Conductivity in solid
$Q_f \in L^2(\Omega_f)$	External input heat flux in the fluid domain
$Q_s \in L^2(\Omega_s)$	External input heat flux in the solid domain
$T_0 \in H^{1/2}(\partial D)$	Input temperature on the total boundary
$h \in L^2(\partial\Omega_T^N)$	Input entering heat flux on the total boundaries
$\mathbf{u} \in H^1(\Omega_s, \mathbb{R}^d)$	Elastic displacement
$\mu \in \mathbb{R}_+^*$	Lamé coefficient
$\lambda \in \mathbb{R}_+^*$	Lamé coefficient
$T_{\text{ref}} \in \mathbb{R}_+^*$	Reference temperature in the solid
$\alpha \in \mathbb{R}_+^*$	Thermal expansion coefficient
$\mathbf{f}_s \in H^1(\Omega_s, \mathbb{R}^d)$	Volume force in the solid phase
$\mathbf{u}_0 \in H^{1/2}(\partial\Omega_s^D, \mathbb{R}^d)$	Prescribed displacement
$\mathbf{g} \in L^2(\partial\Omega_s^N)$	Input traction force

Table 2.2: Physical parameters considered in the weakly coupled model

### 2.2.2 Convection diffusion for the temperature variable

The fluid velocity  $\mathbf{v}$  determines the physical behavior of the temperature  $T$  in the whole domain  $D$ , as a result of convection and diffusion effects inside the fluid domain  $\Omega_f$ , and of pure diffusion inside the solid domain  $\Omega_s$ . Denoting by  $k_f$  and  $k_s$  the thermal conductivity inside  $\Omega_f$  and  $\Omega_s$  respectively, and by  $c_p$  the thermal capacity of the fluid, the temperature field  $T$  is determined by the steady-state

convection-diffusion equations:

$$\left\{ \begin{array}{ll} -\operatorname{div}(k_f \nabla T_f) + \rho c_p \mathbf{v} \cdot \nabla T_f = Q_f & \text{in } \Omega_f \\ -\operatorname{div}(k_s \nabla T_s) = Q_s & \text{in } \Omega_s \\ T = T_0 & \text{on } \partial\Omega_T^D \\ -k_f \frac{\partial T_f}{\partial \mathbf{n}} = h & \text{on } \partial\Omega_T^N \cap \partial\Omega_f \\ -k_s \frac{\partial T_s}{\partial \mathbf{n}} = h & \text{on } \partial\Omega_T^N \cap \partial\Omega_s \\ T_f = T_s & \text{on } \Gamma \\ -k_f \frac{\partial T_f}{\partial \mathbf{n}} = -k_s \frac{\partial T_s}{\partial \mathbf{n}} & \text{on } \Gamma, \end{array} \right. \quad (2.2.2)$$

where we use the subscripts  $f$  and  $s$  for the restrictions  $T_f$  and  $T_s$  of  $T$  to  $\Omega_f$  and  $\Omega_s$  respectively. In (2.2.2),  $Q_f$  and  $Q_s$  are volumic sources inside  $\Omega_f$  and  $\Omega_s$ ; the boundary  $\partial D = \partial\Omega_T^N \cup \partial\Omega_T^D$  is split into a Dirichlet part, where a temperature  $T_0$  is imposed and a Neumann part where a given incoming heat flux  $h$  is applied on  $\partial\Omega_T^N$ . The temperature  $T$  as well as the heat flux are continuous across the interface  $\Gamma$  between  $\Omega_f$  and  $\Omega_s$ .

### 2.2.3 Thermoelasticity with fluid structure interaction for the elastic variable

Finally, the fluid variables  $(\mathbf{v}, p)$  and the temperature  $T$  together determine the displacement  $\mathbf{u}$  of the solid domain  $\Omega_s$ , which is assumed to be an isotropic thermoelastic material with Lamé coefficients  $\lambda, \mu$ , thermal expansion parameter  $\alpha$  and temperature at rest  $T_{\text{ref}}$ . This variable  $\mathbf{u}$  is characterized by the equations of linear thermoelasticity:

$$\left\{ \begin{array}{ll} -\operatorname{div}(\sigma_s(\mathbf{u}, T_s)) = \mathbf{f}_s & \text{in } \Omega_s \\ \mathbf{u} = \mathbf{u}_0 & \text{on } \partial\Omega_s^D \\ \sigma_s(\mathbf{u}, T_s)\mathbf{n} = \mathbf{g} & \text{on } \partial\Omega_s^N \\ \sigma_s(\mathbf{u}, T_s)\mathbf{n} = \sigma_f(\mathbf{v}, p) \cdot \mathbf{n} & \text{on } \Gamma. \end{array} \right. \quad (2.2.3)$$

where the solid stress tensor is given by

$$\sigma_s(\mathbf{u}, T_s) = A e(\mathbf{u}) - \alpha(T_s - T_{\text{ref}})I \text{ with } A e(\mathbf{u}) = 2\mu e(\mathbf{u}) + \lambda \operatorname{Tr}(e(\mathbf{u}))I, \quad (2.2.4)$$

and  $\mathbf{f}_s$  is an applied body force. In (2.2.3), the boundary  $\partial\Omega_s$  is split into respectively a Dirichlet part  $\partial\Omega_s^D$  where a displacement  $\mathbf{u} = \mathbf{u}_0$  is prescribed, a Neumann part  $\partial\Omega_s^N$  where a stress  $\mathbf{g}$  is imposed, and the interface  $\Gamma$  with the fluid domain. On this latter boundary, the solid is submitted to the pressure force imposed by the fluid, which translates into the equality  $\sigma_f(\mathbf{v}, p) \cdot \mathbf{n} = \sigma_s(\mathbf{u}, T_s) \cdot \mathbf{n}$  between the normal fluid and solid stresses.

**Remark 2.2.** The above model is a simplified version of a genuine fluid-solid-thermic coupling between the phases  $\Omega_f$  and  $\Omega_s$ . A more accurate description of fluid-structure interaction would feature a vanishing fluid velocity  $\mathbf{v}$  on the *deformed* interface  $(I + \mathbf{u})(\Gamma)$ , namely:

$$\mathbf{v}(x + \mathbf{u}(x)) = 0, \quad x \in \Gamma, \quad (2.2.5)$$

see e.g. [267], or [319] in an optimization context. Likewise, the equality between normal stresses  $\sigma_s(\mathbf{u}, T_s)\mathbf{n} = \sigma_f(\mathbf{v}, p)\mathbf{n}$  should hold on the deformed interface. In the present work, the displacement  $\mathbf{u}$  of the solid phase is assumed to be small enough so that the influence of the interface deformation on the physical behavior of the fluid can be neglected. Thanks to this simplification, the system (2.2.1) to (2.2.3) is only weakly coupled: its resolution is achieved by solving first the fluid system (2.2.1), then using the resulting fluid velocity  $\mathbf{v}$  in the heat transfer equation (2.2.2), and finally using the fluid stress  $\sigma_f(\mathbf{v}, p)$  and the temperature  $T_s$  to solve (2.2.3).

### 2.2.4 Shape optimization setting

The final goal of the thesis is the resolution of general constrained optimization problems of the form

$$\begin{array}{l} \min_{\Gamma} \quad J(\Gamma, \mathbf{v}(\Gamma), p(\Gamma), T(\Gamma), \mathbf{u}(\Gamma)) \\ \text{s.t.} \quad \begin{cases} g_i(\Gamma, \mathbf{v}(\Gamma), p(\Gamma), T(\Gamma), \mathbf{u}(\Gamma)) = 0, & 1 \leq i \leq p, \\ h_j(\Gamma, \mathbf{v}(\Gamma), p(\Gamma), T(\Gamma), \mathbf{u}(\Gamma)) \leq 0, & 1 \leq j \leq q, \end{cases} \end{array} \quad (2.2.6)$$

where  $J$ ,  $g_i$  and  $h_j$  are *arbitrary* shape objective and constraint functionals set by a user. For its resolution, we shall rely on the null space gradient flow algorithm detailed in [chapter 3](#), which (like any other first order optimization method) requires the knowledge of the shape derivatives of the above functionals. In the context of the method of Hadamard (reviewed in [chapter 1, section 1.2](#)), this means computing the Fréchet derivative of the mapping

$$\boldsymbol{\theta} \mapsto J(\Gamma_{\boldsymbol{\theta}}, \mathbf{v}(\Gamma_{\boldsymbol{\theta}}), p(\Gamma_{\boldsymbol{\theta}}), T(\Gamma, \boldsymbol{\theta}), \mathbf{u}(\Gamma_{\boldsymbol{\theta}})),$$

where we recall that  $\boldsymbol{\theta}$  is a deformation field belonging in general to  $W^{1,\infty}(D, \mathbb{R}^d)$  and  $\Gamma_{\boldsymbol{\theta}} = (I + \boldsymbol{\theta})\Gamma$  (and so on for the constraints  $g_i$  and  $h_j$ ). The setting in the biphasic domain  $D = \Omega_f \cup \Omega_s$  is illustrated on [Figure 2.2](#).

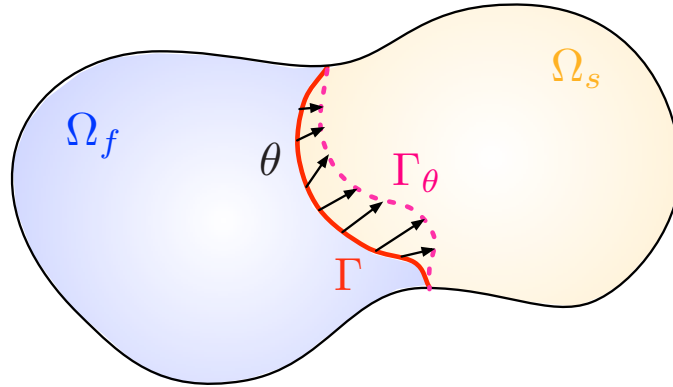


Figure 2.2: Deformation of a partition  $D = \Omega_f \cup \Omega_s$  using Hadamard's method.

For simplicity, we assume in our context that only the interface  $\Gamma$  is subject to optimization, and not the boundary  $\partial D$  of the total domain, therefore the admissible space for deformations  $\boldsymbol{\theta}$  is not  $W^{1,\infty}(D, \mathbb{R}^d)$  but rather the subspace:

$$W_0^{1,\infty}(D, \mathbb{R}^d) = \{\boldsymbol{\theta} \in L^\infty(D, \mathbb{R}^d) \mid \nabla \boldsymbol{\theta} \in L^\infty(D, \mathbb{R}^d \times \mathbb{R}^d) \text{ and } \boldsymbol{\theta} = 0 \text{ on } \partial D\}. \quad (2.2.7)$$

The condition  $\boldsymbol{\theta} = 0$  on  $\partial D$  implies that junction points corresponding to the intersection  $\partial D \cap \Gamma$  are fixed. The analysis to follow can of course be extended to more general shape variations (*e.g.* only  $\boldsymbol{\theta} \cdot \mathbf{n} = 0$  on  $\partial D$  or even  $\boldsymbol{\theta} \in W^{1,\infty}(D, \mathbb{R}^d)$ ), at the expense of including extra terms in shape derivatives accounting for the variations of boundary conditions or for the tangential displacements of the junction points [114].

### 2.3 SHAPE DERIVATIVES FOR A SIMPLIFIED SCALAR FLUID STRUCTURE INTERACTION PROBLEM

One of the salient features of the system (2.2.1) to (2.2.3) is the equality of normal stresses  $\sigma_s(\mathbf{u}, T_s)\mathbf{n} = \sigma_f(\mathbf{v}, p)\mathbf{n}$  imposed on the optimized interface  $\Gamma$ . The calculation of shape derivatives in this context is not a completely standard issue to the best of our knowledge, and the result has an interesting structure: the interface conditions for the adjoint systems are different from those appearing in the state equations (see the expressions (2.4.8) and (2.4.10)). To illustrate this fact, we first consider a simplified scalar problem which gathers with lighter notations all the essential mathematical points.

This preliminary study contains as a particular case the classical Poisson problem considered in [chapter 1, \(1.2.19\)](#), which will also allow us to highlight how the derivation of shape derivatives can be adapted, in quite general contexts, for arbitrary objective functionals.

#### 2.3.1 Presentation of the simplified problem and of its variational formulation

We consider the following setting: the fluid variable  $\mathbf{v}$  is replaced by a scalar variable  $u_f$  solving the Poisson equation (2.3.1) in  $\Omega_f$  with homogeneous Dirichlet boundary conditions on the interface  $\Gamma$ . The elastic variable  $\mathbf{u}$  is replaced by the solution  $u_s$  to another Poisson equation (2.3.2) posed in the complement  $\Omega_s$ . Body sources  $f_f \in H^1(\Omega_f)$ ,  $f_s \in H^1(\Omega_s)$ , with the dimension of a force, are applied in

both cases. Note that such  $H^1$  regularity is needed when computing shape derivatives in the sequel (see the discussion of [chapter 1, section 1.2.3](#)), although weaker regularity is enough for state equations to be well-posed. For the sake of physical units ( $u_f$  and  $u_s$  have respectively the dimension of a velocity and of a displacement), the Laplace operators involved in [\(2.3.1\)](#) and [\(2.3.2\)](#) are scaled with constants  $\nu$  and  $\mu$  that assume respectively the role of a viscosity in the fluid domain, and of a Lamé coefficient in the solid domain. The two systems for  $u_f$  and  $u_s$  are weakly coupled: the equation [\(2.3.1\)](#) for  $u_f$  does not depend on  $u_s$  and in the system [\(2.3.2\)](#) for  $u_s$ , equality of normal fluxes is imposed at the interface  $\Gamma$ , as the exact counterpart in the scalar setting to the continuity of normal stresses  $\sigma_s(\mathbf{u}, T_s)\mathbf{n} = \sigma_f(\mathbf{v}, p)\mathbf{n}$ . For simplicity, homogeneous Dirichlet boundary conditions are imposed on the boundary  $\partial D$  of the total domain  $D = \Omega_f \cup \Omega_s$  in both systems for  $u_f$  and  $u_s$ :

$$\begin{cases} -\nu\Delta u_f = f_f & \text{in } \Omega_f \\ u_f = 0 & \text{on } \Gamma \\ u_f = 0 & \text{on } \partial\Omega_f \setminus \Gamma, \end{cases} \quad (2.3.1)$$

$$\begin{cases} -\mu\Delta u_s = f_s & \text{in } \Omega_s \\ \mu \frac{\partial u_s}{\partial \mathbf{n}} = \nu \frac{\partial u_f}{\partial \mathbf{n}} & \text{on } \Gamma \\ u_s = 0 & \text{on } \partial\Omega_s \setminus \Gamma. \end{cases} \quad (2.3.2)$$

We recall our convention that the normal vector  $\mathbf{n}$  in [\(2.3.2\)](#) is pointing outward  $\Omega_f$  (note that the transmission boundary condition in [\(2.3.2\)](#) reads the same if  $\mathbf{n}$  is pointing outward  $\Omega_s$ ). We consider the minimization problem

$$\min_{\Gamma} J(\Gamma, u_f(\Gamma), u_s(\Gamma)), \quad (2.3.3)$$

where  $J$  is a given cost function upon which we shall impose adequate regularity conditions in due time.

**Remark 2.3.** This simplified setting reduces to the single Poisson problem considered in [chapter 1, section 1.2.3](#) with  $\Omega = \Omega_f$ ,  $\Gamma_N = \emptyset$  in the particular case where the objective function  $J(\Gamma, u_f(\Gamma))$  does not depend on the solid variable (in that case [\(2.3.2\)](#) can be ignored). Although in the present context, no Neumann boundary condition is considered for  $u_f$  for simplicity, we shall still retrieve the results of [chapter 1, propositions 1.8 and 1.9](#) because our way to obtain shape derivatives is rather insensitive to the type of boundary conditions considered.

In the sequel, the dependence of the state variables  $u_f$  and  $u_s$  with respect to  $\Gamma$  is made explicit—using the notations  $u_s(\Gamma)$  and  $u_f(\Gamma)$  as in [\(2.3.3\)](#)—when it is needed (e.g. as in [\(2.3.12\)](#) below). In order to discuss the precise mathematical setting for [\(2.3.1\)](#) [\(2.3.2\)](#), the following spaces of functions on the subdomains  $\Omega_s$  and  $\Omega_f$  are introduced:

$$V_s(\Gamma) = \{v_s \in H^1(\Omega_s) \mid v_s = 0 \text{ on } \partial\Omega_s \setminus \Gamma\}, \quad (2.3.4)$$

$$V_f(\Gamma) = \{v_f \in H^1(\Omega_f) \mid v_f = 0 \text{ on } \partial\Omega_f \setminus \Gamma\}, \quad (2.3.5)$$

$$V_{f,0}(\Gamma) = \{v_f \in H^1(\Omega_f) \mid v_f = 0 \text{ on } \partial\Omega_f\}. \quad (2.3.6)$$

We also consider the subspace  $H_{00}^{1/2}(\Gamma)$  of  $H^{1/2}(\Gamma)$  composed of restrictions to  $\Gamma$  of functions in  $V_f(\Gamma)$ ,

$$H_{00}^{1/2}(\Gamma) = \{v|_{\Gamma} \mid v \in V_f(\Gamma)\}, \quad (2.3.7)$$

and its dual space  $H_{00}^{-1/2}(\Gamma)$ . Roughly speaking, any element  $v \in H_{00}^{1/2}(\Gamma)$  has an extension  $v_f$  to  $\Omega_f$  vanishing on  $\partial\Omega_f \setminus \Gamma$  [\[218\]](#).

In this framework, the state variables  $u_s$  and  $u_f$  in [\(2.3.1\)](#) and [\(2.3.2\)](#) are the unique solutions to the following variational problems:

$$\text{Find } u_f \in V_{f,0}(\Gamma) \text{ such that } \forall v_f \in V_{f,0}(\Gamma), \int_{\Omega_f} \nu \nabla u_f \cdot \nabla v_f dx = \int_{\Omega_f} f_f v_f dx, \quad (2.3.8)$$

$$\text{Find } u_s \in V_s(\Gamma) \text{ such that } \forall v_s \in V_s(\Gamma), \int_{\Omega_s} \mu \nabla u_s \cdot \nabla v_s dx = \int_{\Omega_s} f_s v_s dx - \int_{\Gamma} \nu \frac{\partial u_f}{\partial \mathbf{n}} v_s ds, \quad (2.3.9)$$

where the minus sign in the last term of the right-hand side of [\(2.3.9\)](#) is due to our convention whereby  $\mathbf{n}$  is pointing outward  $\Omega_f$ .

A comment is in order about the meaning of (2.3.9). In general, the normal derivative  $\partial v / \partial \mathbf{n}$  on the interface  $\Gamma$  of an arbitrary function  $v \in V_{f,0}(\Gamma)$  is not defined, because there is no trace theorem for the gradient of functions in  $H^1(\Omega_f)$ . However,  $u_f$  is not a mere function of  $H^1(\Omega_f)$ : (2.3.8) implies that  $\Delta u_f = f_f \in L^2(\Omega_f)$ . Therefore, the last term of (2.3.9) is defined, for any  $v \in H_{00}^{1/2}(\Gamma)$ , by

$$-\int_{\Gamma} \nu \frac{\partial u_f}{\partial \mathbf{n}} v ds := -\int_{\Omega_f} (\nu \Delta u_f \tilde{v} + \nu \nabla u_f \cdot \nabla \tilde{v}) dx = \int_{\Omega_f} (f_f \tilde{v} - \nu \nabla u_f \cdot \nabla \tilde{v}) dx, \quad (2.3.10)$$

where  $\tilde{v} \in V_f(\Gamma)$  is any extension of  $v$  to  $\Omega_f$  satisfying  $\tilde{v} = v$  on  $\Gamma$ . Note that, for smooth  $u_f$ , (2.3.10) is just Green's formula. Since from (2.3.8), the right-hand side of (2.3.10) does not depend on the choice of such extension  $\tilde{v}$ , this identity actually defines  $\partial u_f / \partial \mathbf{n}$  as an element of the dual  $H_{00}^{-1/2}(\Gamma)$ ; see e.g. [218, 176].

The variational formulation (2.3.9) associated to (2.3.2) thus rewrites:

Find  $u_s \in V_s(\Gamma)$  such that  $\forall v_s \in V_s(\Gamma)$ ,

$$\int_{\Omega_s} \mu \nabla u_s \cdot \nabla v_s dx = \int_{\Omega_s} f_s v_s dx + \int_{\Omega_f} f_f \tilde{v}_s dx - \int_{\Omega_f} \nu \nabla u_f \cdot \nabla \tilde{v}_s dx \quad (2.3.11)$$

where  $\tilde{v}_s \in V_f(\Gamma)$  is any extension of  $v_s$  to  $\Omega_f$  such that  $\tilde{v}_s = v_s$  on  $\Gamma$ .

**Remark 2.4.** When  $\Gamma \cap \partial D = \emptyset$ , which happens if for instance  $\Omega_s$  is strictly included in  $D$  ( $\Omega_s \subset\subset D$ ), then  $H_{00}^{1/2}(\Gamma)$  and  $H_{00}^{-1/2}(\Gamma)$  coincide with the more usual fractional Sobolev spaces  $H^{1/2}(\Gamma)$  and  $H^{-1/2}(\Gamma)$  respectively; see [218, 302] about these technicalities.

### 2.3.2 A fully Lagrangian setting for computing shape derivatives of arbitrary objective functionals

Although very common and widely used in the literature (see e.g. [32, 253]), an issue with C ea's method as exposed in chapter 1, section 1.2.3 is that the computation of the shape derivatives depend very much on the assumptions made on the nature of the considered objective functional  $J$ . Different type of functionals may lead to different strong forms for the adjoint equations (see section 2.6.2 where this fact is exemplified), which imposes to redo the analytical derivation whenever the objective function is modified, and to update the numerical implementation accordingly.

In this section, we use a fully Lagrangian setting to derive rigorously the shape derivative of very general objective functionals in the simplified setting of section 2.3.1, in the spirit of the seminal work of Murat and Simon [242]. The shape sensitivities of the state variables  $u_f(\Gamma)$ ,  $u_s(\Gamma)$  are calculated *first* in order to obtain the shape derivative of an arbitrary objective functional in volume form. Then, under sufficient regularity assumptions, the well-known Hadamard structure theorem together with suitable integration by parts yield general shape derivative formulas in the classical form of a boundary integral.

#### A modified objective functional and Lagrangian derivative of the state variables

The starting remark is that the functional  $J$ , although appearing naturally in the formulation of the optimization problem (2.3.3) is not so convenient for the mathematical analysis. Indeed, the domain of definition of  $J(\Gamma, \cdot, \cdot)$  is  $V_{f,0}(\Gamma) \times V_s(\Gamma)$ , a functional space which depends on the first argument  $\Gamma$ . In order to address this issue, the classical idea is to work within a Lagrangian framework rather than a Eulerian one. Therefore, we consider a fixed reference interface  $\Gamma$ , and we introduce a modified functional  $\mathfrak{J}$  obtained by "transporting"  $J$  on a fixed space: for any  $(\boldsymbol{\theta}, \hat{u}_f, \hat{u}_s) \in W_0^{1,\infty}(D, \mathbb{R}^d) \times V_{f,0}(\Gamma) \times V_s(\Gamma)$ , we define

$$\mathfrak{J}(\boldsymbol{\theta}, \hat{u}_f, \hat{u}_s) := J(\Gamma_{\boldsymbol{\theta}}, \hat{u}_f \circ (I + \boldsymbol{\theta})^{-1}, \hat{u}_s \circ (I + \boldsymbol{\theta})^{-1}). \quad (2.3.12)$$

Conversely, this allows to rewrite the objective functional  $J$  as:

$$\begin{aligned} J(\Gamma_{\boldsymbol{\theta}}, u_f(\Gamma_{\boldsymbol{\theta}}), u_s(\Gamma_{\boldsymbol{\theta}})) &= \mathfrak{J}(\boldsymbol{\theta}, u_f(\Gamma_{\boldsymbol{\theta}}) \circ (I + \boldsymbol{\theta}), u_s(\Gamma_{\boldsymbol{\theta}}) \circ (I + \boldsymbol{\theta})) \\ &= \mathfrak{J}(\boldsymbol{\theta}, u_{f,\boldsymbol{\theta}}, u_{s,\boldsymbol{\theta}}), \end{aligned} \quad (2.3.13)$$

where (following the notation of chapter 1, section 1.2) we have denoted  $\Gamma_{\boldsymbol{\theta}} := (I + \boldsymbol{\theta})(\Gamma)$ ,  $u_{s,\boldsymbol{\theta}} := u_s(\Gamma_{\boldsymbol{\theta}}) \circ (I + \boldsymbol{\theta})$  and  $u_{f,\boldsymbol{\theta}} := u_f(\Gamma_{\boldsymbol{\theta}}) \circ (I + \boldsymbol{\theta})$  the transported functions on their reference domains  $\Omega_s$  and  $\Omega_f$ . The identity (2.3.13) is the key motivation for introducing  $\mathfrak{J}$ : indeed, as it is classical in shape optimization, the transported functions  $u_{f,\boldsymbol{\theta}}, u_{s,\boldsymbol{\theta}}$  are differentiable with respect to  $\boldsymbol{\theta}$  without additional regularity assumptions [184, 242]. More precisely, the following lemma holds:

**Lemma 2.1.** *The mappings  $\boldsymbol{\theta} \mapsto u_{f,\boldsymbol{\theta}}$  and  $\boldsymbol{\theta} \mapsto u_{s,\boldsymbol{\theta}}$ , from  $W_0^{1,\infty}(D, \mathbb{R}^d)$  into  $V_{f,0}(\Gamma)$  and  $V_s(\Gamma)$  are Fréchet differentiable at  $\boldsymbol{\theta} = 0$  and their Fréchet derivatives  $\dot{u}_f(\boldsymbol{\theta})$  and  $\dot{u}_s(\boldsymbol{\theta})$  in the direction  $\boldsymbol{\theta}$  are the unique solutions to the following variational problems:*

$$\begin{aligned} & \text{Find } \dot{u}_f(\boldsymbol{\theta}) \in V_{f,0}(\Gamma) \text{ such that } \forall v_f \in V_{f,0}(\Gamma), \\ & \int_{\Omega_f} \nu \nabla \dot{u}_f(\boldsymbol{\theta}) \cdot \nabla v_f dx = \int_{\Omega_f} (\operatorname{div}(f_f \boldsymbol{\theta}) v_f + \nu (\nabla \boldsymbol{\theta} + \nabla \boldsymbol{\theta}^T - \operatorname{div}(\boldsymbol{\theta}) I) \nabla u_f \cdot \nabla v_f) dx, \end{aligned} \quad (2.3.14)$$

$$\begin{aligned} & \text{Find } \dot{u}_s(\boldsymbol{\theta}) \in V_s(\Gamma) \text{ such that } \forall v_s \in V_s(\Gamma), \\ & \int_{\Omega_s} \mu \nabla \dot{u}_s(\boldsymbol{\theta}) \cdot \nabla v_s dx = \int_{\Omega_s} (\operatorname{div}(f_s \boldsymbol{\theta}) v_s + (\nabla \boldsymbol{\theta} + \nabla \boldsymbol{\theta}^T - \operatorname{div}(\boldsymbol{\theta}) I) \mu \nabla u_s \cdot \nabla v_s) dx - \int_{\Gamma} \nu \frac{\partial \dot{u}_f(\boldsymbol{\theta})}{\partial \mathbf{n}} v_s ds, \end{aligned} \quad (2.3.15)$$

where  $-\frac{\partial \dot{u}_f(\boldsymbol{\theta})}{\partial \mathbf{n}} \in H_{00}^{-1/2}(\Gamma)$  is defined for any  $v \in H_{00}^{1/2}(\Gamma)$  by:

$$-\int_{\Gamma} \nu \frac{\partial \dot{u}_f(\boldsymbol{\theta})}{\partial \mathbf{n}} v ds = \int_{\Omega_f} (\operatorname{div}(f_f \boldsymbol{\theta}) \tilde{v} + (\nabla \boldsymbol{\theta} + \nabla \boldsymbol{\theta}^T - \operatorname{div}(\boldsymbol{\theta}) I) \nu \nabla u_f \cdot \nabla \tilde{v}) dx - \int_{\Omega_f} \nu \nabla \dot{u}_f(\boldsymbol{\theta}) \cdot \nabla \tilde{v} dx. \quad (2.3.16)$$

for any extension  $\tilde{v} \in V_f(\Gamma)$  of  $v$  such that  $v = \tilde{v}$  on  $\Gamma$ .

*Proof.* The proof is classical, so we content ourselves with a very brief description of the main ideas, which are very similar to those outlined in [chapter 1, section 1.2.3](#). We first perform a change of variables in the variational formulations (2.3.8) and (2.3.11) so that the integrals involved there are written on fixed domains. Taking  $v_f \circ (I + \boldsymbol{\theta})^{-1} \in V_{f,0}(\Gamma_{\boldsymbol{\theta}})$  as a test function in (2.3.8) for arbitrary  $v_f \in V_{f,0}(\Gamma)$ , this yields a variational formulation for  $u_{f,\boldsymbol{\theta}}$ :

$$\forall v_f \in V_{f,0}(\Gamma), \int_{\Omega_f} \nu A(\boldsymbol{\theta}) \nabla u_{f,\boldsymbol{\theta}} \cdot \nabla v_f dx = \int_{\Omega_f} f_f \circ (I + \boldsymbol{\theta}) v_f | \det(I + \nabla \boldsymbol{\theta}) | dx, \quad (2.3.17)$$

where  $A(\boldsymbol{\theta})$  is the  $d \times d$  matrix  $A(\boldsymbol{\theta}) = | \det(I + \nabla \boldsymbol{\theta}) | (I + \nabla \boldsymbol{\theta})^{-1} (I + \nabla \boldsymbol{\theta})^{-T}$ . Note that  $\Gamma$  corresponds to the reference configuration in (2.3.17) while it is the deformed configuration in (2.3.8).

Now, for a given  $v_s \in V_s(\Gamma)$  and any extension  $\tilde{v}_s \in V_f(\Gamma)$  such that  $\tilde{v}_s = v_s$  on  $\Gamma$ , the function  $\tilde{v}_s \circ (I + \boldsymbol{\theta})^{-1} \in V_f(\Gamma_{\boldsymbol{\theta}})$  is an extension of  $v_s \circ (I + \boldsymbol{\theta})^{-1} \in V_s(\Gamma_{\boldsymbol{\theta}})$  satisfying  $\tilde{v}_s \circ (I + \boldsymbol{\theta})^{-1} = v_s \circ (I + \boldsymbol{\theta})^{-1}$  on  $\Gamma_{\boldsymbol{\theta}}$ . Therefore taking  $v_s \circ (I + \boldsymbol{\theta})^{-1}$  as a test function in (2.3.11) and performing a change of variables yields a variational formulation for  $u_{s,\boldsymbol{\theta}}$ :

$$\begin{aligned} \forall v_s \in V_s(\Gamma), \int_{\Omega_s} \mu A(\boldsymbol{\theta}) \nabla u_{s,\boldsymbol{\theta}} \cdot \nabla v_s dx &= \int_{\Omega_s} f_s \circ (I + \boldsymbol{\theta}) | \det(I + \nabla \boldsymbol{\theta}) | dx + \int_{\Omega_f} f_f \circ (I + \boldsymbol{\theta}) \tilde{v}_s | \det(I + \nabla \boldsymbol{\theta}) | dx \\ &\quad - \int_{\Omega_f} \nu A(\boldsymbol{\theta}) \nabla u_{f,\boldsymbol{\theta}} \cdot \nabla \tilde{v}_s dx \end{aligned} \quad (2.3.18)$$

for any extension  $\tilde{v}_s \in V_f(\Gamma)$  satisfying  $\tilde{v}_s = v_s$  on  $\Gamma$ .

Eventually, a classical use of the implicit function theorem, as in [184], reveals that the mappings  $\boldsymbol{\theta} \mapsto u_{f,\boldsymbol{\theta}}$  and  $\boldsymbol{\theta} \mapsto u_{s,\boldsymbol{\theta}}$ , from  $W_0^{1,\infty}(D, \mathbb{R}^d)$  into  $V_{f,0}(\Gamma)$  and  $V_s(\Gamma)$  respectively, are Fréchet differentiable in the neighborhood of  $\boldsymbol{\theta} = 0$ . Eqns. (2.3.14) and (2.3.15) are then simply obtained by differentiating (2.3.17) and (2.3.18) with respect to  $\boldsymbol{\theta}$ .  $\square$

**Remark 2.5.** The functions  $\dot{u}_f(\boldsymbol{\theta})$  and  $\dot{u}_s(\boldsymbol{\theta})$ , defined by [lemma 2.1](#), are the *Lagrangian* derivatives of  $u_f$  and  $u_s$  with respect to variations of  $\Gamma$ , a notion of derivative which is directly compatible with the variational setting of the PDEs (2.3.1) and (2.3.2). Recall that  $\dot{u}_s(\boldsymbol{\theta})$  and  $\dot{u}_f(\boldsymbol{\theta})$  do not coincide with the perhaps more physical ‘Eulerian’ derivatives  $u'_f(\boldsymbol{\theta})$ ,  $u'_s(\boldsymbol{\theta})$ , that are the differentials of the mappings  $\boldsymbol{\theta} \mapsto u_f(\Gamma_{\boldsymbol{\theta}})$  and  $\boldsymbol{\theta} \mapsto u_s(\Gamma_{\boldsymbol{\theta}})$  (without composition by  $(I + \boldsymbol{\theta})$ ); see e.g. [242, 184, 17] and [chapter 1, section 1.2](#).

### Adjoint system and volume expression of the shape derivative

Assuming that the transported objective function  $\mathfrak{J}$  has continuous partial derivatives at

$$(\boldsymbol{\theta}, \hat{u}_s, \hat{u}_f) = (0, u_s(\Gamma), u_f(\Gamma)),$$

equation (2.3.13) and the chain rule imply that the mapping  $\boldsymbol{\theta} \mapsto J(\Gamma_{\boldsymbol{\theta}}, u_f(\Gamma_{\boldsymbol{\theta}}), u_s(\Gamma_{\boldsymbol{\theta}}))$  is differentiable with respect to  $\boldsymbol{\theta}$  and that its derivative at  $\boldsymbol{\theta} = 0$  reads:

$$\frac{d}{d\boldsymbol{\theta}} \left[ J(\Gamma_{\boldsymbol{\theta}}, u_f(\Gamma_{\boldsymbol{\theta}}), u_s(\Gamma_{\boldsymbol{\theta}})) \right] (\boldsymbol{\theta}) = \frac{d}{d\boldsymbol{\theta}} \left[ \mathfrak{J}(\boldsymbol{\theta}, u_f, \boldsymbol{\theta}, u_s, \boldsymbol{\theta}) \right] (\boldsymbol{\theta}) = \frac{\partial \mathfrak{J}}{\partial \boldsymbol{\theta}} (\boldsymbol{\theta}) + \frac{\partial \mathfrak{J}}{\partial \widehat{u}_f} (\dot{u}_f(\boldsymbol{\theta})) + \frac{\partial \mathfrak{J}}{\partial \widehat{u}_s} (\dot{u}_s(\boldsymbol{\theta})). \quad (2.3.19)$$

Remark that, to keep notations as light as possible, we omit the point  $(0, u_f(\Gamma), u_s(\Gamma))$  where the partial derivatives of  $\mathfrak{J}$  are evaluated in (2.3.19) and below. The occurrence of the Lagrangian derivatives  $\dot{u}_s(\boldsymbol{\theta})$  and  $\dot{u}_f(\boldsymbol{\theta})$  in (2.3.19) are now classically eliminated by introducing adequate adjoint states  $p_f \in V_f(\Gamma)$  and  $p_s \in V_s(\Gamma)$ , defined in weak form by the following variational problems:

$$\text{Find } p_s \in V_s(\Gamma) \text{ such that } \forall v_s \in V_s(\Gamma), \int_{\Omega_s} \mu \nabla p_s \cdot \nabla v_s dx = \frac{\partial \mathfrak{J}}{\partial \widehat{u}_s} (v_s), \quad (2.3.20)$$

$$\text{Find } p_f \in V_f(\Gamma) \text{ such that } p_f = p_s \text{ on } \Gamma \text{ and } \forall v_f \in V_{f,0}(\Gamma), \int_{\Omega_f} \nu \nabla p_f \cdot \nabla v_f dx = \frac{\partial \mathfrak{J}}{\partial \widehat{u}_f} (v_f). \quad (2.3.21)$$

Interestingly, the equality between normal derivatives  $\nu \frac{\partial u_f}{\partial \mathbf{n}} = \mu \frac{\partial u_s}{\partial \mathbf{n}}$  featured in the system (2.3.2) for the state variables translates into the rather surprising boundary condition  $p_f = p_s$  on  $\Gamma$  in (2.3.24) for the adjoint variables  $p_f$ . This boundary condition can be obtained in at least two ways. In the appendix, we provide in section 2.6.2 a derivation of this adjoint boundary condition based on a Lagrangian with C ea's method. In section 2.3.2 below, we also provide an equivalent mixed formulation for the state equations (2.3.1) and (2.3.2) which features different spaces for the 'primal' state variables and the 'adjoint' test functions. This difference of functional spaces is at the root of this change of boundary conditions at the interface.

**Remark 2.6.** The variational problems (2.3.20) and (2.3.21) make sense for general objective functions. Let us remark that they may lead to different adjoint equations when written in strong forms, which is the way they are obtained with C ea's method (this is one of the reasons why the derivation of shape derivatives must be repeated from the beginning when changing the type of the objective functional). For instance, if we consider an objective functional depending on  $u_s(\Gamma)$  and  $u_f(\Gamma)$  through a volume integral (see section 2.6.2 below for the complete derivation of shape sensitivities with the method of C ea):

$$J(\Gamma, u_s(\Gamma), u_f(\Gamma)) := \int_{\Omega_f} j_f(u_f) dx + \int_{\Omega_s} j_s(u_s) dx \quad (2.3.22)$$

for two  $\mathcal{C}^2$  functions  $j_s, j_f : \mathbb{R} \rightarrow \mathbb{R}$  with bounded second order derivatives, then

$$\frac{\partial \mathfrak{J}}{\partial \widehat{u}_s} (v_s) = \int_{\Omega_s} j'_s(u_s) v_s dx, \quad \frac{\partial \mathfrak{J}}{\partial \widehat{u}_f} (v_f) = \int_{\Omega_f} j'_f(u_f) v_f dx$$

which yields the following strong form for the adjoint system (2.3.20) and (2.3.21):

$$\begin{cases} -\mu \Delta p_s = j'_s(u_s) \text{ in } \Omega_s \\ \mu \frac{\partial p_s}{\partial \mathbf{n}} = 0 \text{ on } \Gamma \\ p_s = 0 \text{ on } \partial \Omega_s \setminus \Gamma, \end{cases} \quad (2.3.23)$$

$$\begin{cases} -\nu \Delta p_f = j'_f(u_f) \text{ in } \Omega_f \\ p_f = p_s \text{ on } \Gamma \\ p_f = 0 \text{ on } \partial \Omega_f \setminus \Gamma. \end{cases} \quad (2.3.24)$$

However, if  $J(\Gamma, u_f(\Gamma), u_s(\Gamma))$  now depends on the gradient of one of the variables, for instance

$$J(\Gamma, u_f(\Gamma), u_s(\Gamma)) := \frac{1}{2} \int_{\Omega_s} \mu |\nabla u_s|^2 dx,$$

then

$$\frac{\partial \mathfrak{J}}{\partial \widehat{u}_s} (v_s) = \int_{\Omega_s} \mu \nabla u_s \cdot \nabla v_s dx = - \int_{\Omega_s} \mu \Delta u_s v_s dx + \int_{\partial \Omega_s} \mu \frac{\partial u_s}{\partial \mathbf{n}} v_s ds, \quad \frac{\partial \mathfrak{J}}{\partial \widehat{u}_f} (v_f) = 0,$$



which yields a different strong form for the adjoint equation associated with  $p_s$ :

$$\begin{cases} -\mu\Delta p_s = -\mu\Delta u_s & \text{in } \Omega_s \\ \mu \frac{\partial p_s}{\partial \mathbf{n}} = \mu \frac{\partial u_s}{\partial \mathbf{n}} & \text{on } \Gamma \\ p_s = 0 & \text{on } \partial\Omega_s \setminus \Gamma, \end{cases} \quad (2.3.25)$$

$$\begin{cases} -\nu\Delta p_f = 0 & \text{in } \Omega_f \\ p_f = p_s & \text{on } \Gamma \\ p_f = 0 & \text{on } \partial\Omega_f \setminus \Gamma. \end{cases} \quad (2.3.26)$$

The adjoint variables  $p_s$  and  $p_f$  allow to obtain an expression independent of  $\dot{u}_s(\boldsymbol{\theta})$  and  $\dot{u}_f(\boldsymbol{\theta})$  for the shape derivative of  $J$  as in [chapter 1, proposition 1.4](#):

**Proposition 2.1.** *Assume that the transported objective function  $\mathfrak{J}$  given by (2.3.12) has continuous partial derivatives at  $(\boldsymbol{\theta}, u_s, u_f) = (0, u_s(\Gamma), u_f(\Gamma))$ . Then, the mapping  $\boldsymbol{\theta} \mapsto J(\Gamma_{\boldsymbol{\theta}}, u_f(\Gamma_{\boldsymbol{\theta}}), u_s(\Gamma_{\boldsymbol{\theta}}))$ , from  $W_0^{1,\infty}(D, \mathbb{R}^d)$  into  $\mathbb{R}$ , is Fréchet differentiable at  $\boldsymbol{\theta} = 0$  and its derivative reads:*

$$\begin{aligned} \frac{d}{d\boldsymbol{\theta}} \left[ J(\Gamma_{\boldsymbol{\theta}}, u_f(\Gamma_{\boldsymbol{\theta}}), u_s(\Gamma_{\boldsymbol{\theta}})) \right] (\boldsymbol{\theta}) &= \frac{\partial \mathfrak{J}}{\partial \boldsymbol{\theta}} (\boldsymbol{\theta}) \\ &+ \int_{\Omega_f} \left[ \operatorname{div}(f_f \boldsymbol{\theta}) p_f + (\nabla \boldsymbol{\theta} + \nabla \boldsymbol{\theta}^T - \operatorname{div}(\boldsymbol{\theta}) I) \nu \nabla u_f \cdot \nabla p_f \right] dx \\ &+ \int_{\Omega_s} \left[ \operatorname{div}(f_s \boldsymbol{\theta}) p_s + (\nabla \boldsymbol{\theta} + \nabla \boldsymbol{\theta}^T - \operatorname{div}(\boldsymbol{\theta}) I) \mu \nabla u_s \cdot \nabla p_s \right] dx, \end{aligned} \quad (2.3.27)$$

where  $p_f$  and  $p_s$  are the adjoint states defined by (2.3.20) and (2.3.21).

**Remark 2.7.** Since formula (2.3.27) for the shape derivative of  $J$  involves domain integrals, it is called a volume expression of the shape derivative. In the next subsection, it is shown that it can equivalently be written in terms of surface integrals (which are more obviously satisfying the so-called Hadamard structure theorem stated in [chapter 1, proposition 1.1](#)).

*Proof.* We first insert  $v_s = \dot{u}_s(\boldsymbol{\theta})$  and  $v_f = \dot{u}_f(\boldsymbol{\theta})$  in the adjoint equations (2.3.20) and (2.3.21) to obtain:

$$\frac{\partial \mathfrak{J}}{\partial \widehat{u}_s} (\dot{u}_s(\boldsymbol{\theta})) = \int_{\Omega_s} \mu \nabla p_s \cdot \nabla \dot{u}_s(\boldsymbol{\theta}) dx, \quad (2.3.28)$$

$$\frac{\partial \mathfrak{J}}{\partial \widehat{u}_f} (\dot{u}_f(\boldsymbol{\theta})) = \int_{\Omega_f} \nu \nabla p_f \cdot \nabla \dot{u}_f(\boldsymbol{\theta}) dx. \quad (2.3.29)$$

Then taking  $v_s = p_s$  in the variational formulation (2.3.15) for  $\dot{u}_s(\boldsymbol{\theta})$  yields:

$$\int_{\Omega_s} \mu \nabla p_s \cdot \nabla \dot{u}_s(\boldsymbol{\theta}) dx = \int_{\Omega_s} (\operatorname{div}(f_s \boldsymbol{\theta}) p_s + (\nabla \boldsymbol{\theta} + \nabla \boldsymbol{\theta}^T - \operatorname{div}(\boldsymbol{\theta}) I) \mu \nabla u_s \cdot \nabla p_s) dx - \int_{\Gamma} \nu \frac{\partial \dot{u}_f(\boldsymbol{\theta})}{\partial \mathbf{n}} p_s ds. \quad (2.3.30)$$

Now remarking that  $\tilde{v} = p_f \in V_f(\Gamma)$  is an extension of  $v = p_s = p_f \in H_0^{1/2}(\Gamma)$ , taking  $\tilde{v} = p_f$  in (2.3.16) implies:

$$\int_{\Omega_f} \nu \nabla p_f \cdot \nabla \dot{u}_f(\boldsymbol{\theta}) dx = \int_{\Omega_f} (\operatorname{div}(f_f \boldsymbol{\theta}) p_f + (\nabla \boldsymbol{\theta} + \nabla \boldsymbol{\theta}^T - \operatorname{div}(\boldsymbol{\theta}) I) \nu \nabla u_f \cdot \nabla p_f) dx + \int_{\Gamma} \nu \frac{\partial \dot{u}_f(\boldsymbol{\theta})}{\partial \mathbf{n}} p_s ds. \quad (2.3.31)$$

The desired formula (2.3.27) follows then by summation of (2.3.30) and (2.3.31).  $\square$

### Obtaining the surface expression of the shape derivative

The Hadamard's structure theorem ([chapter 1, proposition 1.1](#)) states that under additional regularity on the optimized interface  $\Gamma$  and on the considered vector fields  $\boldsymbol{\theta}$ , the shape derivative of a sufficiently smooth objective function  $J$  depends only on the normal component  $\boldsymbol{\theta} \cdot \mathbf{n}$  on  $\Gamma$ . In this section, we highlight how this remark allows to obtain a surface expression for the shape derivative of  $J$  from the volume expression (2.3.27) in a simple way. To achieve this, we classically rely on three regularity assumptions, which we assume to be satisfied throughout this section:

1. The considered deformations  $\boldsymbol{\theta}$  are smooth, e.g. of class  $\mathcal{C}^1$ ;
2. The state and adjoint variables  $u_s, u_f, p_s, p_f$  enjoy  $H^2$  regularity in their domain of definition; this is for instance the case when  $D, \Omega_s, \Omega_f$  and  $\Gamma$  are smooth enough; see e.g. [76].
3. The partial derivative  $\frac{\partial \mathfrak{J}}{\partial \boldsymbol{\theta}}$  is sufficiently "regular", in the sense that there exist  $\mathbf{f}_{\mathfrak{J}} \in L^1(D, \mathbb{R}^d)$  and  $\mathbf{g}_{\mathfrak{J}} \in L^1(\Gamma, \mathbb{R}^d)$  such that

$$\forall \boldsymbol{\theta} \in W_0^{1,\infty}(D, \mathbb{R}^d), \quad \frac{\partial \mathfrak{J}}{\partial \boldsymbol{\theta}}(\boldsymbol{\theta}) = \int_D \mathbf{f}_{\mathfrak{J}} \cdot \boldsymbol{\theta} dx + \int_{\Gamma} \mathbf{g}_{\mathfrak{J}} \cdot \boldsymbol{\theta} ds. \quad (2.3.32)$$

The uniqueness of the decomposition (2.3.32), when it exists, is straightforward. The existence of such a structure is typically obtained by performing integration by parts on  $\frac{\partial \mathfrak{J}}{\partial \boldsymbol{\theta}}$  using, in turn, the  $H^2$  regularity of the state and adjoint variables  $u_s, u_f, p_s, p_f$ .

In the following and under these assumptions, we denote by

$$\forall \boldsymbol{\theta} \in W_0^{1,\infty}(D, \mathbb{R}^d), \quad \overline{\frac{\partial \mathfrak{J}}{\partial \boldsymbol{\theta}}}(\boldsymbol{\theta}) := \int_{\Gamma} (\mathbf{g}_{\mathfrak{J}} \cdot \mathbf{n})(\boldsymbol{\theta} \cdot \mathbf{n}) ds, \quad (2.3.33)$$

the part of  $\frac{\partial \mathfrak{J}}{\partial \boldsymbol{\theta}}$  that depends only on  $\boldsymbol{\theta} \cdot \mathbf{n}$ . The main result of this section is the following proposition:

**Proposition 2.2.** *Under the above assumptions, the shape derivative (2.3.27) rewrites as an integral over the boundary  $\Gamma$  involving only the normal component  $\boldsymbol{\theta} \cdot \mathbf{n}$ :*

$$\begin{aligned} \frac{d}{d\boldsymbol{\theta}} \left[ J(\Gamma_{\boldsymbol{\theta}}, u_f(\Gamma_{\boldsymbol{\theta}}), u_s(\Gamma_{\boldsymbol{\theta}})) \right] (\boldsymbol{\theta}) &= \overline{\frac{\partial \mathfrak{J}}{\partial \boldsymbol{\theta}}}(\boldsymbol{\theta}) \\ &+ \int_{\Gamma} \left[ f_f p_f - f_s p_s - \nu \nabla u_f \cdot \nabla p_f + \mu \nabla u_s \cdot \nabla p_s + 2\nu \frac{\partial u_f}{\partial \mathbf{n}} \frac{\partial p_f}{\partial \mathbf{n}} - 2\mu \frac{\partial u_s}{\partial \mathbf{n}} \frac{\partial p_s}{\partial \mathbf{n}} \right] (\boldsymbol{\theta} \cdot \mathbf{n}) ds. \end{aligned} \quad (2.3.34)$$

*Proof.* The regularity assumptions allow to perform integration by parts in the volume expression (2.3.27), which yields:

$$\begin{aligned} \frac{d}{d\boldsymbol{\theta}} \left[ J(\Gamma_{\boldsymbol{\theta}}, u_f(\Gamma_{\boldsymbol{\theta}}), u_s(\Gamma_{\boldsymbol{\theta}})) \right] (\boldsymbol{\theta}) &= \int_{\Gamma} \mathbf{g}_{\mathfrak{J}} \cdot \boldsymbol{\theta} ds + \int_{\Gamma} [f_f p_f (\boldsymbol{\theta} \cdot \mathbf{n}) + \nu (\boldsymbol{\theta} \cdot \nabla u_f) (\nabla p_f \cdot \mathbf{n}) + \nu (\mathbf{n} \cdot \nabla u_f) (\nabla p_f \cdot \boldsymbol{\theta}) - \nu (\nabla u_f \cdot \nabla p_f) (\boldsymbol{\theta} \cdot \mathbf{n})] ds \\ &- \int_{\Gamma} [f_s p_s (\boldsymbol{\theta} \cdot \mathbf{n}) + \mu (\boldsymbol{\theta} \cdot \nabla u_s) (\nabla p_s \cdot \mathbf{n}) + \mu (\mathbf{n} \cdot \nabla u_s) (\nabla p_s \cdot \boldsymbol{\theta}) - \mu (\nabla u_s \cdot \nabla p_s) (\boldsymbol{\theta} \cdot \mathbf{n})] ds \\ &+ \int_D \boldsymbol{\Lambda} \cdot \boldsymbol{\theta} dx, \end{aligned} \quad (2.3.35)$$

for some function  $\boldsymbol{\Lambda} \in L^1(D, \mathbb{R}^d)$  which does not need to be written explicitly. Hadamard's structure theorem implies that (2.3.35) must vanish for vector fields  $\boldsymbol{\theta}$  compactly supported in  $\Omega_s$  or in  $\Omega_f$ , or for vector fields  $\boldsymbol{\theta}$  which are tangential to  $\Gamma$ . Therefore,  $\boldsymbol{\Lambda} = 0$ , and one obtains (2.3.34) by computing (2.3.35) with  $\boldsymbol{\theta}$  normal to  $\Gamma$ .  $\square$

Formula (2.3.34) is called a surface expression of the shape derivative of  $J$ . It is retrieved in section 2.6.2 for a particular case of an objective functional  $J$  with Céa's method. An asset of the above Lagrangian method is that it depends neither on the nature of the objective function  $J$  (e.g. if it depends on  $u_f, u_s$  or their gradients), nor on the type of boundary conditions satisfied by the state variables  $u_f$  and  $u_s$ . Both expressions (2.3.27) and (2.3.34) are convenient to implement because they require minimal inputs from the user: namely, the expression of the partial derivatives of  $\mathfrak{J}$  with respect to  $\boldsymbol{\theta}$  (for (2.3.27) and (2.3.34)) and to the state variables  $u_s, u_f$  (for solving the adjoint system (2.3.14) and (2.3.15)).

**Remark 2.8.** Propositions 2.1 and 2.2 respectively reduce to propositions 1.8 and 1.9 in the context of the Poisson problem (1.2.19) considered in chapter 1, by replacing  $u$  and  $p$  by  $u_f$  and  $p_f$  and setting  $u_s = 0, p_s = 0$ . Notice that in the derivation of (2.3.34), we never used the zero Dirichlet boundary condition of  $u_f, p_f$  and  $u_s, p_s$  on respectively  $\partial\Omega_f \setminus \Gamma$  and  $\partial\Omega_s \setminus \Gamma$ . Therefore, it can be seen that (2.3.34) is also valid in the context of proposition 1.9 (featuring Neumann boundary conditions,  $\Gamma_N \neq \emptyset$ ).

**Remark 2.9.** In practice, one can obtain very quickly the surface expression (2.3.34) of the shape derivative from the volume expression (2.3.27) by

1. replacing volume integrals by surface integrals, adding a minus sign when these involve  $\Omega_s$  in order to account the orientation of the normal  $\mathbf{n}$  outward to  $\Omega_f$ ;
2. replacing  $\nabla\boldsymbol{\theta}$  by  $\mathbf{nn}^T(\boldsymbol{\theta}\cdot\mathbf{n})$ .

Let us now illustrate these results with the calculation of the shape derivative of the ‘solid compliance’, which is not of the form (2.3.22) considered previously:

$$J(\Gamma, u_f(\Gamma), u_s(\Gamma)) = \int_{\Omega_s} \mu |\nabla u_s|^2 dx. \quad (2.3.36)$$

The associated transported objective function  $\mathfrak{J}$  via (2.3.12) reads:

$$\mathfrak{J}(\boldsymbol{\theta}, \widehat{u}_f, \widehat{u}_s) = \int_{(I+\boldsymbol{\theta})\Omega_s} \mu |\nabla(\widehat{u}_s \circ (I+\boldsymbol{\theta})^{-1})|^2 dx = \int_{\Omega_s} \mu |(I+\nabla\boldsymbol{\theta})^{-T} \nabla \widehat{u}_s|^2 |\det(I+\nabla\boldsymbol{\theta})| dx. \quad (2.3.37)$$

**Lemma 2.2.** *The functional  $\mathfrak{J}$  defined in (2.3.37) has continuous partial derivatives at*

$$(\boldsymbol{\theta}, \widehat{u}_s, \widehat{u}_f) = (0, u_s(\Gamma), u_f(\Gamma)),$$

which are given by:

$$\frac{\partial \mathfrak{J}}{\partial \widehat{u}_f}(v_f) = 0, \quad \frac{\partial \mathfrak{J}}{\partial \widehat{u}_s}(v_s) = \int_{\Omega_s} 2\mu \nabla u_s \cdot \nabla v_s dx, \quad (2.3.38)$$

$$\frac{\partial \mathfrak{J}}{\partial \boldsymbol{\theta}}(\boldsymbol{\theta}) = \int_{\Omega_s} (\mu |\nabla u_s|^2 \operatorname{div}(\boldsymbol{\theta}) - 2\mu \nabla u_s \cdot \nabla \boldsymbol{\theta} \cdot \nabla u_s) dx, \quad (2.3.39)$$

$$\overline{\frac{\partial \mathfrak{J}}{\partial \boldsymbol{\theta}}}(\boldsymbol{\theta}) = \int_{\Gamma} \left[ -\mu |\nabla u_s|^2 + 2\mu \left| \frac{\partial u_s}{\partial \mathbf{n}} \right|^2 \right] (\boldsymbol{\theta} \cdot \mathbf{n}) ds. \quad (2.3.40)$$

Therefore the solid compliance (2.3.36)  $\boldsymbol{\theta} \mapsto J(\Gamma_{\boldsymbol{\theta}}, u_s(\Gamma_{\boldsymbol{\theta}}), u_f(\Gamma_{\boldsymbol{\theta}}))$  is differentiable with respect to  $\boldsymbol{\theta} \in W_0^{1,\infty}(D, \mathbb{R}^d)$ , and the shape derivative is given by [proposition 2.1](#) in volume form or [proposition 2.2](#) in surface form.

**Remark 2.10.** Note that in this particular case the adjoint state  $p_s$  defined by (2.3.20) satisfies  $p_s = 2u_s$ , but no such property holds for  $p_f$  (the problem is not self-adjoint). Unlike the usual case of the standard Poisson equation, there is no obvious way to write the solid compliance (2.3.36) in a form (2.3.22) involving only volumic integrals in  $u_s, u_f$  without occurrence of their gradient as in [section 2.6.2](#). For example, the formulation as the work of external forces,

$$J(\Gamma, u_s(\Gamma), u_f(\Gamma)) = \int_{\Omega_s} f_s u_s dx - \int_{\Gamma} \nu \frac{\partial u_f}{\partial \mathbf{n}} u_s ds, \quad (2.3.41)$$

involves the gradient of  $u_f$ ; it is therefore less convenient to handle than (2.3.36), since the partial derivative

$$\frac{\partial \mathfrak{J}}{\partial u_f} : v \mapsto - \int_{\Gamma} \nu \frac{\partial v}{\partial \mathbf{n}} u_s ds$$

is not defined on  $V_{f,0}$ , but merely in a subspace of smoother functions (a treatment could be considered however, by considering different solution spaces, see e.g. [294, 216]). We shall also discuss this point when calculating the shape derivative of the lift functional in [chapter 6, section 6.2.3](#).

### A mixed variational formulation for the state and adjoint problem

We conclude this section with a remark which may shed some light on the a priori surprising boundary condition  $p_s = p_f$  on  $\Gamma$  for the adjoint systems (2.3.20) and (2.3.21). The key observation is that the systems (2.3.1) and (2.3.2) may be formally described by means of a variational formulation for the couple  $(u_s, u_f)$  which features different functional spaces for the solution and test functions (a so-called

Petrov-Galerkin variational formulation, see e.g. [144]). The latter is simply obtained by summing (2.3.8) and (2.3.11):

$$\text{Find } (u_f, u_s) \in H_0^1(\Omega_f) \times H^1(\Omega_s) \text{ such that } \forall v \in H^1(D),$$

$$\int_{\Omega_s} \mu \nabla u_s \cdot \nabla v dx + \int_{\Omega_f} \nu \nabla u_f \cdot \nabla v dx = \int_{\Omega_s} f_s v dx + \int_{\Omega_f} f_f v dx. \quad (2.3.42)$$

Problem (2.3.42) implicitly encloses the transmission condition  $\mu \frac{\partial u_s}{\partial \mathbf{n}} = \nu \frac{\partial u_f}{\partial \mathbf{n}}$  on  $\Gamma$ , and the need to resort to extensions of test functions defined on  $\Omega_s$  in the variational formulation (2.3.11) for  $u_s$  is reflected here in that the test function  $v$  belongs to  $H^1(D)$ .

As is customary (see for instance section 2.6.2), the adjoint system is obtained by formally taking the linear transpose of the mixed variational formulation (2.3.42) (with a different right-hand side), which exchanges the roles of the functional spaces for the sought solution and the test functions:

$$\text{Find } p \in H^1(D) \text{ such that } \forall (v_f, v_s) \in H_0^1(\Omega_f) \times H^1(\Omega_s),$$

$$\int_{\Omega_s} \nu \nabla p \cdot \nabla v_s dx + \int_{\Omega_f} \nu \nabla p \cdot \nabla v_f dx = \frac{\partial \mathfrak{J}}{\partial \hat{u}_s}(v_s) + \frac{\partial \mathfrak{J}}{\partial \hat{u}_f}(v_f). \quad (2.3.43)$$

The above equation is in turn equivalent to the triangular system (2.3.20) and (2.3.21) for the restrictions  $p_s$  and  $p_f$  of  $p$  on  $\Omega_s$  and  $\Omega_f$ , where the transmission condition  $p_s = p_f$  on  $\Gamma$  is implicitly contained in the requirement that  $p$  be an element of  $H^1(D)$ .

**Remark 2.11.** The above argument is only formal because it is not obvious that the variational problem (2.3.42) be well-posed. It can however be made rigorous by changing the functional spaces featured in there, more precisely, by searching for  $(u_f, u_s)$  in

$$\{(u_f, u_s) \in H_0^1(\Omega_f) \times V_s(\Gamma), \mathbb{1}_{\Omega_f} \nu \nabla u_f + \mathbb{1}_{\Omega_s} \nu \nabla u_s \in H(\text{div}, D)\},$$

where  $\mathbb{1}_{\Omega_f}$  (resp.  $\mathbb{1}_{\Omega_s}$ ) is the characteristic function of  $\Omega_f$  (resp.  $\Omega_s$ ), and by searching for  $v$  in  $L^2(D)$ . It is then possible to prove that the *inf-sup* condition of the Banach-Necas-Babuska theorem holds (see [144]) in the case of this new version of (2.3.42), which guarantees its well-posedness.

## 2.4 SHAPE DERIVATIVES FOR THE THREE-PHYSIC PROBLEM

We now briefly describe how the methodology presented in section 2.3 applies to the weakly coupled, multiphysics system (2.2.1) to (2.2.3). Let us introduce the functional spaces which are required, respectively, for the Navier-Stokes equations

$$V_{v,p}(\Gamma) = \{(\mathbf{w}, q) \in H^1(\Omega_f, \mathbb{R}^d) \times L^2(\Omega_f)/\mathbb{R} \mid \mathbf{w} = 0 \text{ on } \partial\Omega_f\},$$

for the thermal equation

$$V_T(\Gamma) = \{S \in H^1(D) \mid S = 0 \text{ on } \partial\Omega_T^D\},$$

for the thermo-mechanical equations

$$V_{\mathbf{u}}(\Gamma) = \{\mathbf{r} \in H^1(\Omega_s, \mathbb{R}^d) \mid \mathbf{r} = 0 \text{ on } \partial\Omega_s^D\}.$$

Note that, as is customary in the theory of the Navier-Stokes equations, the quotient space  $L^2(\Omega_f)/\mathbb{R}$ , associated to the pressure field, gathers square integrable functions defined up to an additive constant. We consider as well the affine spaces associated to the non-homogeneous Dirichlet boundary data  $\mathbf{v}_0 \in H^{1/2}(\partial\Omega_f^D, \mathbb{R}^d)$ ,  $\mathbf{u}_0 \in H^{1/2}(\partial\Omega_s^D, \mathbb{R}^d)$  and  $T_0 \in H^{1/2}(\partial\Omega_T^D)$  featured in (2.2.1) to (2.2.3):

$$\begin{aligned} \mathbf{v}_0 + V_{v,p}(\Gamma) &= \{(\mathbf{w}, q) \in H^1(\Omega_f, \mathbb{R}^d) \times L^2(\Omega_f)/\mathbb{R} \mid \mathbf{w} = \mathbf{v}_0 \text{ on } \partial\Omega_f^D \text{ and } \mathbf{w} = 0 \text{ on } \Gamma\}, \\ T_0 + V_T(\Gamma) &= \{S \in H^1(D) \mid S = T_0 \text{ on } \partial\Omega_T^D\}, \\ \mathbf{u}_0 + V_{\mathbf{u}}(\Gamma) &= \{\mathbf{r} \in H^1(\Omega_s, \mathbb{R}^d) \mid \mathbf{r} = \mathbf{u}_0 \text{ on } \partial\Omega_s^D\}. \end{aligned}$$

Finally, we shall make use of the trace space

$$H_{00}^{1/2}(\Gamma, \mathbb{R}^d) = \{\mathbf{r}|_{\Gamma} \mid \mathbf{r} \in H^1(\Omega_f, \mathbb{R}^d) \text{ and } \mathbf{r} = 0 \text{ on } \partial\Omega_f \setminus \Gamma\}, \quad (2.4.1)$$

and its dual  $H_{00}^{-1/2}(\Gamma, \mathbb{R}^d)$ . The state variables  $\mathbf{v}, p, T, \mathbf{u}$  are the solutions to the following variational problems: for the Navier-Stokes equations (2.2.1), find  $(\mathbf{v}, p) \in \mathbf{v}_0 + V_{\mathbf{v},p}(\Gamma)$  such that

$$\forall (\mathbf{w}, q) \in V_{\mathbf{v},p}(\Gamma) \quad \int_{\Omega_f} [\sigma_f(\mathbf{v}, p) : \nabla \mathbf{w} + \rho \mathbf{w} \cdot \nabla \mathbf{v} \cdot \mathbf{v} - q \operatorname{div}(\mathbf{v})] dx = \int_{\Omega_f} \mathbf{f}_f \cdot \mathbf{w} dx; \quad (2.4.2)$$

for the thermal equation (2.2.2), find  $T \in T_0 + V_T(\Gamma)$  such that, for any  $S \in V_T(\Gamma)$ ,

$$\int_{\Omega_s} k_s \nabla T \cdot \nabla S dx + \int_{\Omega_f} (k_f \nabla T \cdot \nabla S + \rho c_p S \mathbf{v} \cdot \nabla T) dx = \int_{\Omega_s} Q_s S dx + \int_{\Omega_f} Q_f S dx + \int_{\partial \Omega_f^N} h S ds; \quad (2.4.3)$$

for the thermo-mechanical equations (2.2.3), find  $\mathbf{u} \in \mathbf{u}_0 + V_{\mathbf{u}}(\Gamma)$  such that

$$\forall \mathbf{r} \in V_{\mathbf{u}}(\Gamma), \quad \int_{\Omega_s} \sigma_s(\mathbf{u}, T_s) : \nabla \mathbf{r} dx = \int_{\Omega_s} \mathbf{f}_s \cdot \mathbf{r} dx + \int_{\partial \Omega_s^N} \mathbf{g} \cdot \mathbf{r} ds - \int_{\Gamma} \mathbf{r} \cdot \sigma_f(\mathbf{v}, p) \cdot \mathbf{n} ds. \quad (2.4.4)$$

Let us comment on the well-posedness of the coupled system of variational problems (2.4.2), (2.4.3) and (2.4.4). As in section 2.3, the volumic source terms are assumed to enjoy  $H^1$  regularity in their domain:  $\mathbf{f}_f \in H^1(\Omega_f, \mathbb{R}^d)$ ,  $\mathbf{f}_s \in H^1(\Omega_s, \mathbb{R}^d)$ ,  $Q_f \in H^1(\Omega_f)$ ,  $Q_s \in H^1(\Omega_s)$ . The surface fluxes  $h, \mathbf{g}$  are assumed to belong to  $L^2$  spaces. At first, the classical theory for the Navier-Stokes equation states that (2.4.2) is well-posed as soon as the Reynolds numbers  $\operatorname{Re} := \|\mathbf{v}_0\|_{H^{1/2}(\partial \Omega_f^P, \mathbb{R}^d)} \rho / \nu$  is sufficiently small; see [304].

The variational formulation (2.4.3) of the thermal problem is not well-posed in utter generality because of the lack of coercivity induced by the advection term  $\int_{\Omega_f} \rho c_p S \mathbf{v} \cdot \nabla T dx$  and of the presence of inhomogeneous Dirichlet boundary conditions. However, in usual applications [255, 225], it is customary to impose a Dirichlet boundary condition  $T = T_{0,f}$  at the inlet of the computational domain (i.e. where  $\mathbf{v} \cdot \mathbf{n} < 0$ ) and a Neumann boundary condition  $-k_f \nabla T \cdot \mathbf{n} = 0$  at the outlet ( $\mathbf{v} \cdot \mathbf{n} > 0$ ). This together with the incompressibility condition  $\operatorname{div}(\mathbf{v}) = 0$  is easily shown to imply the coercivity of the bilinear form featured in (2.4.3); see e.g. [75].

Eventually, the well-posedness of the linear elasticity problem (2.4.4) results from the Lax-Milgram theorem, the only subtle point is that, as in section 2.3,  $\sigma_f(\mathbf{v}, p) \cdot \mathbf{n}$  is an element of the dual space of  $H_{00}^{1/2}(\Gamma, \mathbb{R}^d)$ : if  $\mathbf{v}, p$  were regular enough, the following integration by parts would hold true:

$$\forall \mathbf{r} \in V_{\mathbf{u}}(\Gamma), \quad - \int_{\Gamma} \mathbf{r} \cdot \sigma_f(\mathbf{v}, p) \cdot \mathbf{n} ds = \int_{\Omega_f} (-\operatorname{div}(\sigma_f(\mathbf{v}, p)) \cdot \mathbf{r} - \sigma_f(\mathbf{v}, p) : \nabla \mathbf{r}) dx. \quad (2.4.5)$$

Hence the normal stress  $\sigma_f(\mathbf{v}, p) \mathbf{n}$  can be understood mathematically as an element of  $H_{00}^{-1/2}(\Gamma, \mathbb{R}^d)$ , defined by

$$\forall \mathbf{r} \in H_{00}^{1/2}(\Gamma, \mathbb{R}^d), \quad - \int_{\Gamma} \mathbf{r} \cdot \sigma_f(\mathbf{v}, p) \cdot \mathbf{n} ds = \int_{\Omega_f} (\mathbf{f}_f \cdot \tilde{\mathbf{r}} - \rho \tilde{\mathbf{r}} \cdot \nabla \mathbf{v} \cdot \mathbf{v} - \sigma_f(\mathbf{v}, p) : \nabla \tilde{\mathbf{r}} + \tilde{q} \operatorname{div}(\mathbf{v})) dx \quad (2.4.6)$$

for any extension  $(\tilde{\mathbf{r}}, \tilde{q}) \in H^1(\Omega_f, \mathbb{R}^d) \times (L^2(\Omega_f)/\mathbb{R})$  satisfying  $\tilde{\mathbf{r}} = \mathbf{r}$  on  $\Gamma$ . Note that although it is not fully necessary, we consider also an extension  $\tilde{q}$  of the pressure field to maintain a complete analogy with (2.3.10). This turns to be convenient in the calculation of the shape derivative performed in section 2.6.1.

Throughout this section, we assume that the above conditions for the well-posedness of the coupled system of variational problems (2.4.2), (2.4.3) and (2.4.4) are fulfilled.

In the above context, we aim at solving the minimization problem (2.2.6) where the velocity  $\mathbf{v}(\Gamma)$ , pressure  $p(\Gamma)$ , temperature  $T(\Gamma)$  and elastic displacement  $\mathbf{u}(\Gamma)$  are the solutions to (2.2.1) to (2.2.3).

In order to compute shape derivatives with respect to variations of a given interface  $\Gamma$ , we introduce as in section 2.3.2 a transported functional  $\mathfrak{J}$  defined on the fixed functional space

$$W_0^{1,\infty}(D, \mathbb{R}^d) \times H^1(\Omega_f, \mathbb{R}^d) \times (L^2(\Omega_f)/\mathbb{R}) \times H^1(D) \times H^1(\Omega_s, \mathbb{R}^d)$$

by:

$$\begin{aligned} \forall \boldsymbol{\theta} \in W_0^{1,\infty}(D, \mathbb{R}^d), (\hat{\mathbf{v}}, \hat{p}, \hat{T}, \hat{\mathbf{u}}) &\in H^1(\Omega_f, \mathbb{R}^d) \times (L^2(\Omega_f)/\mathbb{R}) \times H^1(D) \times H^1(\Omega_s, \mathbb{R}^d), \\ \mathfrak{J}(\boldsymbol{\theta}, \hat{\mathbf{v}}, \hat{p}, \hat{T}, \hat{\mathbf{u}}) &= J(\Gamma_{\boldsymbol{\theta}}, \hat{\mathbf{v}} \circ (I + \boldsymbol{\theta})^{-1}, \hat{p} \circ (I + \boldsymbol{\theta})^{-1}, \hat{T} \circ (I + \boldsymbol{\theta})^{-1}, \hat{\mathbf{u}} \circ (I + \boldsymbol{\theta})^{-1}). \end{aligned} \quad (2.4.7)$$

The only requirement made on  $J$  is that the associated functional  $\mathfrak{J}$  has continuous partial derivatives at  $(\boldsymbol{\theta}, \hat{\mathbf{v}}, \hat{p}, \hat{T}, \hat{\mathbf{u}}) = (0, \mathbf{v}(\Gamma), p(\Gamma), T(\Gamma), \mathbf{u}(\Gamma))$ . Under this assumption, arguing as in section 2.3.2 and

section 2.6.2, we define the adjoint variables  $\mathbf{w}, q, S, \mathbf{r}$  as follows. The elasticity adjoint variable  $\mathbf{r} \in V_{\mathbf{u}}(\Gamma)$  is the solution of

$$\int_{\Omega_s} A\mathbf{e}(\mathbf{r}) : \nabla \mathbf{r}' dx = \frac{\partial \mathfrak{J}}{\partial \widehat{\mathbf{u}}}(\mathbf{r}') \quad \forall \mathbf{r}' \in V_{\mathbf{u}}(\Gamma). \quad (2.4.8)$$

The thermal adjoint variable  $S \in V_T(\Gamma)$  is the solution of

$$\int_{\Omega_s} k_s \nabla S \cdot \nabla S' dx + \int_{\Omega_f} (k_f \nabla S \cdot \nabla S' + \rho c_p \mathbf{v} \cdot \nabla S') dx = \int_{\Omega_s} \alpha \operatorname{div}(\mathbf{r}) S' dx + \frac{\partial \mathfrak{J}}{\partial \widehat{T}}(S) \quad \forall S' \in V_T(\Gamma). \quad (2.4.9)$$

The fluid adjoint variables  $(\mathbf{w}, q) \in H^1(\Omega_f, \mathbb{R}^d) \times (L^2(\Omega_f)/\mathbb{R})$  are the solution of

$$\begin{aligned} \mathbf{w} = \mathbf{r} \text{ on } \Gamma \text{ and } \forall (\mathbf{w}', q') \in V_{\mathbf{v},p}(\Gamma) \\ \int_{\Omega_f} \left( \sigma_f(\mathbf{w}, q) : \nabla \mathbf{w}' + \rho \mathbf{w} \cdot \nabla \mathbf{w}' \cdot \mathbf{v} + \rho \mathbf{w} \cdot \nabla \mathbf{v} \cdot \mathbf{w}' - q' \operatorname{div}(\mathbf{w}) \right) dx = \\ \int_{\Omega_f} -\rho c_p S \nabla T \cdot \mathbf{w}' dx + \frac{\partial \mathfrak{J}}{\partial (\mathbf{v}', p')}(\mathbf{w}', q'), \end{aligned} \quad (2.4.10)$$

and we recall our convention whereby the point  $(0, \mathbf{v}(\Gamma), p(\Gamma), T(\Gamma), \mathbf{u}(\Gamma))$  where the partial derivatives of  $J$  are evaluated is omitted.

As expected, the cascade dependency  $(\mathbf{v}, q) \rightarrow T \rightarrow \mathbf{u}$  in the state variables (the variable on the right of the arrow depends on that on the left) is reversed into  $\mathbf{r} \rightarrow S \rightarrow (\mathbf{w}, q)$  for the adjoint variables, which reflects the fact that the adjoint system is formally the linearized transpose of the state problem.

**Remark 2.12.** The existence and uniqueness of a solution  $(\mathbf{w}, q, S, \mathbf{r})$  in  $H^1(\Omega_f, \mathbb{R}^d) \times (L^2(\Omega_f)/\mathbb{R}) \times H^1(D) \times H^1(\Omega_s, \mathbb{R}^d)$  to the adjoint system (2.4.8) to (2.4.10) follows from the same considerations as in the case of the state system (2.4.2) to (2.4.4), except when it comes to the *linearized* Navier-Stokes equation (2.4.10). The latter is well-posed provided the Reynolds number is sufficiently small; see [171], Chap. IV about this point.

**Remark 2.13.** Let us consider the case of a particular objective functional of the form

$$J(\Gamma, \mathbf{v}(\Gamma), p(\Gamma), T(\Gamma), \mathbf{u}(\Gamma)) = \int_{\Omega_f} j_f(\mathbf{v}(\Gamma), p(\Gamma), T(\Gamma)) dx + \int_{\Omega_s} j_s(\mathbf{u}(\Gamma), T(\Gamma)) dx, \quad (2.4.11)$$

where  $j_f : \mathbb{R}^d \times \mathbb{R}_p \times \mathbb{R}_T \rightarrow \mathbb{R}$  and  $j_s : \mathbb{R}^d \times \mathbb{R}_T \rightarrow \mathbb{R}$  are smooth and satisfy adequate growth conditions. The adjoint equations (2.4.8) to (2.4.10) rewrite respectively in strong form, for the elasticity system,

$$\begin{cases} -\operatorname{div}(A\mathbf{e}(\mathbf{r})) = \frac{\partial j_s}{\partial \mathbf{u}} & \text{in } \Omega_s \\ \mathbf{r} = 0 & \text{on } \partial\Omega_s^D \\ A\mathbf{e}(\mathbf{r})\mathbf{n} = 0 & \text{on } \partial\Omega_s^N \cup \Gamma, \end{cases}$$

for the thermal equation,

$$\left\{ \begin{array}{ll} -\operatorname{div}(k_s \nabla S_s) = \alpha \operatorname{div}(\mathbf{r}) + \frac{\partial j_s}{\partial T_s} & \text{in } \Omega_s \\ -\operatorname{div}(k_f \nabla S_f) - \rho c_p \mathbf{v} \cdot \nabla S_f = \frac{\partial j_f}{\partial T_f} & \text{in } \Omega_f \\ S = 0 & \text{on } \partial\Omega_T^D \\ k_s \frac{\partial S_s}{\partial \mathbf{n}} = 0 & \text{on } \partial\Omega_T^N \cap \partial\Omega_s \\ k_f \frac{\partial S_f}{\partial \mathbf{n}} + \rho c_p (\mathbf{v} \cdot \mathbf{n}) S_f = 0 & \text{on } \partial\Omega_T^N \cap \partial\Omega_f \\ S_s = S_f & \text{on } \Gamma \\ -k_s \frac{\partial S_s}{\partial \mathbf{n}} = -k_f \frac{\partial S_f}{\partial \mathbf{n}} & \text{on } \Gamma, \end{array} \right.$$

for the Navier-Stokes equations,

$$\left\{ \begin{array}{ll} -\operatorname{div}(\sigma_f(\mathbf{w}, q)) + \rho(\nabla \mathbf{v}^T \mathbf{w} - \nabla \mathbf{w} \mathbf{v}) = -\rho c_p S \nabla T_f + \frac{\partial j_f}{\partial \mathbf{v}} & \text{in } \Omega_f \\ -\operatorname{div}(\mathbf{w}) = \frac{\partial j_f}{\partial p} & \text{in } \Omega_f \\ \mathbf{w} = 0 & \text{on } \partial \Omega_f^D \\ \sigma_f(\mathbf{w}, q) \mathbf{n} + \rho(\mathbf{v} \cdot \mathbf{n}) \mathbf{w} = 0 & \text{on } \partial \Omega_f^N \\ \mathbf{w} = \mathbf{r} & \text{on } \Gamma. \end{array} \right. \quad (2.4.12)$$

Note the ‘‘surprising’’ fact in the linearized adjoint system (2.4.12) for the Navier-Stokes equations that the interface condition for the velocity variable  $\mathbf{w}$  is of Dirichlet type on  $\Gamma$ , while it was of Neumann type for the direct elasticity problem (2.2.3).

A very similar analysis to that of section 2.3 yields the shape derivative of  $J$  in the present physical context; proofs are postponed to section 2.6.1.

**Proposition 2.3.** *Assume that the transported objective function  $\mathfrak{J}$ , defined by (2.4.7), has continuous partial derivatives at  $(\boldsymbol{\theta}, \widehat{\mathbf{v}}, \widehat{p}, \widehat{T}, \widehat{\mathbf{u}}) = (0, \mathbf{v}(\Gamma_\boldsymbol{\theta}), p(\Gamma_\boldsymbol{\theta}), T(\Gamma_\boldsymbol{\theta}), \mathbf{u}(\Gamma_\boldsymbol{\theta}))$ . Then the objective function  $J$ , considered in (2.2.6), is differentiable with respect to  $\boldsymbol{\theta} \in W_0^{1,\infty}(D, \mathbb{R}^d)$  and the derivative reads*

$$\begin{aligned} & \frac{d}{d\boldsymbol{\theta}} \left[ J(\Gamma_\boldsymbol{\theta}, \mathbf{v}(\Gamma_\boldsymbol{\theta}), p(\Gamma_\boldsymbol{\theta}), T(\Gamma_\boldsymbol{\theta}), \mathbf{u}(\Gamma_\boldsymbol{\theta})) \right] (\boldsymbol{\theta}) \\ &= \frac{\partial \mathfrak{J}}{\partial \boldsymbol{\theta}} (\boldsymbol{\theta}) + \int_{\Omega_f} [\mathbf{w} \cdot \operatorname{div}(\mathbf{f}_f \otimes \boldsymbol{\theta}) - (\sigma_f(\mathbf{v}, p) : \nabla \mathbf{w} + \rho \mathbf{w} \cdot \nabla \mathbf{v} \cdot \mathbf{v}) \operatorname{div}(\boldsymbol{\theta})] dx \\ & \quad + \int_{\Omega_f} [\sigma_f(\mathbf{v}, p) : (\nabla \mathbf{w} \nabla \boldsymbol{\theta}) + \sigma_f(\mathbf{w}, q) : (\nabla \mathbf{v} \nabla \boldsymbol{\theta}) + \rho \mathbf{w} \cdot (\nabla \mathbf{v} \nabla \boldsymbol{\theta}) \cdot \mathbf{v}] dx \\ & \quad - \int_{\Omega_s} \operatorname{div}(\boldsymbol{\theta})(k_s \nabla T \cdot \nabla S) dx - \int_{\Omega_f} \operatorname{div}(\boldsymbol{\theta})(k_f \nabla T \cdot \nabla S + \rho c_p (\mathbf{v} \cdot \nabla T) S) dx \\ & \quad + \int_{\Omega_s} k_s (\nabla \boldsymbol{\theta} + \nabla \boldsymbol{\theta}^T) \nabla T \cdot \nabla S dx + \int_{\Omega_f} [k_f (\nabla \boldsymbol{\theta} + \nabla \boldsymbol{\theta}^T) \nabla T \cdot \nabla S + \rho c_p \mathbf{v} \cdot (\nabla \boldsymbol{\theta}^T \nabla T) S] dx \\ & \quad \quad \quad + \int_{\Omega_s} \operatorname{div}(Q_s \boldsymbol{\theta}) S dx + \int_{\Omega_f} \operatorname{div}(Q_f \boldsymbol{\theta}) S dx \\ & \quad + \int_{\Omega_s} [-\operatorname{div}(\boldsymbol{\theta}) \sigma_s(\mathbf{u}, T) : \nabla \mathbf{r} + \sigma_s(\mathbf{u}, T) : (\nabla \mathbf{r} \nabla \boldsymbol{\theta}) + A e(\mathbf{r}) : (\nabla \mathbf{u} \nabla \boldsymbol{\theta}) + \mathbf{r} \cdot \operatorname{div}(\mathbf{f}_s \otimes \boldsymbol{\theta})] dx, \end{aligned} \quad (2.4.13)$$

where  $\mathbf{r}, S, \mathbf{w}, q$  are the adjoint states defined by (2.4.8) to (2.4.10).

**Proposition 2.4.** *If in addition the state and adjoint variables  $\mathbf{v}, T_s, T_f, \mathbf{u}, \mathbf{w}, S, \mathbf{r}$  (resp.  $p, q$ ) have  $H^2$  (resp  $H^1$ ) regularity in their domain of definition, and if the partial derivative  $\frac{\partial \mathfrak{J}}{\partial \boldsymbol{\theta}}$  has a decomposition of the form (2.3.32), then (2.4.13) rewrites as an integral over the interface  $\Gamma$  depending only on the normal component  $\boldsymbol{\theta} \cdot \mathbf{n}$  of  $\boldsymbol{\theta}$ :*

$$\begin{aligned} & \frac{d}{d\boldsymbol{\theta}} \left[ J(\Gamma_\boldsymbol{\theta}, \mathbf{v}(\Gamma_\boldsymbol{\theta}), p(\Gamma_\boldsymbol{\theta}), T(\Gamma_\boldsymbol{\theta}), \mathbf{u}(\Gamma_\boldsymbol{\theta})) \right] (\boldsymbol{\theta}) \\ &= \overline{\frac{\partial \mathfrak{J}}{\partial \boldsymbol{\theta}}} (\boldsymbol{\theta}) + \int_{\Gamma} (\mathbf{f}_f \cdot \mathbf{w} - \sigma_f(\mathbf{v}, p) : \nabla \mathbf{w} + \mathbf{n} \cdot \sigma_f(\mathbf{w}, q) \nabla \mathbf{v} \cdot \mathbf{n} + \mathbf{n} \cdot \sigma_f(\mathbf{v}, p) \nabla \mathbf{w} \cdot \mathbf{n}) (\boldsymbol{\theta} \cdot \mathbf{n}) ds \\ & \quad + \int_{\Gamma} \left( k_s \nabla T_s \cdot \nabla S_s - k_f \nabla T_f \cdot \nabla S_f + Q_f S - Q_s S_s - 2k_s \frac{\partial T_s}{\partial \mathbf{n}} \frac{\partial S_s}{\partial \mathbf{n}} + 2k_f \frac{\partial T_f}{\partial \mathbf{n}} \frac{\partial S_f}{\partial \mathbf{n}} \right) (\boldsymbol{\theta} \cdot \mathbf{n}) ds \\ & \quad + \int_{\Gamma} (\sigma_s(\mathbf{u}, T_s) : \nabla \mathbf{r} - \mathbf{f}_s \cdot \mathbf{r} - \mathbf{n} \cdot A e(\mathbf{r}) \nabla \mathbf{u} \cdot \mathbf{n} - \mathbf{n} \cdot \sigma_s(\mathbf{u}, T_s) \nabla \mathbf{r} \cdot \mathbf{n}) (\boldsymbol{\theta} \cdot \mathbf{n}) ds, \end{aligned} \quad (2.4.14)$$

where  $\overline{\frac{\partial \mathfrak{J}}{\partial \boldsymbol{\theta}}}$  denotes the part of  $\partial \mathfrak{J} / \partial \boldsymbol{\theta}$  that depends only on the normal trace  $\boldsymbol{\theta} \cdot \mathbf{n}$  (see (2.3.33)).

**Remark 2.14.** Formula (2.4.13) is a volume expression of the shape derivative, while formula (2.4.14) is a surface expression of the same derivative. Formula (2.4.14) can be simplified a little by using the

following identities on  $\Gamma$ , which arise as consequences of the boundary conditions featured in (2.2.1) to (2.2.3):

$$\mathbf{n} \cdot \sigma_f(\mathbf{w}, q) \nabla \mathbf{v} \cdot \mathbf{n} = \sigma_f(\mathbf{w}, q) : \nabla \mathbf{v}, \quad (2.4.15)$$

$$\mathbf{n} \cdot \sigma_f(\mathbf{v}, p) \nabla \mathbf{w} \cdot \mathbf{n} - \mathbf{n} \cdot \sigma_s(\mathbf{u}, T_s) \nabla \mathbf{r} \cdot \mathbf{n} = \sigma_f(\mathbf{v}, p) : \nabla \mathbf{w} - \sigma_s(\mathbf{v}, p) : \nabla \mathbf{r}. \quad (2.4.16)$$

The above equations (2.4.13) and (2.4.14) generalize more classical shape derivatives expressions for each of the physics considered individually: the elastic, thermic, and fluid parts coincide with expressions stated in e.g. [32, 168] for the elasticity, [314] for pure thermoelasticity, [253] for the thermal conductivity terms, and [259, 100, 112] for the Navier Stokes equations. However some terms of formula (2.4.14) vanishing for particular instances of the objective function  $J$  (for example,  $Ae(\mathbf{r}) \cdot \mathbf{n} = 0$  for objective functions written as a volume integral depending only on  $\mathbf{u}$  and not on its gradient), may be missing in previous works.

The shape derivative formulas (2.4.13) and (2.4.14) may seem rather complex at first glance, however they can be implemented once and for all. In a practical implementation, they allow then to automate very conveniently the numerical assembly of shape derivatives, since only analytical expressions of the partial derivatives of the objective functionals are required. These are usually obtained very easily in weak form, which turns to be also convenient for the resolution of the adjoint equations (2.4.8) to (2.4.10).

**Remark 2.15.** In the case where the Stokes system is considered instead of the Navier-Stokes equations in (2.2.1), the above formulas hold by setting  $\rho = 0$  in the terms involving the fluid state and adjoint variables  $\mathbf{v}$  and  $\mathbf{w}$ .

## 2.5 NUMERICAL TEST CASES

This section is devoted to the presentation of several 2-d test cases which in particular allow to verify numerically the shape derivatives formulas of propositions 2.3 and 2.4. This verification is based on the assumption that their correct implementation should make objective functions decrease and constraints become gradually satisfied, in accordance with the expected behavior of our null space algorithm for constrained optimization (detailed in chapter 3).

Here, we demonstrate on various multiphysics examples how the previous ideas can be effectively implemented in order to address a wide range of topology optimization problems.

This part deviates from the published work [153] in several aspects. First, supplementary test cases are treated in sections 2.5.2 to 2.5.5. Second, the test case of sections 2.5.6 and 2.5.7 that were originally treated in [153] have been respectively revised (in order to illustrate differences between the use of volume and surface expressions of the shape derivative) or improved (constraints are treated explicitly with our null space gradient flow optimization algorithm).

### 2.5.1 A few details about the numerical implementation

Our numerical implementation follows algorithm 1.1 outlined in chapter 1 with the level set based mesh evolution method of [25] summarized in algorithm 1.2. We rely on the open-source FreeFem++ environment for the resolution of Finite Element problems [183] (see [34] for its use in the context of structural optimization and [112] about its use in the context of fluid flow optimization). Since much more details shall be provided in chapter 6, we content ourselves to provide here only a brief overview of our implementation.

The Navier-Stokes system (2.4.2) is solved by using a mixed formulation, where the space  $V_{\mathbf{v}, p}$  is discretized with  $\mathbb{P}_1$ -bubble  $\times$   $\mathbb{P}_1$  elements. The equations (2.2.2) and (2.2.3) for the temperature  $T$  and the elastic displacement  $\mathbf{u}$  are solved with  $\mathbb{P}_1$  finite elements. In all the considered examples, the Reynolds number and the velocity  $\mathbf{v}$  of the fluid are sufficiently small so that the convergence and stability of our numerical schemes are guaranteed without the need for more complex numerical strategies, e.g. upwinding methods. Note that, when solving the fluid-structure interaction problem (2.2.3), a numerical estimate of the normal stress  $\sigma_f(\mathbf{v}, p)\mathbf{n}$  is used as a boundary load in the discretization of the variational formulation (2.4.4). Naturally, a mixed formulation analogous to (2.3.42) could be implemented for the triplet  $(\mathbf{v}, p, \mathbf{u})$ , which would avoid computing normal derivatives at the boundary.

Both surface and volume expressions of shape derivatives are considered for the computation of a descent direction as described in chapter 1, section 1.4.1. Once the shape derivative  $DJ(\Gamma)$  of the considered objective function  $J$  (or the constraints) is assembled, the identification problem (1.4.4) (also called extension and regularization of the interface velocity) is solved using FreeFEM. We rely on the inner



product (1.2.24) when using the surface expression (2.4.14) for the shape derivative, and on (1.2.25) when using the volume expression (2.4.13). In several cases, we observed without explanation that the use of the volume expression was beneficial, either for obtaining better shapes or smoother convergence curves for the objective function and constraints. A few comparisons illustrating these facts shall be provided hereafter. Note that this matter has been discussed in a number of works, see e.g. [188, 168].

The open-source library `mmg` [108] is used when it comes to the isosurface discretization and quality-oriented remeshing operations outlined in chapter 1, section 1.4.2. Let us mention that the meshes obtained are in general non symmetric even if the discretized shape is symmetric. Since this can affect slightly the quality of the optimized shape, the level set function associated to the shape of the next iteration is symmetrized before remeshing in test cases where symmetry is to be expected (see chapter 6 for the details). This allows to avoid large symmetry loss induced by the accumulation of small numerical errors. The computation of the signed distance function to a meshed domain (see section 1.3.1) is performed by using the open-source algorithm `mshdist`; see [111].

In all our examples, the considered shape and topology optimization problems feature equality or inequality constraints, for instance on the volume of one of the two phases  $\Omega_f, \Omega_s$ . These constraints are gradually enforced and maintained thanks to the null space gradient flow algorithm detailed in chapter 3. For now, let us content ourselves with saying that optimization trajectories are calculated by mean of a gradient flow (a dynamical system based on the gradient of the objective and constraint functions) that is able to detect and to handle equality and inequality constraints.

In the following, we treat a variety of test cases which are all subcases of the full three physics model. The first three examples are benchmark test cases of the literature featuring only one physics. The next three problems involve two physics: pure fluid-structure interaction, thermoelasticity, and heat convection. The final example involves all physics simultaneously.

## 2.5.2 Cantilever beam in linearized elasticity

Our first example is the very classical cantilever beam test case which has been commonly used to illustrate many topology optimization algorithms, e.g. with the Hadamard method [32, 34] or the SIMP method [43]. We consider the classical compliance minimization problem of finding a structure  $\Omega_s \subset D$  enclosed in a domain  $D = [0, 2] \times [0, 1]$  fixed at top and bottom left boundaries and subject to a traction load  $\mathbf{g}$  on the middle of the right-hand boundary (the setting is illustrated on Figure 2.3). The objective

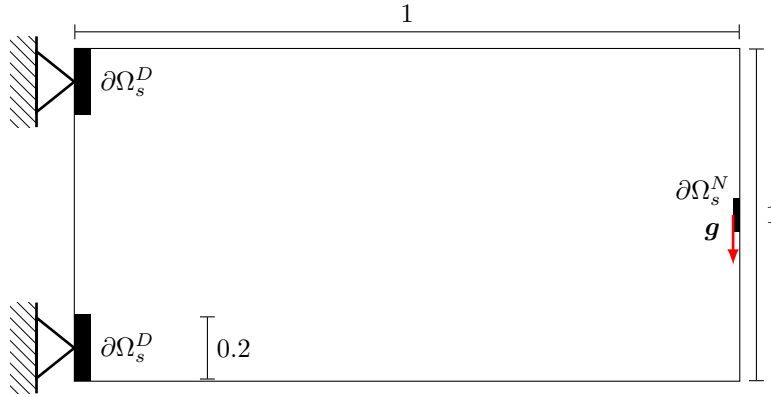


Figure 2.3: Setting of the cantilever optimization problem of section 2.5.2.

is to minimize the compliance of the structure subject to a volume constraint:

$$\begin{aligned} \min_{\Omega_s \subset D} \quad & J(\Omega_s, \mathbf{u}(\Omega_s)) := \int_{\Omega_s} Ae(\mathbf{u}) : e(\mathbf{u}) dx \\ \text{s.t.} \quad & \text{Vol}(\Omega_s) := \int_{\Omega_s} dx = V_{target}. \end{aligned} \quad (2.5.1)$$

where  $V_{target}$  is a target volume set to 0.6 in our implementation.

In this setting, the formulas (2.3.38) and (2.3.39) needed for the shape derivative of  $J$  and the definition of the adjoint systems read:

$$\frac{\partial \mathfrak{J}}{\partial \boldsymbol{\theta}}(\boldsymbol{\theta}) = \int_{\Omega_s} (Ae(\mathbf{u}) : e(\mathbf{u}) \text{div}(\boldsymbol{\theta}) - 2Ae(\mathbf{u}) : (\nabla \mathbf{u} \nabla \boldsymbol{\theta})) dx, \quad (2.5.2)$$

$$\overline{\frac{\partial \mathfrak{J}}{\partial \boldsymbol{\theta}}}(\boldsymbol{\theta}) = \int_{\Gamma} (-Ae(\mathbf{u}) : e(\mathbf{u}) + 2\mathbf{n} \cdot (Ae(\mathbf{u}))\nabla \mathbf{u} \cdot \mathbf{n})(\boldsymbol{\theta} \cdot \mathbf{n}) ds, \quad (2.5.3)$$

$$\frac{\partial \mathfrak{J}}{\partial \widehat{\mathbf{u}}}(\mathbf{r}') = \int_{\Omega_s} 2Ae(\mathbf{u}) : e(\mathbf{r}') dx. \quad (2.5.4)$$

The shape derivative of the volume functional  $\text{Vol}(\Omega_s)$  reads simply

$$\frac{d}{d\boldsymbol{\theta}} \text{Vol}(\Omega_s, \boldsymbol{\theta}) = \int_{\Omega_s} \text{div}(\boldsymbol{\theta}) dx = - \int_{\Gamma} \boldsymbol{\theta} \cdot \mathbf{n} ds, \quad (2.5.5)$$

where the minus sign comes from our convention that  $\mathbf{n}$  is pointing inside  $\Omega_s$ . For this example, the values of the Lamé parameters are set to  $\lambda = 12.96$  and  $\mu = 5.56$ .

Intermediate shapes obtained with our implementation by using the *surface* expression for the shape derivative are displayed on [Figure 2.4](#). An instance of a mesh associated to one of the intermediate shapes where a topological change occurs is visible on [Figure 2.5](#), which further illustrates a key feature of our numerical method: an explicitly meshed, black and white description of the shape is available during the whole optimization process. The history curves for the objective and constraint function values are displayed on [Figure 2.6](#). These illustrate the ability of our optimization algorithm to decrease the objective function while maintaining the constraint satisfied. Observe that very good shapes are obtained within approximately a thirty iterations, and that the oscillations visible for the objective function  $J$  at subsequent iterations are the result of numerical noise around the optimum.

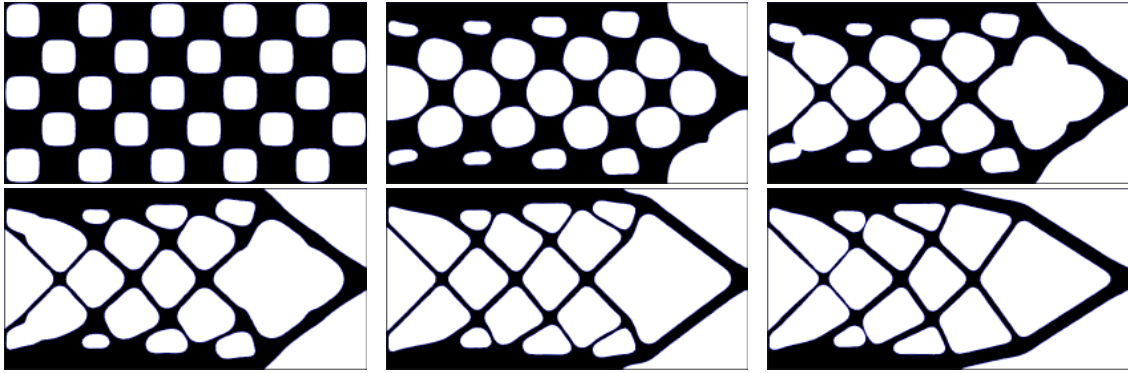


Figure 2.4: Intermediate iterations 0, 8, 15, 19, 30 and 110 for the cantilever test case of [section 2.5.2](#).

### 2.5.3 Optimal shapes for pure heat conduction

We now address a different problem where only thermal effects are considered, for which the only state variable active is the temperature  $T$ . This test case has been treated by Marck and Privat [225] with a density based topology optimization method. We consider a domain  $D = [0.1, 0.1]$  with two phases  $\Omega_s$  and  $\Omega_f$  representing two conductive materials with respective conductivity constants  $k_s = 1$  and  $k_f = 100$ . The whole domain is heated with a uniform volume heat source  $Q_s = Q_f = 10^4$ . The most conductive phase  $\Omega_f$  is connected to a “cold” Dirichlet boundary  $\partial\Omega_T^D$  featuring  $T = T_0$  (with  $T_0 = 0$  in our implementation). All other parts of the boundary are adiabatic (no heat flux escaping the domain). The setting is represented on [Figure 2.7](#).

The objective is to find the distribution of material  $\Omega_f$  which minimizes the average temperature of the whole domain using a limited amount of volume  $V_{target}$  (set to  $V_{target} = 0.15$  in our case):

$$\begin{aligned} \min_{\Gamma} \quad & J(\Gamma, T(\Gamma)) = \int_D T dx \\ \text{s.t.} \quad & \text{Vol}(\Omega_f) \leq V_{target}. \end{aligned} \quad (2.5.6)$$

For this problem, the partial derivatives of the objective function  $J$  read simply

$$\frac{\partial \mathfrak{J}}{\partial \boldsymbol{\theta}}(\boldsymbol{\theta}) = \overline{\frac{\partial \mathfrak{J}}{\partial \boldsymbol{\theta}}}(\boldsymbol{\theta}) = 0 \quad (2.5.7)$$

$$\frac{\partial \mathfrak{J}}{\partial T}(S') = \int_D S' dx. \quad (2.5.8)$$

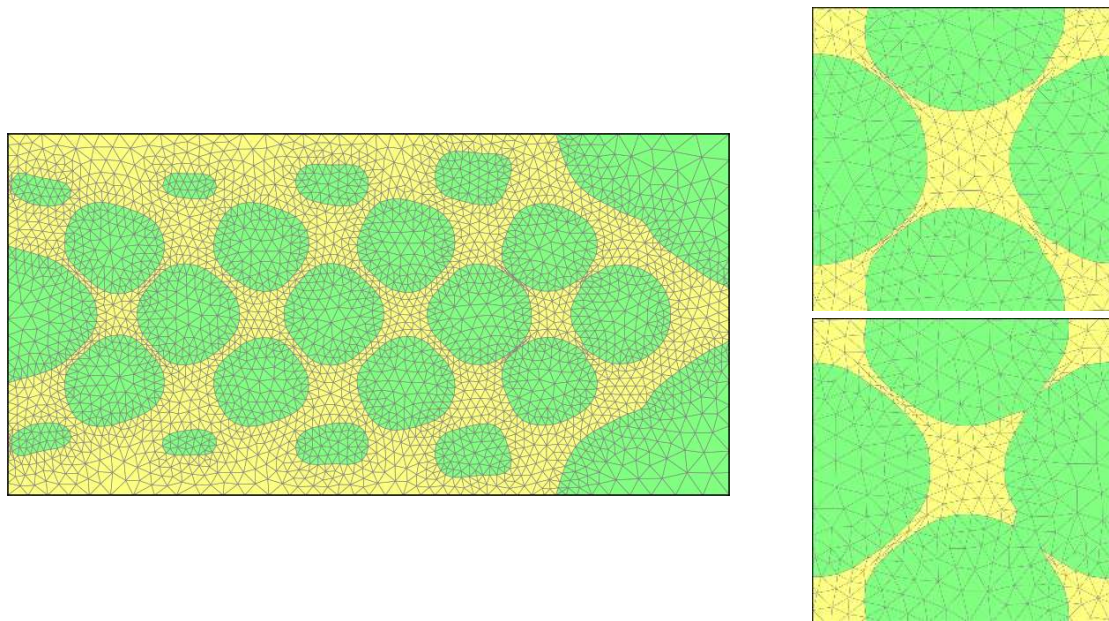


Figure 2.5: Mesh of an intermediate optimization iteration (on the *left*) and zoom on a region featuring a change of topology (on the *right*) for the cantilever test case of section 2.5.2.

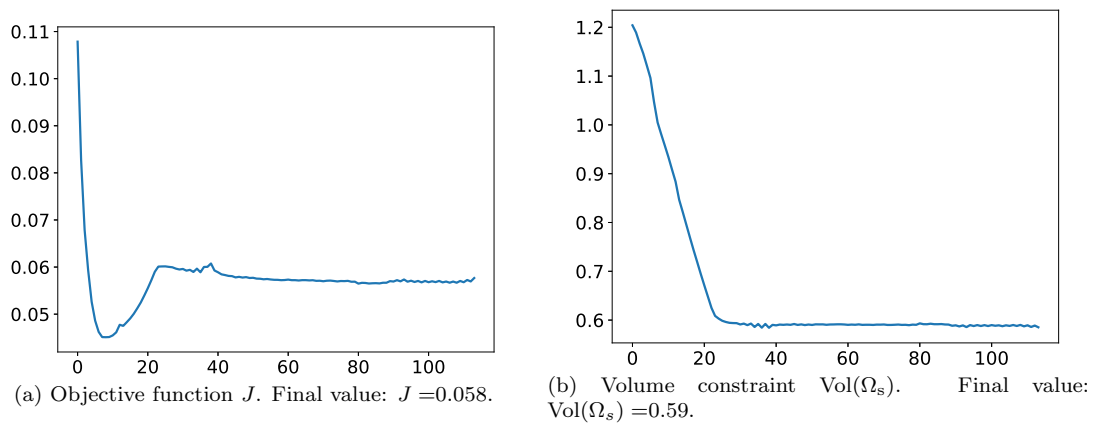


Figure 2.6: Convergence history for the cantilever test case of section 2.5.2.

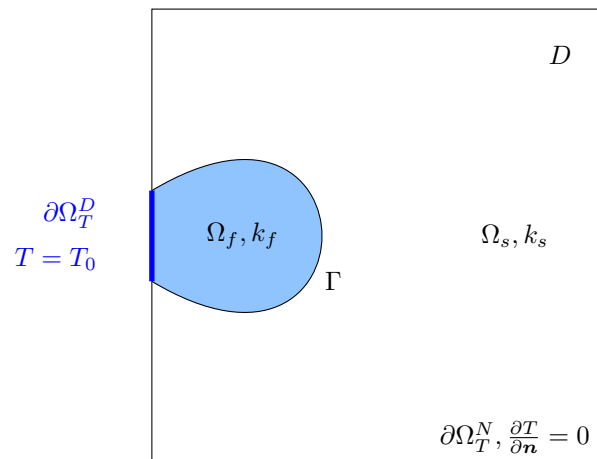


Figure 2.7: Setting of the optimal heat conduction problem of section 2.5.3.

Initial and final shapes obtained by using the surface expression (2.4.14) of the shape derivative are displayed on Figure 2.8. A few intermediate iterations are visible on Figure 2.10. It is interesting to compare them to those obtained with the *volume* expression (2.4.13) of the shape derivative. Both final shapes are similar, although the optimization path is a little bit different with smaller details appearing more quickly with the volume expression of the shape derivative, and a final shape slightly better for a similar number of iterations. The convergence curves for the objective and constraint functionals are plotted on Figure 2.13. It is interesting to observe that although the volume constraint in (2.5.6) is an *inequality* constraint. The process converges without oscillations characterizing many other optimization algorithms using active set strategies [297]. The details of the method allowing to obtain this smoothness of convergence are explained in chapter 3, section 3.4.1.

A mesh of the final shape is shown on Figure 2.12. Let us mention that local parameters were prescribed in the remeshing tool `mmg` in order to obtain a (small) constant mesh edge size for the discretization of the optimized boundary  $\Gamma$ .

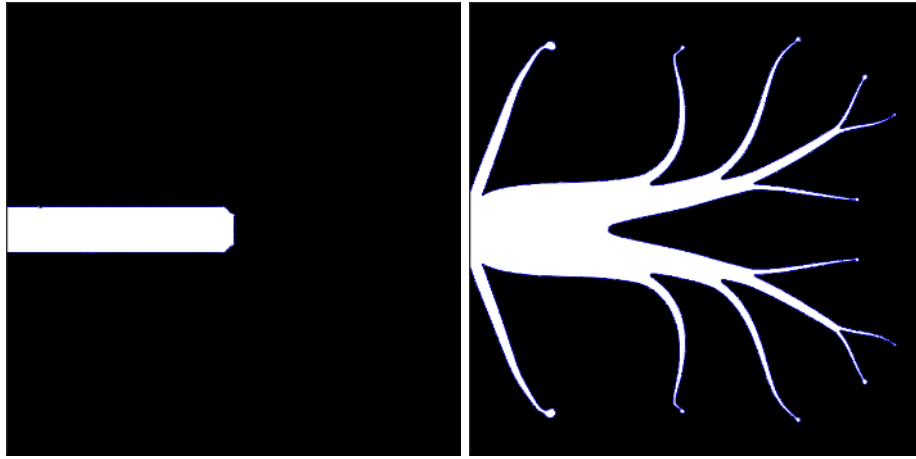


Figure 2.8: Initial and final configurations for the pure heat conduction test case of section 2.5.3. The fluid material is represented in white.

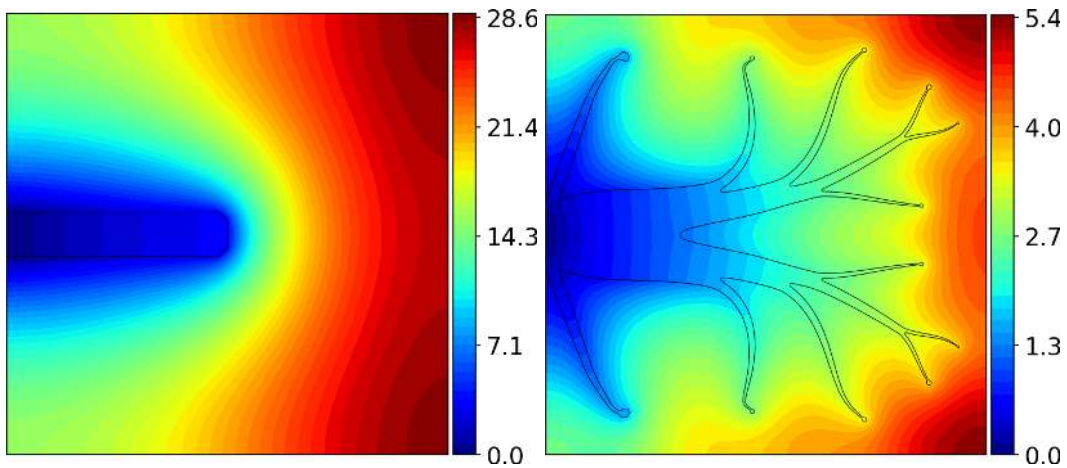


Figure 2.9: Initial and final temperature fields for the pure heat conduction test case of section 2.5.3.

#### 2.5.4 Optimal drag profiles for Stokes and Navier-Stokes flows

The final example featuring only one physics is concerned with the classical problem of finding the shape of an obstacle minimizing the drag force induced by a surrounding flow [259, 172]. Here we follow the setting proposed in [113]. The global domain is a rectangle  $D = [0, 1] \times [0, 1]$ . A flow  $\mathbf{v}$  is entering on the left side with a velocity  $\mathbf{v} = \mathbf{e}_x$  where  $\mathbf{e}_x$  is the unit horizontal direction. The viscosity and density parameters of the flow are given by  $\nu = 5e - 3$  and  $\rho = 1$ , which corresponds to a Reynolds number  $\text{Re} = 200$ . In order to limit the effects induced by the bottom and top walls, a slip boundary condition

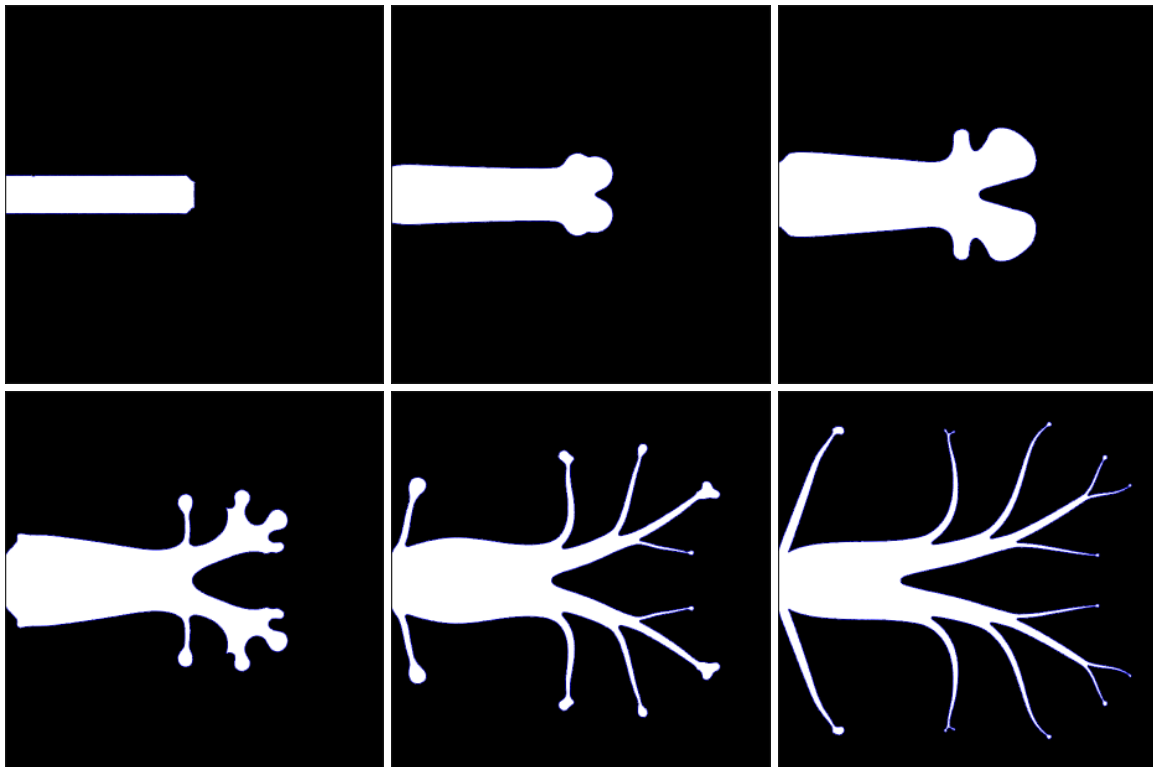


Figure 2.10: Intermediate iterations 0, 10, 25, 40, 90 and 200 for the pure heat conduction test case of [section 2.5.3](#) using the *surface* expression (2.4.14) of the shape derivative. Final objective function value:  $J = 2.8$ .

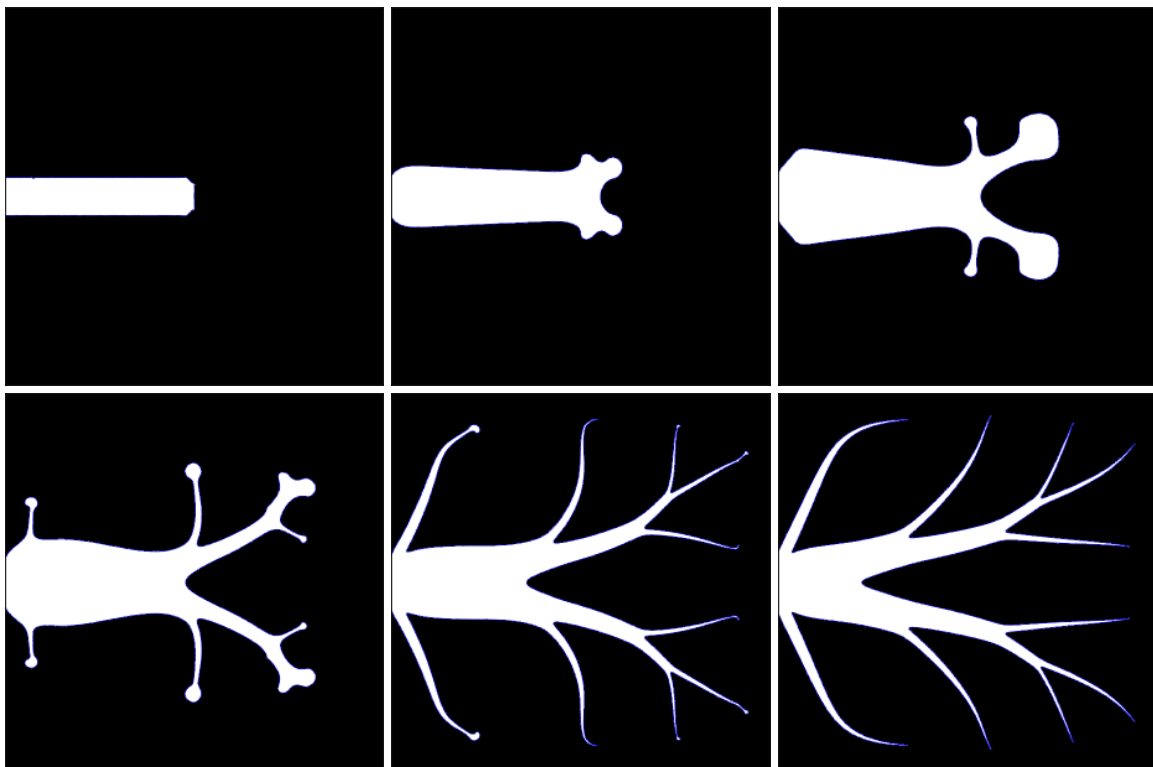


Figure 2.11: Intermediate iterations 0, 10, 25, 40, 90 and 190 for the pure heat conduction test case of [section 2.5.3](#) using the *volume* expression of the shape derivative. Final objective function value:  $J = 2.6$ .

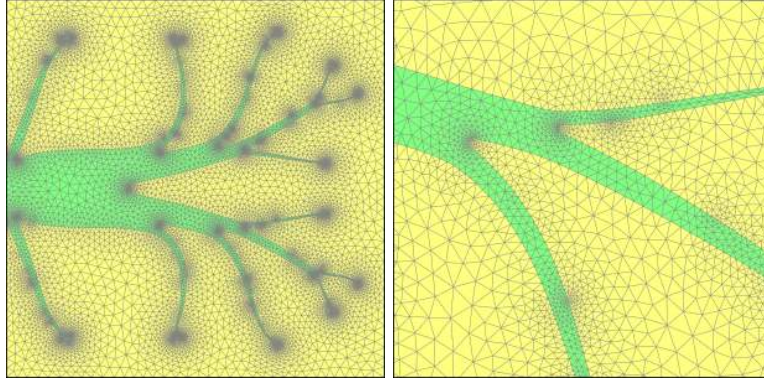


Figure 2.12: Final mesh and zoom on a part of it for the heat conduction test case of section 2.5.3.

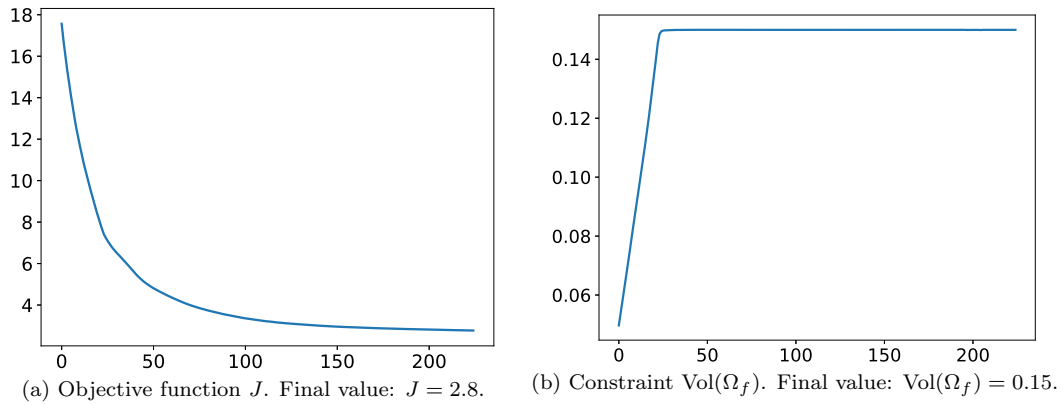


Figure 2.13: Convergence history for the heat conduction test case of section 2.5.3.

$\mathbf{v} \cdot \mathbf{n} = 0$  is assumed on these boundaries<sup>1</sup>. The flow leaves the domain on the right boundary with a zero normal stress boundary condition:  $\sigma_f(\mathbf{v}, p) \cdot \mathbf{n} = 0$ . The setting is illustrated on Figure 2.14).

The objective is to minimize the energy dissipation of the fluid generated by the solid structure, subject to a constraint on its volume, as well as on its center of mass fixed to the center  $\mathbf{x}_0$  of the domain:

$$\begin{aligned} \min_{\Gamma} \quad & J(\Gamma, \mathbf{v}(\Gamma)) = \int_{\Omega_f} 2\nu e(\mathbf{v}) : e(\mathbf{v}) dx \\ \text{s.t.} \quad & \begin{cases} \text{Vol}(\Omega_f) = V_{\text{target}} \\ \mathbf{X}(\Omega_s) := \frac{1}{|\Omega_s|} \int_{\Omega_s} \mathbf{x} dx = \mathbf{x}_0. \end{cases} \end{aligned} \quad (2.5.9)$$

We have set  $V_{\text{target}} = 0.03$  and  $\mathbf{x}_0 = (0.5, 0.5)$  in our implementation. The partial derivatives of the objective function allowing to obtain its shape derivative are given by

$$\frac{\partial \mathfrak{J}}{\partial \boldsymbol{\theta}}(\boldsymbol{\theta}) = \int_{\Omega_f} [-2\nu e(\mathbf{v}) : \nabla \mathbf{v} \nabla \boldsymbol{\theta} + 2\nu e(\mathbf{v}) : e(\mathbf{v}) \text{div}(\boldsymbol{\theta})] dx, \quad (2.5.10)$$

$$\frac{\partial \overline{\mathfrak{J}}}{\partial \boldsymbol{\theta}}(\boldsymbol{\theta}) = - \int_{\Gamma} 2\nu e(\mathbf{v}) : e(\mathbf{v}) (\boldsymbol{\theta} \cdot \mathbf{n}) ds, \quad (2.5.11)$$

$$\frac{\partial \mathfrak{J}}{\partial (\mathbf{v}, p)}(\mathbf{w}', q') = \int_{\Omega_f} 4\nu e(\mathbf{v}) : e(\mathbf{w}') dx. \quad (2.5.12)$$

The reader may verify that these equations plugged in the formulas (2.4.13) and (2.4.13) yield shape derivative formulas identical to those derived for the same problem e.g. in [112]. The derivative of the center of mass  $\mathbf{X}(\Omega_s)$  is given by

$$\frac{d}{d\boldsymbol{\theta}} \mathbf{X}(\Omega_{s,\boldsymbol{\theta}})(\boldsymbol{\theta}) = \frac{1}{|\Omega_s|} \int_{\Omega_s} (\mathbf{X}(\Omega_s) - \mathbf{x}) \text{div}(\boldsymbol{\theta}) dx = \frac{1}{|\Omega_s|} \int_{\Gamma} (\mathbf{X}(\Omega_s) - \mathbf{s}) \boldsymbol{\theta} \cdot \mathbf{n} ds.$$

<sup>1</sup>This boundary condition is not exactly the one considered in (2.2.1) but the computation of shape derivatives remains identical up to the account of the same boundary condition  $\mathbf{w} \cdot \mathbf{n} = 0$  for the adjoint variable on this boundary

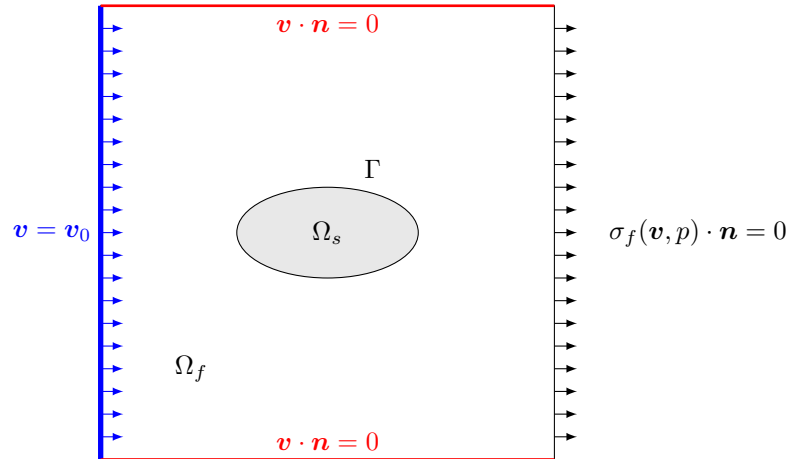


Figure 2.14: Setting of the minimum drag flow problem of section 2.5.4.

Final optimized shapes are shown on Figure 2.15 below in both situations where the fluid is characterized by the Stokes or Navier-Stokes equations. The celebrated rugby ball shape can be recognized for the Stokes problem [259]. Velocity fields are plotted on Figure 2.16. Convergence curves for objective and constraint functions are visible on Figure 2.19. We note that our constrained optimization algorithm was able to compute these shapes without the oscillations of the constraints obtained in [110] with the Augmented Lagrangian Method. A few intermediate iterations and the mesh of the final shape for the Navier Stokes case are respectively shown on Figs. 2.17 and 2.18.

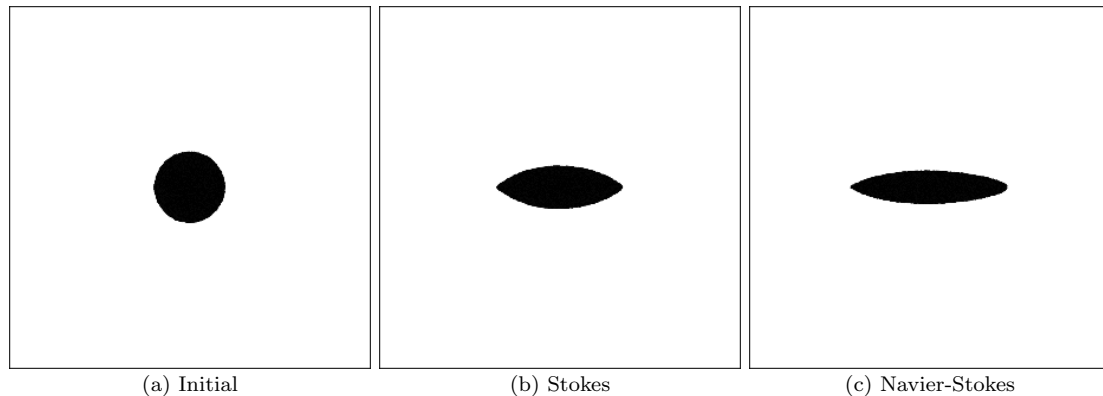


Figure 2.15: Initial and final configurations for the minimum drag test case of section 2.5.4.

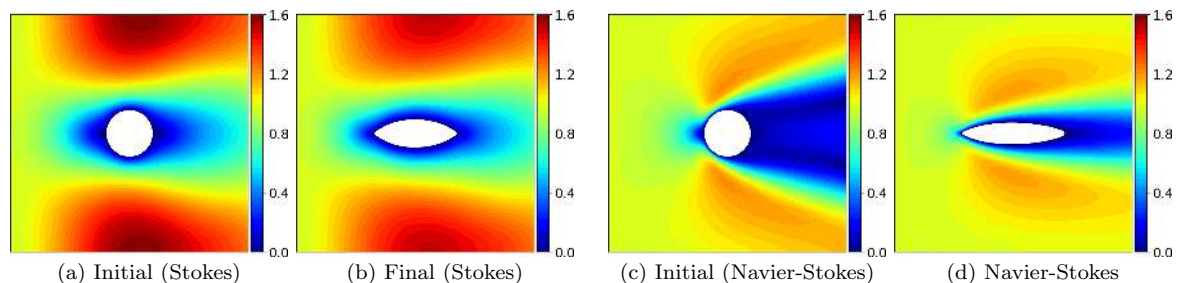


Figure 2.16: Initial and final norm fields for the fluid velocity.

### 2.5.5 Minimum compliance problem in thermoelasticity

In this paragraph, we reproduce the test case of Xia and Wang [314] for compliance minimization in thermoelasticity. A structure  $\Omega_s \subset D$  to be found in the rectangular domain  $D = [0, 2] \times [0, 1]$  is

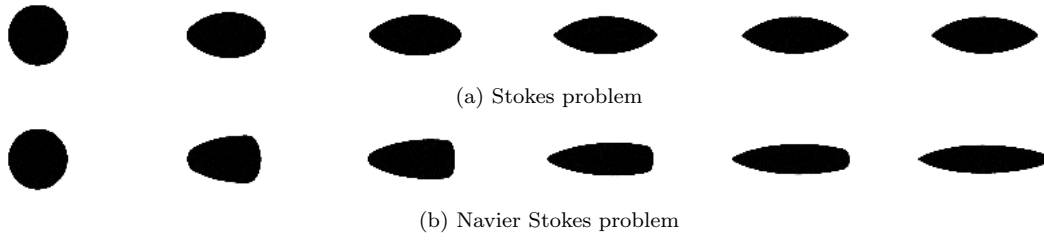


Figure 2.17: Intermediate iterations 0, 5, 10, 20, 30 and 70 for the minimum drag test case of section 2.5.4.

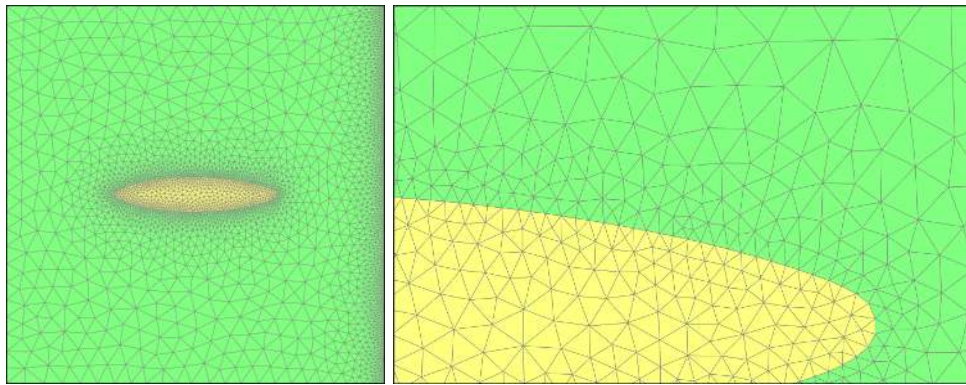


Figure 2.18: Mesh of the *Navier-Stokes* minimum drag profile and zoom on a part of the mesh. The mesh was refined near the boundary of the obstacle and the outlet boundary to capture the flow variability.

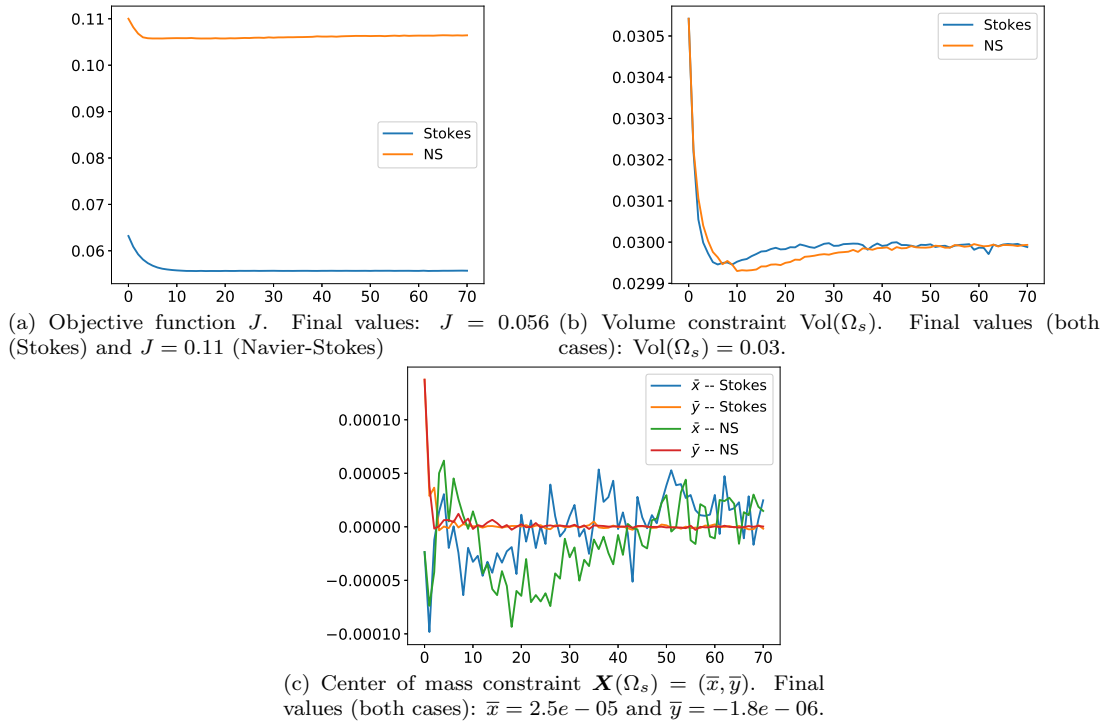


Figure 2.19: Convergence history for the test case of section 2.5.3.



clamped on the left and right sides of the domain, and subjected to a traction load on the middle of the bottom boundary. It is made of an elastic material characterized by Lamé parameters  $\lambda = 11510$ ,  $\mu = 7673$ , thermoelastic coefficient  $\alpha = 0.77$  and reference temperature  $T_{ref} = 0$ .

A constant temperature field  $T = T_{ref} + \Delta T$  is applied on the whole structure, which induces thermal expansion. The setting is reproduced on [Figure 2.20](#). We aim to solve the minimum compliance minimization problem subject to a volume inequality constraint:

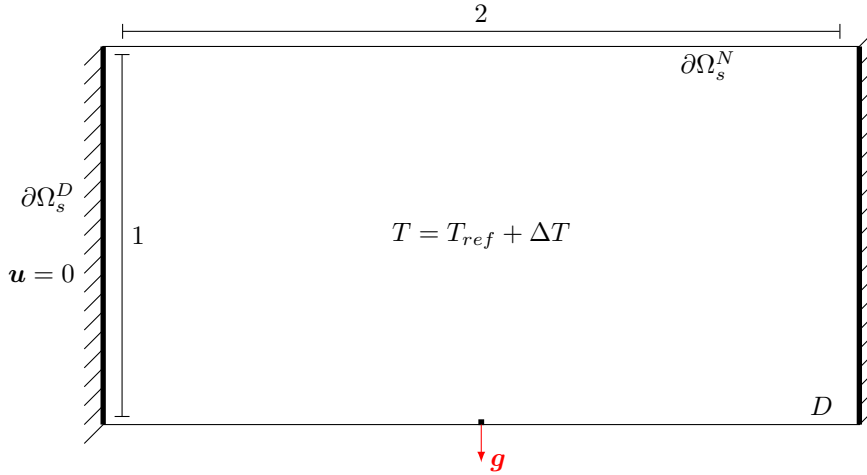


Figure 2.20: Setting for the thermoelastic compliance minimization problem of [section 2.5.5](#), issued from [\[314\]](#).

$$\begin{aligned} \min_{\Omega_s \subset D} \quad & J(\Omega_s, \mathbf{u}(\Omega_s)) := \int_{\Omega_s} A e(\mathbf{u}) : e(\mathbf{u}) dx \\ \text{s.t.} \quad & \text{Vol}(\Omega_s) := \int_{\Omega_s} dx \leq V_{target}. \end{aligned}$$

Strictly speaking, the problem still involves only one physics (the heat conduction problem [\(2.2.2\)](#) does not need to be solved since the temperature field is prescribed). We could have made it multiphysics by applying e.g. heat sources  $Q_f$  and  $Q_s$ , however we keep the setting of Xia and Wang [\[314\]](#) for the sake of comparison. The force applied in [\[314\]](#) has a value  $F = 1$ , which is set in our implementation by prescribing a traction force density  $\mathbf{g} = -1/\varepsilon$  on a small portion of size  $\varepsilon = 0.0125$  on the boundary. The upper bound volume is set to  $V_{target} = 0.4$ .

Following [\[314\]](#), we solve the optimization problem for four values of  $\Delta T$  ( $\Delta T = 0, 5, 10$  or  $20$ ). On [Figure 2.21](#), we plot the convergence history curves for the objective function and the volume constraint for each test case. Very interestingly, and as observed in [\[314\]](#), we retrieve the fact that the volume constraint  $\text{Vol}(\Omega_s) \leq V_{target}$  is saturated for the first two test cases  $\Delta T = 0$  and  $\Delta T = 5$ , and is not saturated otherwise. Notice as well in our case the smooth convergence of these curves obtained with our null space gradient flow, in contrast to those of the original paper which relied on a variant of the Augmented Lagrangian Method. In particular, our algorithm is able to quickly detect if the volume constraint needs to remain saturated or not.

Optimized shapes including intermediate iterations are shown on [Figure 2.22](#), and the corresponding final compliance and volume values are shown in [Table 2.3](#). Note that our numerical values do not coincide exactly with those of Wang because their original physical parameters were multiplied by nondimensionalization constants more compatible with our setting. However we clearly retrieve very similar optimized shapes.

## 2.5.6 A steady-state fluid-structure interaction problem

We now address a true multiphysics problem involving two physics. A fluid is flowing through a pipe, where it is pushing on a vertical beam of solid clamped at its bottom; see [Figure 2.23](#) for a schematic of the test case. Thermal effects are neglected (namely, [\(2.2.2\)](#) is ignored), so that [\(2.2.3\)](#) reduces to a standard linear elasticity system with the pressure load induced by the fluid. The objective is to minimize

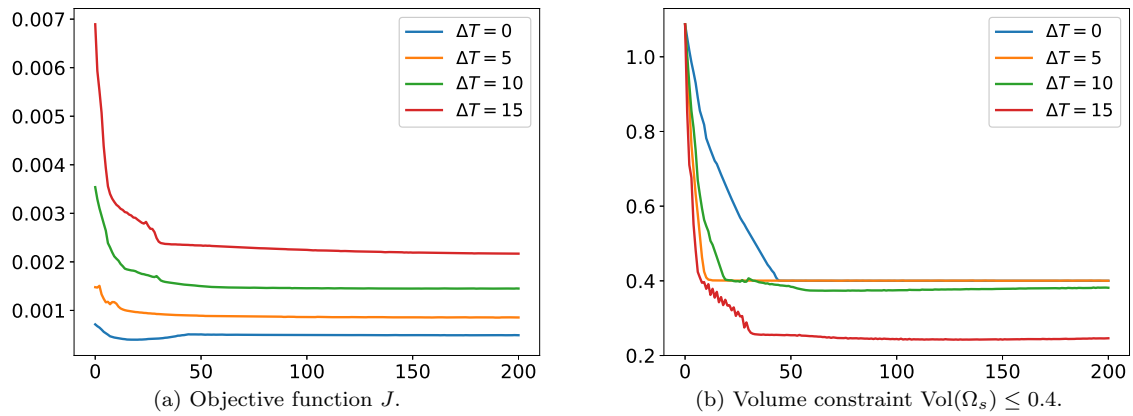


Figure 2.21: Convergence history for the thermoelasticity test case of section 2.5.5.

$\Delta T$	Final $J$	Final $\text{Vol}(\Omega_s)$
0	0.00049	0.4
5	0.00085	0.4
10	0.0015	0.38
15	0.0022	0.25

Table 2.3: Optimized compliance and volume values for the thermoelasticity test case of section 2.5.5. The results are analogous to those of [311].

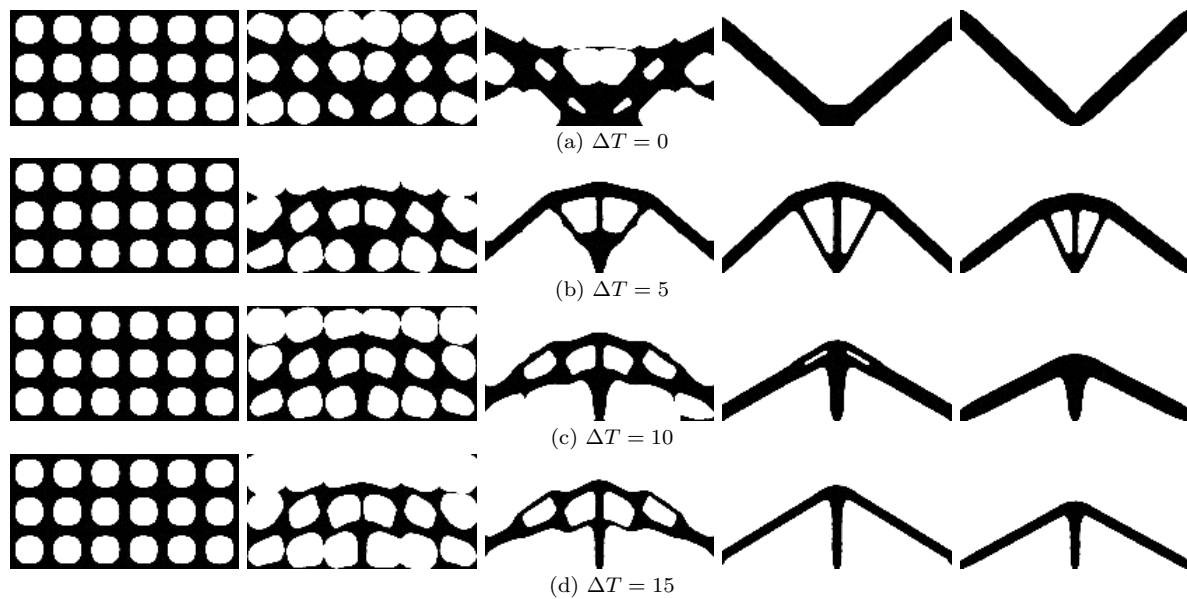


Figure 2.22: Intermediate iterations 0, 4, 13, 50 and 200 for each case of the thermoelasticity problem of section 2.5.5.

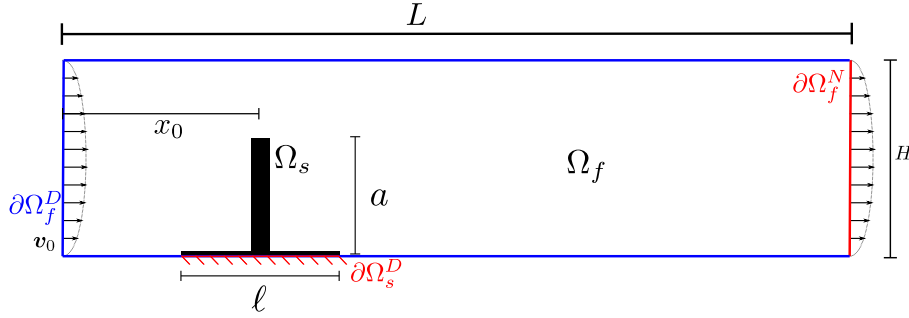


Figure 2.23: Physical setting of the fluid-structure optimization problem of [section 2.5.6](#). During optimization, the black domain cannot be reduced but only enlarged by adding reinforcements.

$L$	$H$	$\ell$	$a$	$\nu$	$\rho$	$\lambda$	$\mu$
2	0.5	0.4	0.3	0.005	1	0.00529	0.0476

Table 2.4: Numerical values of the physical parameters for the fluid-structure problem of [section 2.5.6](#)

the compliance of the solid phase  $\Omega_s$  subject to a volume constraint, that is:

$$\begin{aligned} \min_{\Gamma} \quad & J(\Gamma, \mathbf{u}(\Gamma)) = \int_{\Omega_s} A\epsilon(\mathbf{u}) : \epsilon(\mathbf{u}) dx \\ \text{s.t.} \quad & \text{Vol}(\Omega_s) = V_{target}. \end{aligned} \quad (2.5.13)$$

This example was previously considered by Yoon [319] with a different density-based (SIMP) method. In our case, we set  $V_{target} = 0.025|D|$  where  $|D| = 1$  is the volume of the total domain.

The numerical values of the considered physical parameters are given in [Table 2.4](#). The velocity profile  $\mathbf{v}_0$  imposed at the entrance of the pipe is parabolic, with maximum amplitude  $v_{\max} = 1$ , and value 0 at the upper and lower walls. The Reynolds number is equal to  $\text{Re} = \frac{\rho h v_{\max}}{\nu} = 60$ . The elastic displacement is set to 0 on a horizontal segment of length  $\ell$  supporting the beam.

In this part, we slightly elaborate from our study previously published in [153] in that we highlight some numerical differences between the use of the surface [\(2.4.14\)](#) and volume expression [\(2.4.13\)](#) when computing the shape derivative. The optimization problem [\(2.5.13\)](#) is solved for both cases.

The initial and optimized shapes featuring the norm of the fluid velocity fields  $\mathbf{v}$  are displayed on [Figure 2.24](#). The corresponding elastic deformations induced by the fluid on the solid structure (obtained by moving mesh nodes along  $\mathbf{u}$  for visualization purposes) are shown on [Figure 2.25](#). [Figure 2.26](#) shows parts of the initial and final meshes for the case using the surface expression of the shape derivative. Again, we point out that we used a feature of the library `mmg` to selectively enforce a fine mesh resolution near the interface  $\Gamma$ , while allowing larger triangles far from this boundary in order to save computational effort. Note that our final design are different from that in [319] because our Reynolds number is much larger and the location of the beam is different.

The evolution of the objective function and of the volume fraction are reported on [Figure 2.28](#). Note that in the first part of the optimization,  $J$  increases sometimes substantially (especially in the case relying on the volume expression of the shape derivative) due to the fact that the volume constraint is not yet satisfied, or due to sudden discontinuities at topological changes.

Several intermediate shapes are displayed on [Figure 2.29](#) when using the surface expression of the shape derivative, and on [Figure 2.30](#) when relying on the volume expression. The optimization path is surprisingly very different when using the volume expression; very thin parts of the structure exist for a longer time. Furthermore, although the final designs are somewhat similar, the one obtained with the volume expression has a slightly better performance for a very similar volume constraint. Note in particular the small solid bump near the top of the vertical beam which appears at the end of the optimization with the volume expression of the shape derivative. In our published work [153], we checked that this bump leads indeed to a better design; however it was obtained with the surface expression of the shape derivative at the cost of much more resolved meshes: the first design had a prescribed minimum edge length size  $h_{\min} = 0.001$  with 28,000 vertices in [153] versus  $h_{\min} = 0.003$  for approx. 11,000

vertices in the present work. The meshes of the optimized shape published in [153] are reproduced on Figure 2.27. Interestingly, in the present lower resolution case, only the volume expression of the shape derivative seems able to retrieve this bump.

The optimization ran in approximately half an hour on a laptop equipped with Intel(R) Core(TM) i7-4702MQ @ 2.20 GHz. Note that for this example, no optimization of the implementation was performed for reducing the total computational time, e.g. with the use of preconditioners or parallelism when solving finite element problems that represent most of the computational effort for these 2-d examples. Such ingredients become unavoidable for 3-d examples and are discussed in chapter 6.

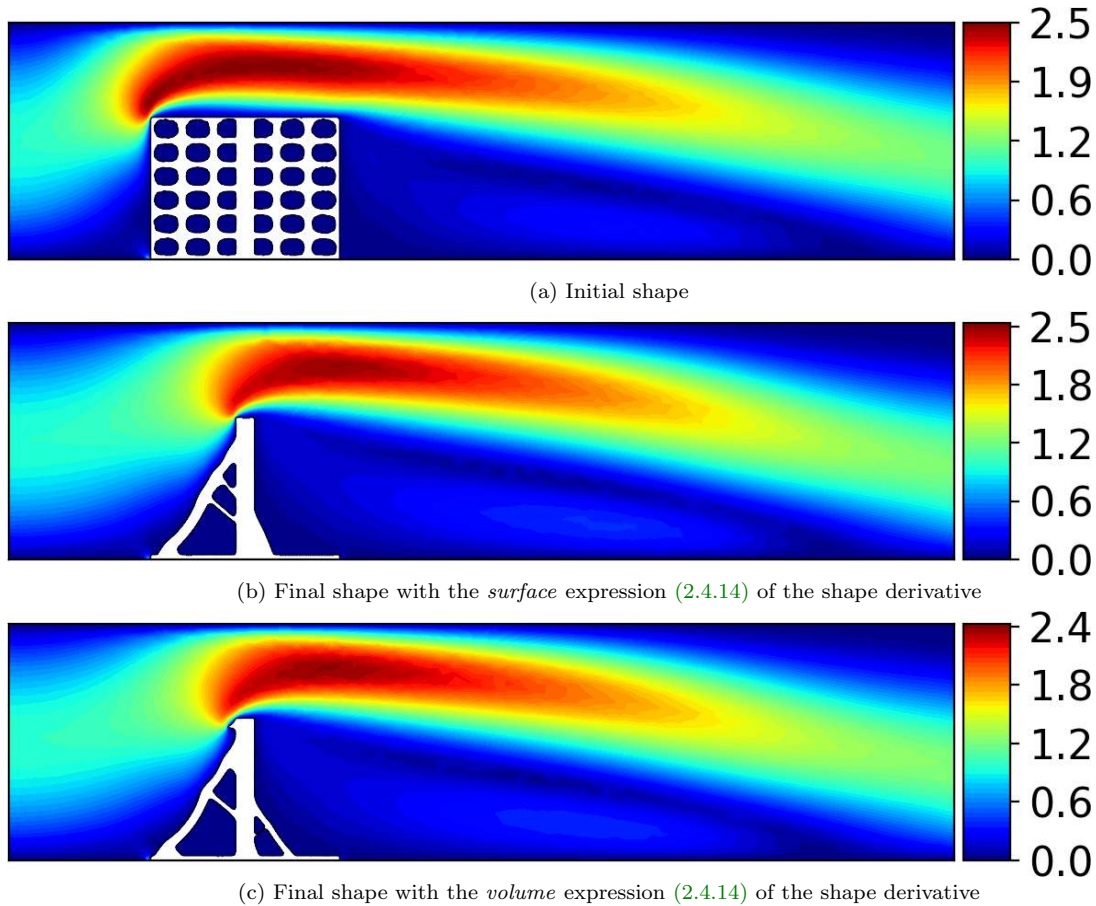


Figure 2.24: Initial and final shapes and velocity norm fields for the fluid-structure interaction test case of section 2.5.6. The optimized solid structure is depicted in white.

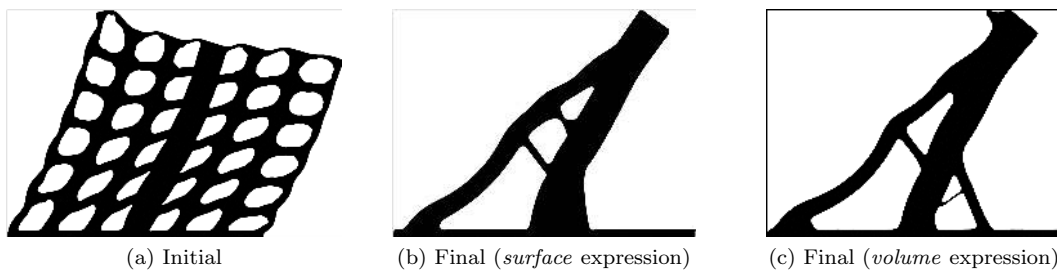


Figure 2.25: Computed elastic deformations of the initial and optimized solid structures subjected to the pressure induced by an inlet flow.

### 2.5.7 Convective heat transfer

Our second example involving two different physics is concerned with a coupling of the flow and heat equations (2.2.1) and (2.2.2), i.e. the elastic behavior of the complement  $\Omega_s$  of the optimized fluid phase

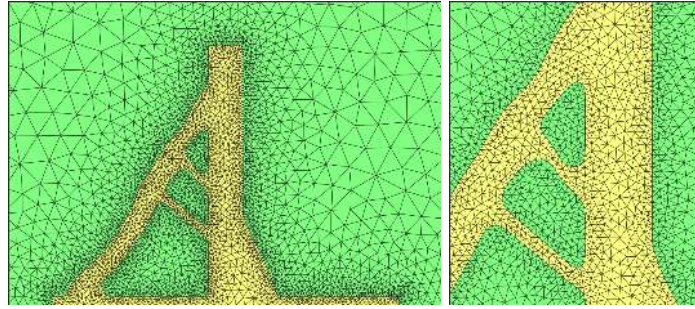


Figure 2.26: Zooms on the meshed solid structure of the final design obtained with the surface expression of the shape derivative. The minimum mesh size was set to  $h_{\min} = 0.003$ .

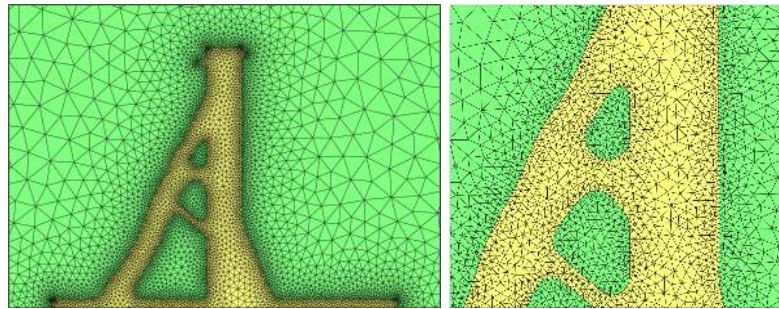


Figure 2.27: Zooms on the meshed solid structure of the final design from our published work [153] with the surface expression of the shape derivative. The minimum mesh size was set to  $h_{\min} = 0.001$ . The optimal shape is slightly different with the “bump” at the top left of the structure obtained with the volume expression.

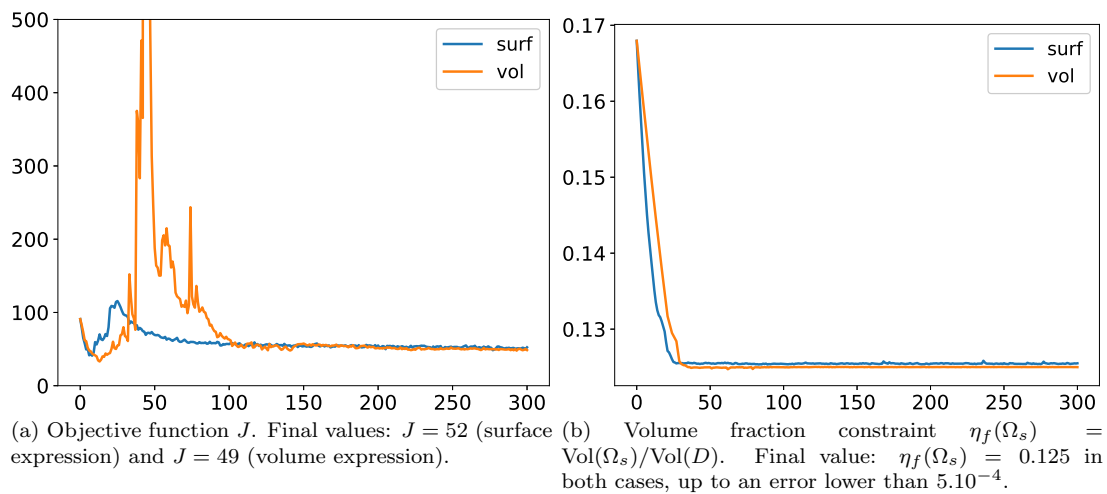


Figure 2.28: Convergence history for the fluid-structure test case of section 2.5.6. The mentions “surf” and “vol” refer respectively to the use of the surface and volume expressions of the shape derivatives.

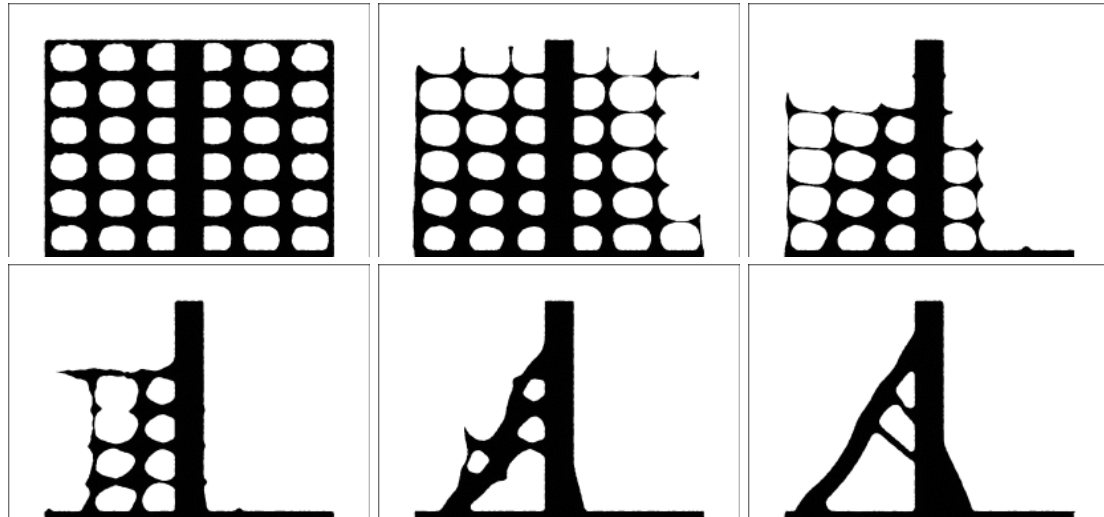


Figure 2.29: Intermediate iterations 0, 6, 15, 26, 56 and 300 for the fluid-structure test case of section 2.5.6 using the *surface* expression (2.4.14) of the shape derivative.

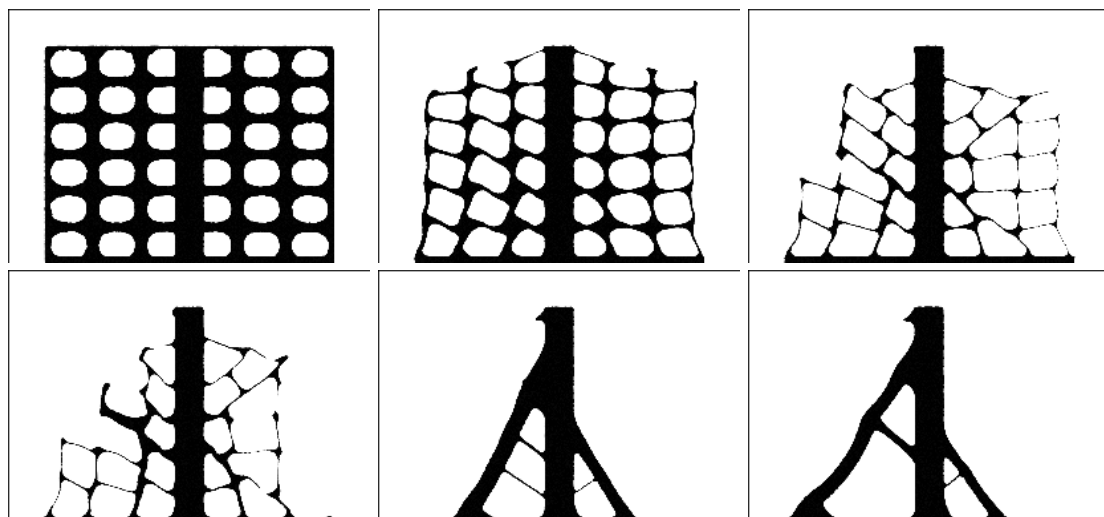


Figure 2.30: Intermediate iterations 0, 15, 40, 56, 150 and 300 for the fluid-structure test case of section 2.5.6 using the *volume* expression (2.4.14) of the shape derivative.

$L$	$\rho$	$c_p$	$\nu$	$v_{\max}$	$k_s$	$k_f$	$\omega$	$T_{in}$	$T_{low}$	$T_{up}$ (Ex. 1)	$T_{up}$ (Ex. 2)
0.1	10	100	0.005	1	10	1	0.4	0	10	-5	10

Table 2.5: Numerical values of the physical parameters in the convective heat transfer problem of section 2.5.7.

$\Omega_f$  is not taken into account (equation (2.2.3) is ignored). This test case features a cavity where a fluid is entering with a parabolic profile (with maximal velocity  $v_{\max} = 1$ ) and an inlet temperature  $T_{in}$ . The setting is similar to that in [225] (where a density based method is used), although we rely on different parameters values; namely, a higher Reynold number ( $Re = \frac{\rho L v_{\max}}{5\nu} = 40$  while  $Re=3$  in [225]) and different prescribed temperatures  $T_{up}$  and  $T_{low}$  for the lower and upper walls respectively. The other regions of the boundary of the cavity are insulated from the outside, i.e. zero normal fluxes boundary conditions hold for the temperature. The setting is represented on Figure 2.31 and numerical values of the parameters involved are reported in Table 2.5.

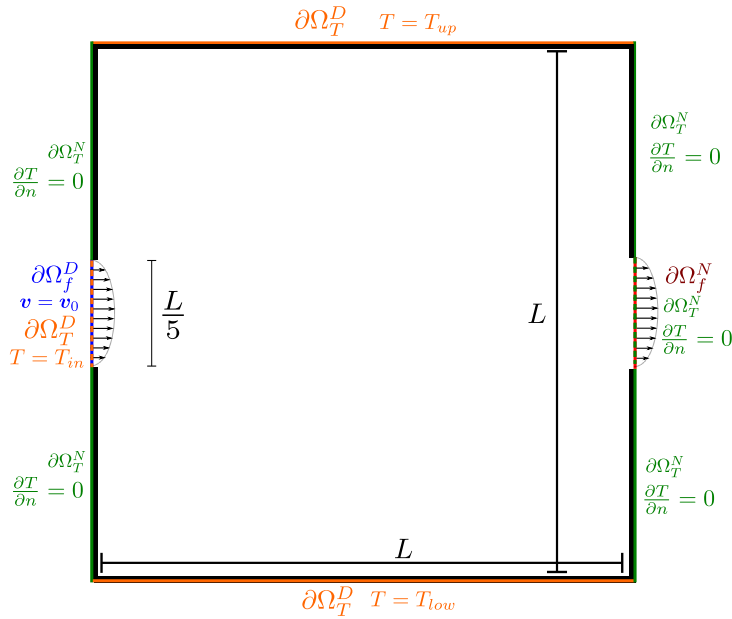


Figure 2.31: Setting of the convective heat transfer test case of section 2.5.7. The black layers at the walls stand for solid, non optimizable boundaries.

Our aim is maximize the heat transferred by the fluid subject to an upper bound on the output pressure drop and a volume constraint:

$$\begin{aligned}
 \min_{\Gamma} \quad & J(\Gamma, \mathbf{v}(\Gamma), T(\Gamma)) := - \int_{\Omega_f} \rho c_p \mathbf{v} \cdot \nabla T dx \\
 \text{s.t.} \quad & \begin{cases} \text{DP}(p(\Gamma)) := \int_{\partial\Omega_f^D} p ds - \int_{\partial\Omega_f^N} p ds \leq \text{DP}_{static} \\ \text{Vol}(\Omega_f) = V_{target}. \end{cases} \quad (2.5.14)
 \end{aligned}$$

Note that the objective function  $J$  rewrites indeed as the opposite of heat transferred from the inlet  $\partial\Omega_f^D$  to the outlet  $\partial\Omega_f^N$  upon integration by parts:

$$\int_{\Omega_f} \rho c_p \mathbf{v} \cdot \nabla T dx = \int_{\partial\Omega_f} \rho c_p T (\mathbf{v} \cdot \mathbf{n}) ds = \int_{\partial\Omega_f^N} \rho c_p T (\mathbf{v} \cdot \mathbf{n}) ds + \int_{\partial\Omega_f^D} \rho c_p T_0 (\mathbf{v}_0 \cdot \mathbf{n}) ds$$

where the second term is a constant depending on the inlet data. The upper bound constraint on the static pressure drop is set to  $\text{DP}_{static} = 11.4$  and the volume target to  $V_{target} = 0.2|D|$ .

**Remark 2.16.** This test case improves the one treated in our published work [153], sec. 6.3., in that we now handle a fully constrained problem (2.5.14) instead of a penalized version of it. The viscous energy

dissipation (as considered in the drag minimization test case of [section 2.5.4](#) and in the test case of [\[153\]](#)) has also been replaced by the static pressure drop DP, which has a more straightforward interpretation in industrial applications and also considered in topology optimization works, e.g. [\[255\]](#).

**Remark 2.17.** The constraint function DP does not *a priori* make sense because the pressure  $p$  belongs a priori to  $L^2(\Omega_f)/\mathbb{R}$  and has no well defined trace on the boundary, unless possibly under higher regularity assumptions. However, it is believed interesting to observe that our shape derivative formulas seem to work despite the assumptions violation.

For this example, the formulas needed for the calculation of the shape sensitivity of the objective function are:

$$\frac{\partial \mathfrak{J}}{\partial \boldsymbol{\theta}}(\boldsymbol{\theta}) = \int_{\Omega_f} [-\rho c_p \mathbf{v} \cdot \nabla T] \operatorname{div}(\boldsymbol{\theta}) + \rho c_p \mathbf{v} \cdot \nabla \boldsymbol{\theta}^T \nabla T dx, \quad \overline{\frac{\partial \mathfrak{J}}{\partial \boldsymbol{\theta}}}(\boldsymbol{\theta}) = 0 \quad (2.5.15)$$

$$\frac{\partial \mathfrak{J}}{\partial(\widehat{\mathbf{v}}, \widehat{p})}(\mathbf{w}', q') = - \int_{\Omega_f} \rho c_p \mathbf{w}' \cdot \nabla T dx, \quad \frac{\partial \mathfrak{J}}{\partial \widehat{T}}(S') = - \int_{\Omega_f} \rho c_p \mathbf{v} \cdot \nabla S' dx. \quad (2.5.16)$$

Furthermore, the shape derivative of the pressure drop constraint DP ([remark 2.17](#)) can be calculated (at least formally) with the following partial derivatives:

$$\frac{\partial \text{DP}}{\partial \boldsymbol{\theta}}(\boldsymbol{\theta}) = 0, \quad \overline{\frac{\partial \text{DP}}{\partial \boldsymbol{\theta}}}(\boldsymbol{\theta}) = 0, \quad (2.5.17)$$

$$\frac{\partial \text{DP}}{\partial(\widehat{\mathbf{v}}, \widehat{p})}(\mathbf{w}', q') = \int_{\partial \Omega_f^D} q' ds - \int_{\partial \Omega_f^N} q' ds. \quad (2.5.18)$$

We consider two possible configurations for the applied temperature  $T_{\text{up}}$  at the upper wall of the system. The first example corresponds to the configuration considered in [\[224\]](#); the upper and lower wall temperatures are equal and higher than the inlet temperature:  $T_{\text{in}} < T_{\text{up}} = T_{\text{low}}$ . In a second example (Ex.2), we consider the case where the upper wall temperature is lower than the inlet temperature:  $T_{\text{up}} < T_{\text{in}} < T_{\text{low}}$ .

The optimized shapes, associated with temperature fields and fluid velocity fields are represented respectively on [Figs. 2.32](#) to [2.33](#). Several intermediate shapes arising in the course of the optimization process are depicted on [Figure 2.36](#). In the first case featuring  $T_{\text{up}} = T_{\text{low}}$ , we retrieve an optimized shape analogous to those presented in [\[224\]](#). It consists of two main pipes connecting the inlet to the outlet with an additional vertical tubular bar of fluid domain where the velocity is almost zero. As already noticed in [\[224\]](#), this vertical inclusion takes advantage of the low diffusivity of the fluid material in order to insulate thermally the main pipes from the cold input left boundary  $T_{\text{in}}$ . Note that in contrast with the results presented in [\[153\]](#), the final shape is symmetric despite the use of non-symmetric meshes because the symmetry of the level set function was enforced at every iteration (see further details in [chapter 6](#) about this point).

In the second case featuring  $T_{\text{up}} < T_{\text{in}}$ , an analogous behavior is observed: the optimized shape consists of a main pipe accumulating the hot temperature from the bottom wall with an additional outgrowth of the fluid domain where the velocity is almost zero which insulates the main pipe from the cold upper wall.

The convergence curves for the objective and constraint functions are depicted on [Figure 2.37](#). Note that the slight increase of objective function at the end in case of the second configuration  $T_{\text{low}} < T_{\text{in}} < T_{\text{up}}$  is related to the disappearance of the very thin outgrowths of fluid at the top of the domain: these are beneficial for the performance of the shape but may tend to disappear because they are thinner than the prescribed mesh resolution.

## 2.5.8 Optimization of a compliant thermoelastic solid with fluid-structure interaction

We finally turn to a shape optimization example in the full three-physic setting presented in [section 2.2](#). A fluid is flowing from the left to the right of a two-dimensional pipe; at the center of this pipe, a solid body is attached to the boundary of a small non optimizable square  $\omega$  of side length  $c$ . The flow is entering the pipe at the inlet with a parabolic profile (with maximal velocity  $v_{\text{max}} = 1$ ), and a prescribed temperature  $T_{\text{in}}$  and the solid body receives a thermal flux  $h$  applied at the boundary  $\partial \omega$  of the square. The reference temperature of the solid material is equal to the fluid inlet temperature:  $T_{\text{ext}} = T_{\text{in}}$ . All the other boundaries in this device are insulated from the outside: zero Neumann boundary conditions  $\frac{\partial T}{\partial \mathbf{n}} = 0$  hold for the temperature; see [Figure 2.38](#) for a schematic of the problem.



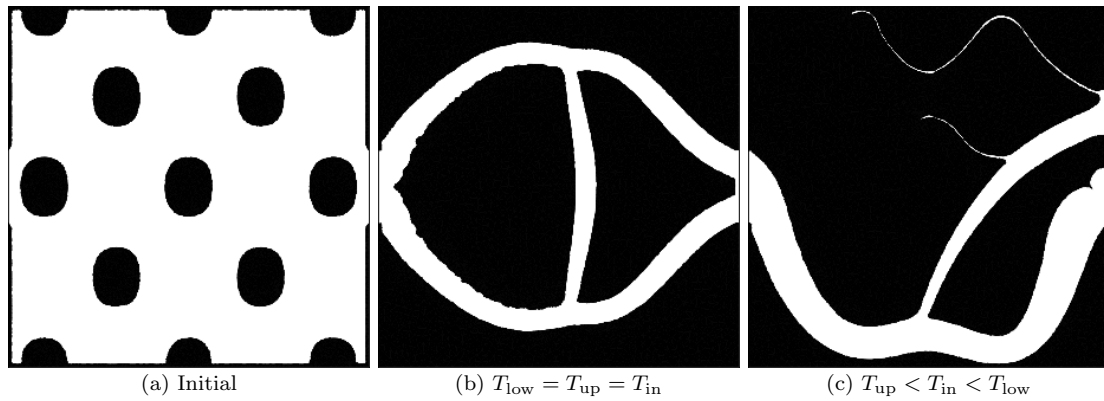


Figure 2.32: Initial and final configurations for the convective heat transfer test case of section 2.5.7.

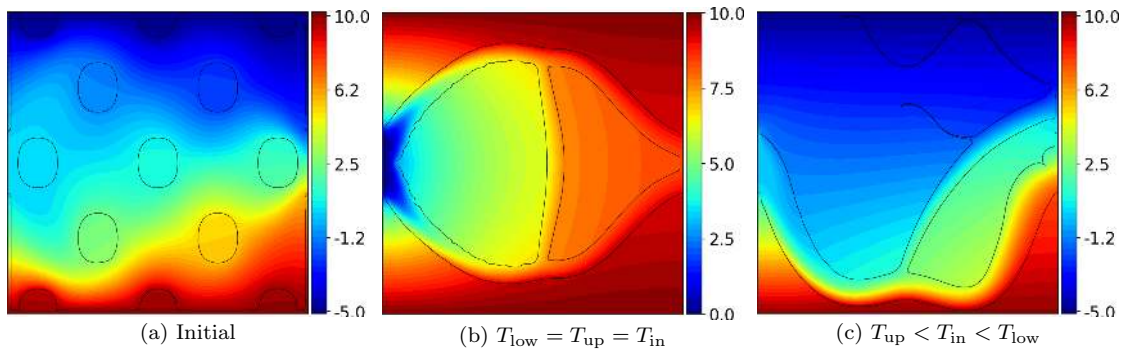


Figure 2.33: Initial and final temperature fields for the convective heat transfer test case of section 2.5.7.

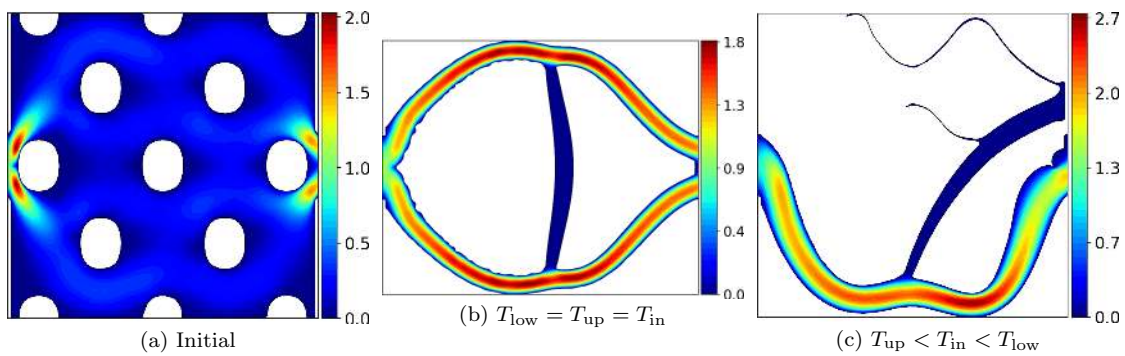


Figure 2.34: Initial and final norm fields of the velocity ( $v, p$ ) for the convective heat transfer test case of section 2.5.7.

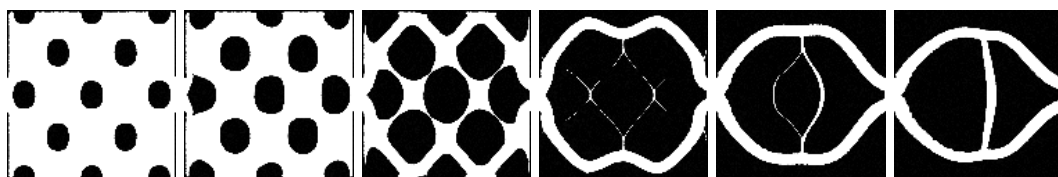


Figure 2.35: Intermediate iterations 0, 10, 30, 55, 188 and 300 for the convective heat transfer test case of section 2.5.7 with  $T_{low} = T_{up} = T_{in}$ .

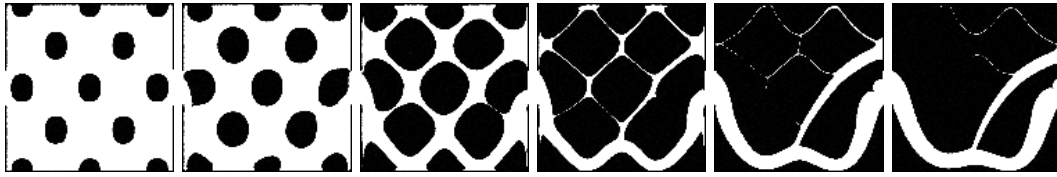


Figure 2.36: Intermediate iterations 0, 10, 30, 55, 188 and 300 for the convective heat transfer test case of section 2.5.7 with  $T_{up} < T_{in} < T_{low}$ .

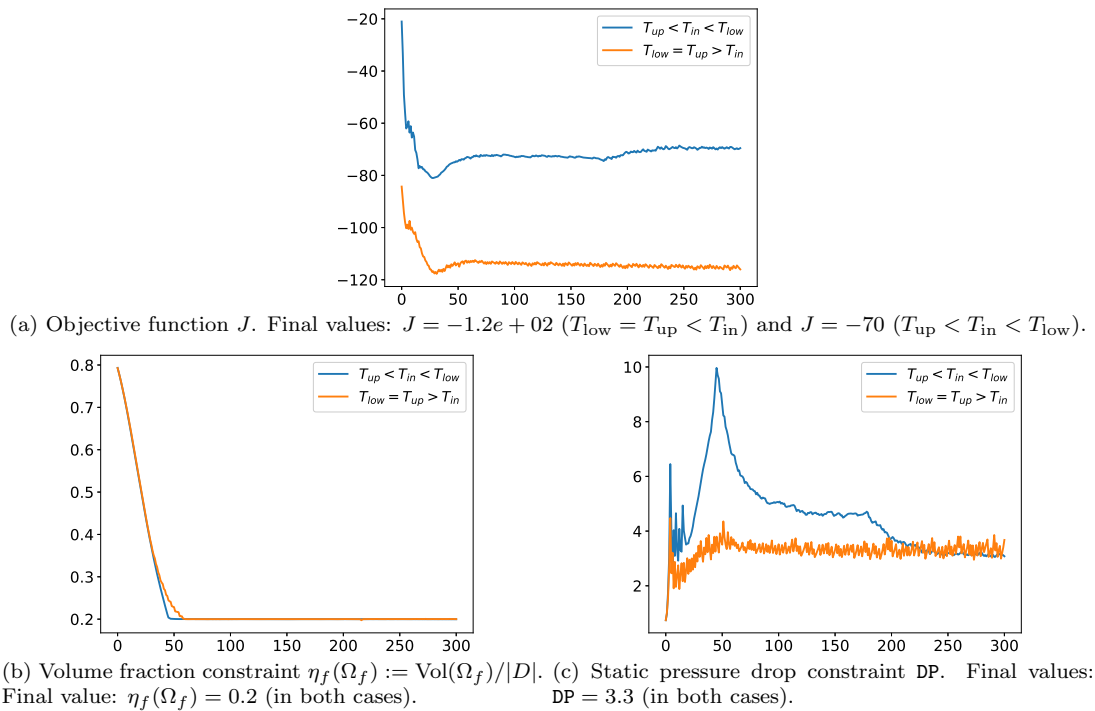


Figure 2.37: Convergence history for the convective heat transfer test case of section 2.5.7.

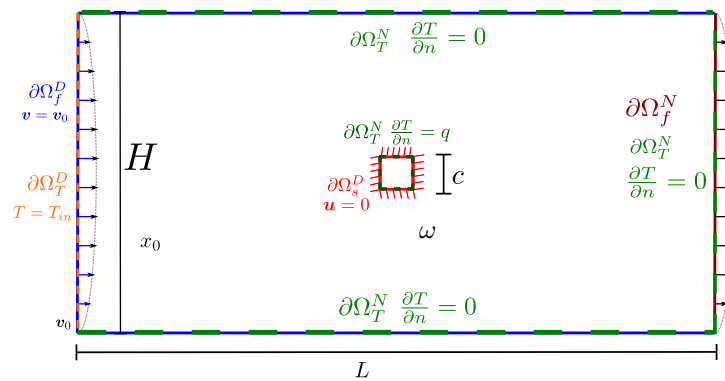


Figure 2.38: Setting of the thermoelastic fluid-structure problem of section 2.5.8.

$L$	$H$	$c$	$\rho$	$c_p$	$\nu$	$v_{\max}$	$k_s$	$k_f$	$T_{\text{in}}$	$T_{\text{ref}}$	$\lambda$	$\mu$	$\alpha$	$h$
2	1	0.1	1	0.5	0.01	1	10	1	0	0	12.96	5.55	3	$\pm 250$

 Table 2.6: Parameter values for the thermoelastic fluid-structure test case of [section 2.5.8](#).

We consider in total four different physical configurations corresponding to two possible signs of the thermal flux—either  $h > 0$  or  $h < 0$ —and to two different systems for the physical behavior of the flow: either the Navier-Stokes system (2.2.1), or its linear Stokes counterpart. In the case where  $h > 0$ , the square boundary  $\partial\omega$  plays the role of a thermal source: the high temperature in the solid body induces thermal expansion. In the latter case, where  $h < 0$ ,  $\partial\omega$  plays the role of a thermal sink: the lower temperature in the solid body induces thermal contraction. In both cases, the role of the fluid is to mitigate the temperature variations induced in  $\Omega_s$  by the thermal source term  $h$ .

The Reynolds number is set to  $\text{Re} = \frac{\rho L v_{\max}}{\nu} = 200$  and the Péclet number is  $\text{Pe} = \frac{L v_{\max}}{k_f} = 1000$ . The volume of the solid phase is imposed to be equal to  $\text{Vol}(\Omega_s) = 0.03\text{Vol}(D)$  where  $\text{Vol}(D) = 2$  is the volume of the total domain. In all the considered regimes, a sufficiently high value of the thermal dilation coefficient  $\alpha$  is used so as to make the thermoelastic effect dominant. The various numerical values for the physical parameters of the problem are summarized in [Table 2.6](#).

Our aim is to minimize the mechanical efforts induced in the solid structure  $\Omega_s$  by thermal dilation effects and the stress imposed by the fluid subject to the volume constraint:

$$\begin{aligned} \min_{\Gamma} \quad & J(\Gamma, \mathbf{u}(\Gamma)) := \int_{\Omega_s} \sigma_s(\mathbf{u}, T_s) : \nabla \mathbf{u} \, dx = \int_{\Omega_s} (Ae(\mathbf{u}) : e(\mathbf{u}) - \alpha(T_s - T_{\text{ref}})\text{div}(\mathbf{u})) \, dx, \\ \text{s.t.} \quad & \text{Vol}(\Omega_s) = V_{\text{target}}. \end{aligned} \quad (2.5.19)$$

The objective functional  $J(\Gamma, \mathbf{u}(\Gamma))$  corresponds to the internal energy stored inside the structure. Its sensitivities with respect to the shape  $\Gamma$  are calculated thanks to the following formulas:

$$\frac{\partial \mathfrak{J}}{\partial \boldsymbol{\theta}}(\boldsymbol{\theta}) = \int_{\Omega_s} (\sigma_s(\mathbf{u}, T_s) : \nabla \mathbf{u} \, \text{div}(\boldsymbol{\theta}) - 2Ae(\mathbf{u}) : \nabla \mathbf{u} \nabla \boldsymbol{\theta} + \alpha(T_s - T_{\text{ref}})\text{Tr}(\nabla \mathbf{u} \nabla \boldsymbol{\theta})) \, dx, \quad (2.5.20)$$

$$\overline{\frac{\partial \mathfrak{J}}{\partial \boldsymbol{\theta}}}(\boldsymbol{\theta}) = - \int_{\Gamma} [\sigma_s(\mathbf{u}, T_s) : \nabla \mathbf{u} - 2\mathbf{n} \cdot Ae(\mathbf{u}) \nabla \mathbf{u} \cdot \mathbf{n} + \alpha(T_s - T_{\text{ref}})\mathbf{n} \cdot \nabla \mathbf{u} \cdot \mathbf{n}] (\boldsymbol{\theta} \cdot \mathbf{n}) \, ds, \quad (2.5.21)$$

$$\frac{\partial \mathfrak{J}}{\partial \widehat{\mathbf{u}}}(\mathbf{r}') = \int_{\Omega_s} (2Ae(\mathbf{u}) : e(\mathbf{r}') - \alpha(T_s - T_{\text{ref}})\text{div}(\mathbf{r}')) \, dx, \quad (2.5.22)$$

$$\frac{\partial \mathfrak{J}}{\partial \widehat{T}}(S') = - \int_{\Omega_s} \alpha S' \text{div}(\mathbf{u}) \, dx. \quad (2.5.23)$$

The optimized shapes in the four situations are displayed in [Figure 2.39](#). Note that for the Stokes flow with  $h > 0$ , we used an initial shape perforated with holes in order to improve the convergence. For the other test cases, we used a disk shape initial domain. The convergence histories for the objective function and deviation to the volume constraint are shown on [Figure 2.40](#). Notably, our optimization algorithm is able to decrease the objective function while keeping constant the volume fraction in the solid phase. Several intermediate shapes are represented in [Figure 2.41](#) in the situation where  $h > 0$  and the fluid behavior is driven either by the Stokes or the full Navier-Stokes equations (2.2.1). For this latter case, the state variables  $\mathbf{v}, T$  and  $\mathbf{u}$  are additionally depicted on [Figure 2.42](#). In all four cases, the solid part  $\Omega_s$  tends to have a large contact surface with the fluid, so as to mitigate the effect of the thermal source (recall that  $T_{\text{in}} = T_{\text{ref}}$ ). The optimized shapes in the cases  $h > 0$  and  $h < 0$  are dramatically different, and are quite unintuitive from the mechanical viewpoint. Finally, the optimized shapes for a common value of  $h$  are noticeably different between the Stokes and Navier-Stokes cases, which could be expected due to the non negligible value of the Reynolds number. We have indeed checked that the optimized shape in the case of a Stokes flow has worse performance when evaluated in the context of a Navier-Stokes flow than the optimized shape in this setting (and vice-versa).

In this example, the objective function  $J$  turns out to be very sensitive with respect to very small variations of the shape. Recall that we do not resort to any upwinding scheme in our implementation. Therefore, we used a very fine mesh resolution (the minimum edge length is  $\text{hmin}=0.001$ ) as well as the

volume expression (2.3.27) of the shape derivative, which both enhanced the quality of the optimization process; see Figure 2.43.

On average, each intermediate mesh of  $D$  has approximately 20,000 vertices. A typical 300 iteration run of any of the aforementioned test cases (including Newton iterations for the numerical resolution of the Navier-Stokes equations) took approximately 6 hours on a 2.50 GHz Intel(R) Xeon(R) CPU.

As we have already mentioned, these are preliminary results. The ongoing work presented in chapter 6 focuses on more realistic 3-d test cases.

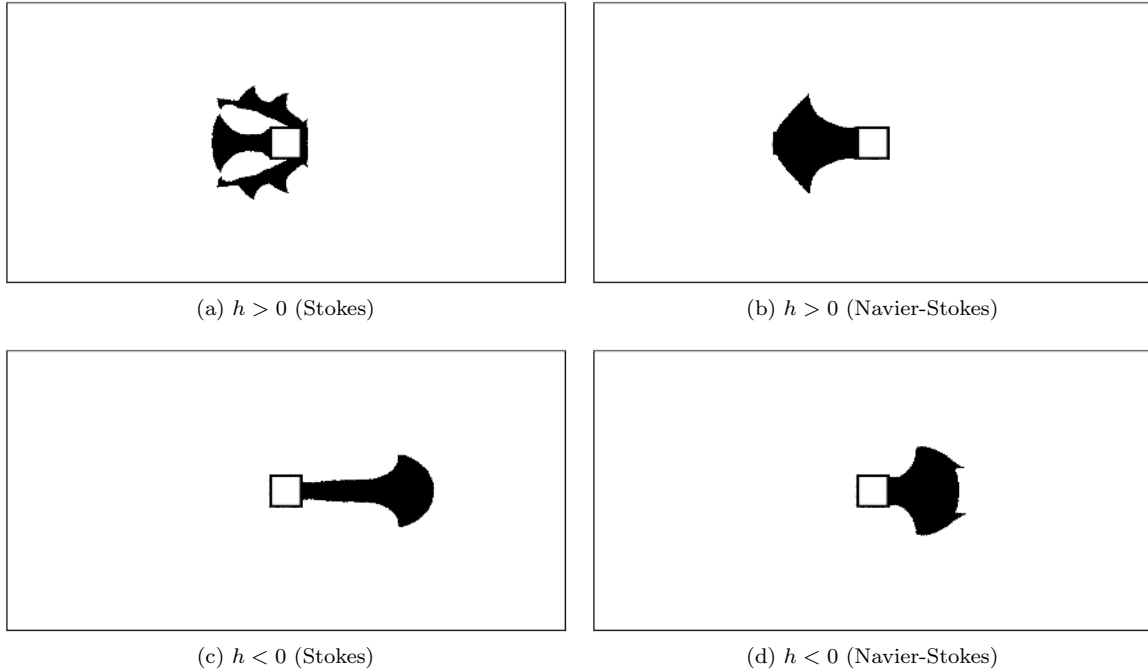


Figure 2.39: Optimized shapes for the three-physic test case of section 2.5.8 in the four considered physical situations.

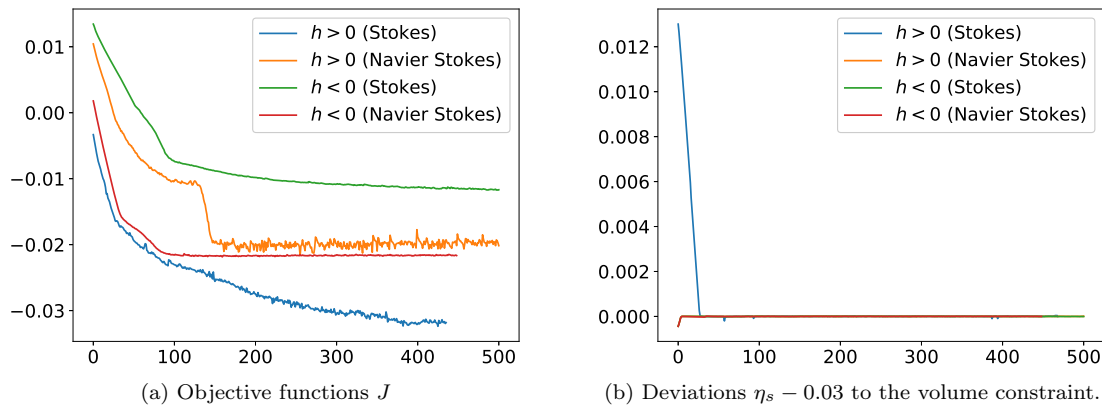


Figure 2.40: Convergence histories of the objective and constraint functions in the three-physic test case of section 2.5.8.

## 2.6 APPENDIX

### 2.6.1 Proof of propositions 2.3 and 2.4

We provide in this appendix a proof of propositions 2.3 and 2.4, or equivalently of (2.4.13) and (2.4.14), which is a mere adaptation of the arguments involved in section 2.3. Using classical arguments based on the implicit function theorem (see e.g. [184]), one proves that under the condition that the linearized

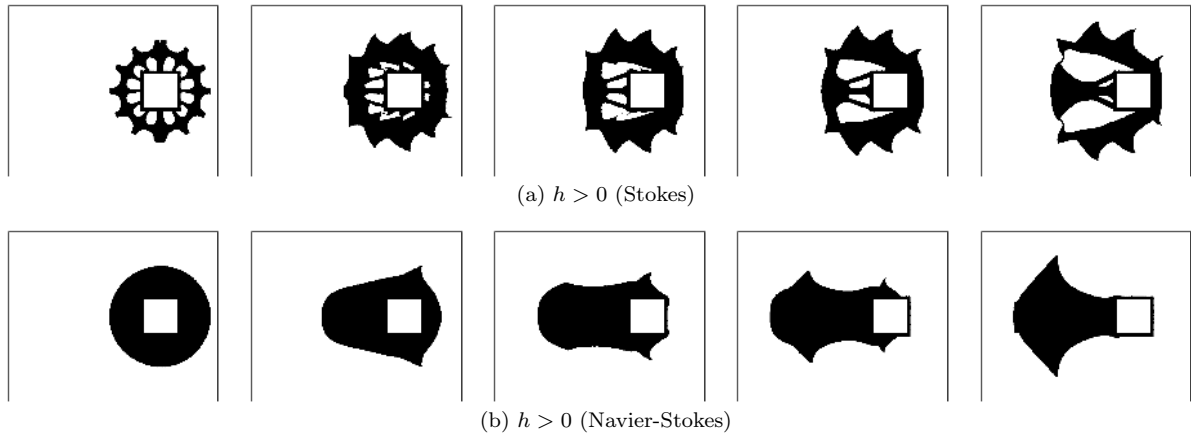


Figure 2.41: From left to right and top to bottom: iterations 1, 35, 80, 120 and 300 of the optimization process in the three-physic context of section 2.5.8 for Stokes and Navier Stokes flow where  $h > 0$ .

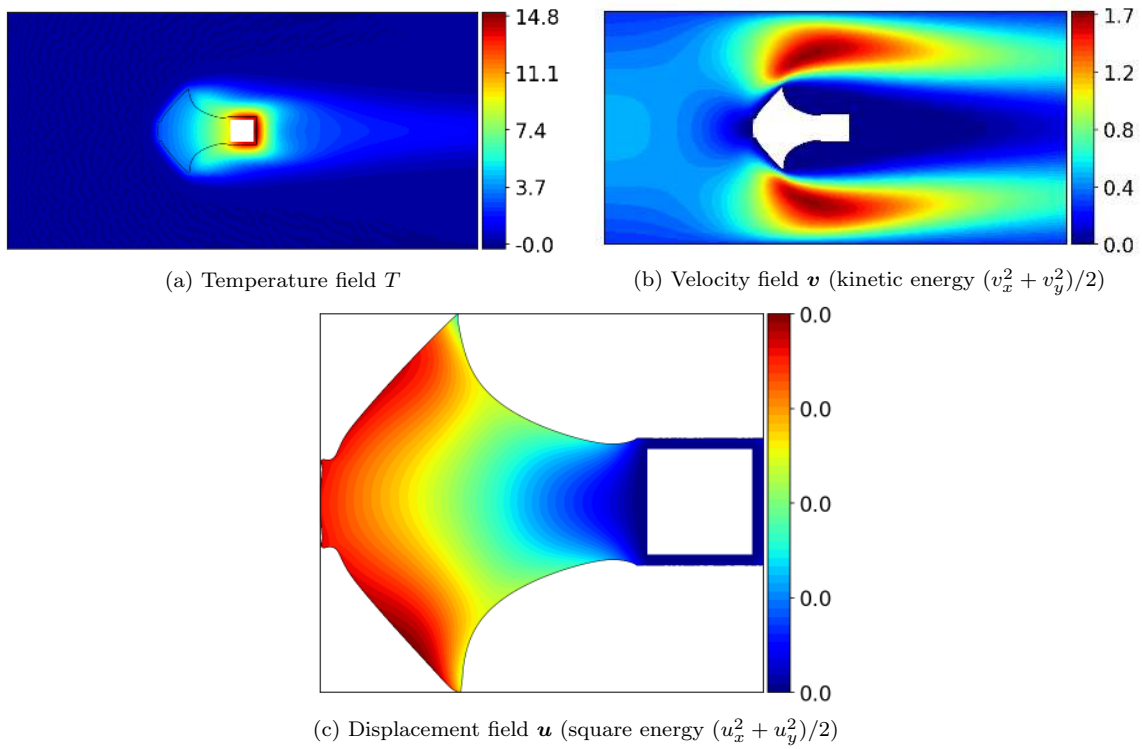


Figure 2.42: State variables  $\mathbf{v}$ ,  $T$  and  $\mathbf{u}$  for the optimized configuration of the three-physic shape optimization problem of section 2.5.8, in the situation  $h > 0$  and solved with the Navier-Stokes equations.

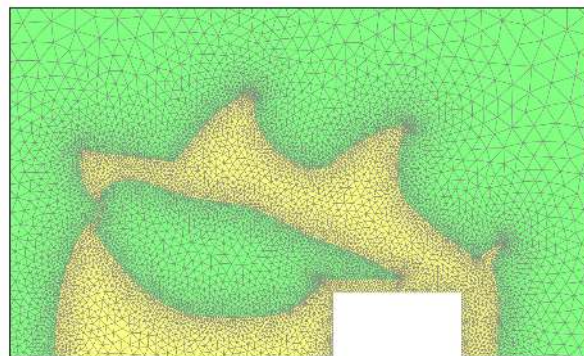


Figure 2.43: Zoom on the mesh for the final configuration of the Stokes case of Ex. 1.

version of the state equations (2.2.1) to (2.2.3) are well posed (see remark 2.12), the mappings  $\mathbf{v}(\Gamma_\theta) \circ (I + \theta)$ ,  $p(\Gamma_\theta) \circ (I + \theta)$ ,  $T(\Gamma_\theta) \circ (I + \theta)$ , and  $\mathbf{u}(\Gamma_\theta) \circ (I + \theta)$  are differentiable with respect to  $\theta$ . Differentiating the variational formulations (2.4.2) to (2.4.4), one finds that the Fréchet derivatives  $\dot{\mathbf{v}}(\theta)$ ,  $\dot{p}(\theta)$ ,  $\dot{T}(\theta)$  and  $\dot{\mathbf{u}}(\theta)$  at  $\theta = 0$  solve the following variational problems:

Find  $(\dot{\mathbf{v}}(\theta), \dot{p}(\theta)) \in V_{\mathbf{v},p}(\Gamma)$  such that  $\forall (\mathbf{w}', q') \in V_{\mathbf{v},p}(\Gamma)$ ,

$$\begin{aligned} & \int_{\Omega_f} [\sigma_f(\dot{\mathbf{v}}(\theta), \dot{p}(\theta)) : \nabla \mathbf{w}' + \rho \mathbf{w}' \cdot \nabla \mathbf{v} \cdot \dot{\mathbf{v}}(\theta) + \rho \mathbf{w}' \cdot \nabla \dot{\mathbf{v}}(\theta) \cdot \mathbf{v} - q' \operatorname{div}(\dot{\mathbf{v}}(\theta))] dx \\ &= \int_{\Omega_f} [\mathbf{w}' \cdot \operatorname{div}(\mathbf{f}_f \otimes \theta) - (\sigma_f(\mathbf{v}, p) : \nabla \mathbf{w}' + \rho \mathbf{w}' \cdot \nabla \mathbf{v} \cdot \mathbf{v}) \operatorname{div}(\theta)] dx \\ & \quad + \int_{\Omega_f} (\sigma_f(\mathbf{v}, p) : (\nabla \mathbf{w}' \nabla \theta) + \sigma_f(\mathbf{w}', q') : (\nabla \mathbf{v} \nabla \theta) + \rho \mathbf{w}' \cdot \nabla \mathbf{v} \nabla \theta \cdot \mathbf{v}) dx, \end{aligned} \quad (2.6.1)$$

Find  $\dot{T}(\theta) \in V_T(\Gamma)$  such that  $\forall S' \in V_T(\Gamma)$ ,

$$\begin{aligned} & \int_{\Omega_s} k_s \nabla \dot{T}(\theta) \cdot \nabla S' dx + \int_{\Omega_f} (k_f \nabla \dot{T}(\theta) \cdot \nabla S + \rho c_p S' \dot{\mathbf{v}}(\theta) \cdot \nabla T) dx = - \int_{\Omega_f} \rho c_p S' \mathbf{v} \cdot \nabla \dot{T}(\theta) dx \\ & \quad + \int_{\Omega_s} [\operatorname{div}(Q_s \theta) S' + k_s (\nabla \theta + \nabla \theta^T - \operatorname{div}(\theta) I) \nabla T \cdot \nabla S'] dx \\ & \quad + \int_{\Omega_f} [\operatorname{div}(Q_f \theta) S' + k_f (\nabla \theta + \nabla \theta^T - \operatorname{div}(\theta) I) \nabla T \cdot \nabla S'] dx \\ & \quad + \int_{\Omega_f} (-\rho c_p S' \mathbf{v} \cdot \nabla T \operatorname{div}(\theta) + \rho c_p S' \mathbf{v} \cdot \nabla \theta^T \nabla T) dx, \end{aligned} \quad (2.6.2)$$

Find  $\dot{\mathbf{u}}(\theta) \in V_{\mathbf{u}}(\Gamma)$  such that  $\forall \mathbf{r}' \in V_{\mathbf{u}}(\Gamma)$ ,

$$\begin{aligned} & \int_{\Omega_s} A e(\dot{\mathbf{u}}(\theta)) : \nabla \mathbf{r}' dx = \int_{\Omega_s} \alpha \dot{T}(\theta) \operatorname{div}(\mathbf{r}') dx + \int_{\Omega_s} [-\operatorname{div}(\theta) \sigma_s(\mathbf{u}, T_s) : \nabla \mathbf{r}' + \operatorname{div}(\mathbf{f}_s \otimes \theta) \cdot \mathbf{r}'] dx \\ & \quad + \int_{\Omega_s} (\sigma_s(\mathbf{u}, T_s) : (\nabla \mathbf{r}' \nabla \theta) + A e(\mathbf{r}') : (\nabla \mathbf{u} \nabla \theta)) dx - \int_{\Gamma} \mathbf{r}' \cdot \sigma_f(\dot{\mathbf{v}}(\theta), \dot{p}(\theta)) \cdot \mathbf{n} ds, \end{aligned} \quad (2.6.3)$$

where  $\sigma_f(\dot{\mathbf{v}}(\theta), \dot{p}(\theta)) \cdot \mathbf{n}$  is an element of the dual space  $H_{00}^{-1/2}(\Gamma, \mathbb{R}^d)$  of  $H_{00}^{1/2}(\Gamma, \mathbb{R}^d)$  whose action is given by (differentiating (2.4.5) with respect to  $\theta$ ):

$$\begin{aligned} & \forall \mathbf{r}' \in H_{00}^{1/2}(\Gamma, \mathbb{R}^d), - \int_{\Gamma} \mathbf{r}' \cdot \sigma_f(\dot{\mathbf{v}}(\theta), \dot{p}(\theta)) \cdot \mathbf{n} ds \\ &= \int_{\Omega_f} (\operatorname{div}(\mathbf{f}_f \otimes \theta) \cdot \tilde{\mathbf{r}} - (\rho \tilde{\mathbf{r}} \cdot \nabla \mathbf{v} \cdot \mathbf{v} + \sigma_f(\mathbf{v}, p) : \nabla \tilde{\mathbf{r}}) \operatorname{div}(\theta)) dx \\ & \quad + \int_{\Omega_f} (\rho \tilde{\mathbf{r}} \cdot \nabla \mathbf{v} \nabla \theta \cdot \mathbf{v} + \sigma_f(\mathbf{v}, p) : (\nabla \tilde{\mathbf{r}} \nabla \theta) + \sigma_f(\tilde{\mathbf{r}}, \tilde{q}) : (\nabla \mathbf{v} \nabla \theta)) dx \\ & \quad - \int_{\Omega_f} (\sigma_f(\dot{\mathbf{v}}(\theta), \dot{p}(\theta)) : \nabla \tilde{\mathbf{r}} + \rho \tilde{\mathbf{r}} \cdot \nabla \mathbf{v} \cdot \dot{\mathbf{v}}(\theta) + \rho \tilde{\mathbf{r}} \cdot \nabla \dot{\mathbf{v}}(\theta) \cdot \mathbf{v} - \tilde{q} \operatorname{div}(\dot{\mathbf{v}}(\theta))) dx, \end{aligned} \quad (2.6.4)$$

for any extension  $(\tilde{\mathbf{r}}, \tilde{q}) \in V_{\mathbf{v},p}(\Gamma)$  satisfying  $\tilde{\mathbf{r}} = \mathbf{r}'$  on  $\Gamma$ . Note that the above expression is independent of the chosen extension because of (2.6.1) with  $\mathbf{w}' = \tilde{\mathbf{r}}$  and  $q' = \tilde{q}$ . Then by the definition (2.4.7) of  $\mathfrak{J}$ :

$$\begin{aligned} & J(\Gamma_\theta, \mathbf{v}(\Gamma_\theta), p(\Gamma_\theta), T(\Gamma_\theta), \mathbf{u}(\Gamma_\theta)) \\ &= \mathfrak{J}(\theta, \mathbf{v}(\Gamma_\theta) \circ (I + \theta), p(\Gamma_\theta) \circ (I + \theta), T(\Gamma_\theta) \circ (I + \theta), \mathbf{u}(\Gamma_\theta) \circ (I + \theta)), \end{aligned} \quad (2.6.5)$$

whence the chain rule yields:

$$\frac{d}{d\theta} \left[ J(\Gamma_\theta, \mathbf{v}(\Gamma_\theta), p(\Gamma_\theta), T(\Gamma_\theta), \mathbf{u}(\Gamma_\theta)) \right] (\theta) = \frac{\partial \mathfrak{J}}{\partial \theta} (\theta) + \frac{\partial \mathfrak{J}}{\partial (\mathbf{v}, p)} (\dot{\mathbf{v}}(\theta), \dot{p}(\theta)) + \frac{\partial \mathfrak{J}}{\partial T} (\dot{T}(\theta)) + \frac{\partial \mathfrak{J}}{\partial \mathbf{u}} (\dot{\mathbf{u}}(\theta)). \quad (2.6.6)$$

One then uses the adjoint equations (2.4.8) to (2.4.10) with  $\mathbf{r}' = \dot{\mathbf{u}}(\boldsymbol{\theta})$ ,  $S' = \dot{T}(\boldsymbol{\theta})$ ,  $\mathbf{w}' = \dot{\mathbf{v}}(\boldsymbol{\theta})$ ,  $q' = \dot{p}(\boldsymbol{\theta})$  as test functions to obtain:

$$\frac{\partial \mathfrak{J}}{\partial \mathbf{u}}(\dot{\mathbf{u}}(\boldsymbol{\theta})) = \int_{\Omega_s} Ae(\mathbf{r}) : \nabla \dot{\mathbf{u}}(\boldsymbol{\theta}) dx = \int_{\Omega_s} Ae(\dot{\mathbf{u}}(\boldsymbol{\theta})) : \nabla \mathbf{r} dx, \quad (2.6.7)$$

$$\frac{\partial \mathfrak{J}}{\partial T}(\dot{T}(\boldsymbol{\theta})) = \int_{\Omega_s} k_s \nabla S \cdot \nabla \dot{T}(\boldsymbol{\theta}) dx + \int_{\Omega_f} (k_f \nabla S \cdot \nabla \dot{T}(\boldsymbol{\theta}) + \rho c_p S \mathbf{v} \cdot \nabla \dot{T}(\boldsymbol{\theta})) dx - \int_{\Omega_s} \alpha \dot{T}(\boldsymbol{\theta}) \operatorname{div}(\mathbf{r}) dx, \quad (2.6.8)$$

$$\begin{aligned} & \frac{\partial \mathfrak{J}}{\partial (\mathbf{v}, p)}(\dot{\mathbf{v}}(\boldsymbol{\theta}), \dot{p}(\boldsymbol{\theta})) \\ &= \int_{\Omega_f} (\sigma_f(\mathbf{w}, q) : \nabla \dot{\mathbf{v}}(\boldsymbol{\theta}) + \rho \mathbf{w} \cdot \nabla \mathbf{v} \cdot \dot{\mathbf{v}}(\boldsymbol{\theta}) + \rho \mathbf{w} \cdot \nabla \dot{\mathbf{v}}(\boldsymbol{\theta}) \cdot \mathbf{v} - \dot{p}(\boldsymbol{\theta}) \operatorname{div}(\mathbf{w})) dx + \int_{\Omega_f} \rho c_p S \nabla T \cdot \dot{\mathbf{v}}(\boldsymbol{\theta}) dx \\ &= \int_{\Omega_f} (\sigma_f(\dot{\mathbf{v}}(\boldsymbol{\theta}), \dot{p}(\boldsymbol{\theta})) : \nabla \mathbf{w} + \rho \mathbf{w} \cdot \nabla \dot{\mathbf{v}}(\boldsymbol{\theta}) \cdot \mathbf{v} + \rho \mathbf{w} \cdot \nabla \mathbf{v} \cdot \dot{\mathbf{v}}(\boldsymbol{\theta}) - q \operatorname{div}(\dot{\mathbf{v}}(\boldsymbol{\theta}))) dx + \int_{\Omega_f} \rho c_p S \nabla T \cdot \dot{\mathbf{v}}(\boldsymbol{\theta}) dx. \end{aligned} \quad (2.6.9)$$

Using now equations (2.6.2) and (2.6.3) with  $\mathbf{r}' = \mathbf{r}$ ,  $S' = S$  as test functions and (2.6.4) with  $(\tilde{\mathbf{r}}, \tilde{q}) = (\mathbf{w}, q)$  as an extension of  $\mathbf{r}' = \mathbf{r} \in H_{00}^{1/2}(\Gamma, \mathbb{R}^d)$  to eliminate the bilinear terms, the above three equations rewrite:

$$\begin{aligned} \frac{\partial \mathfrak{J}}{\partial \mathbf{u}}(\dot{\mathbf{u}}(\boldsymbol{\theta})) &= - \int_{\Gamma} \mathbf{r} \cdot \sigma_f(\dot{\mathbf{v}}(\boldsymbol{\theta}), \dot{p}(\boldsymbol{\theta})) \cdot \mathbf{n} ds + \int_{\Omega_s} \alpha \dot{T}(\boldsymbol{\theta}) \operatorname{div}(\mathbf{r}) dx \\ &\quad + \int_{\Omega_s} [-\operatorname{div}(\boldsymbol{\theta}) \sigma_s(\mathbf{u}, T_s) : \nabla \mathbf{r} + \operatorname{div}(\mathbf{f}_s \otimes \boldsymbol{\theta}) \cdot \mathbf{r}] dx \\ &\quad + \int_{\Omega_s} (\sigma_s(\mathbf{u}, T_s) : (\nabla \mathbf{r} \nabla \boldsymbol{\theta}) + Ae(\mathbf{r}) : (\nabla \mathbf{u} \nabla \boldsymbol{\theta})) dx, \end{aligned} \quad (2.6.10)$$

$$\begin{aligned} \frac{\partial \mathfrak{J}}{\partial T}(\dot{T}(\boldsymbol{\theta})) &= - \int_{\Omega_s} \alpha \dot{T}(\boldsymbol{\theta}) \operatorname{div}(\mathbf{r}) dx - \int_{\Omega_f} \rho c_p S \dot{\mathbf{v}}(\boldsymbol{\theta}) \cdot \nabla T dx \\ &\quad + \int_{\Omega_s} [\operatorname{div}(Q_s \boldsymbol{\theta}) S + k_s (\nabla \boldsymbol{\theta} + \nabla \boldsymbol{\theta}^T - \operatorname{div}(\boldsymbol{\theta}) I) \nabla T \cdot \nabla S] dx \\ &\quad + \int_{\Omega_f} [\operatorname{div}(Q_f \boldsymbol{\theta}) S + k_f (\nabla \boldsymbol{\theta} + \nabla \boldsymbol{\theta}^T - \operatorname{div}(\boldsymbol{\theta}) I) \nabla T \cdot \nabla S] dx \\ &\quad + \int_{\Omega_f} (-\rho c_p S \mathbf{v} \cdot \nabla T \operatorname{div}(\boldsymbol{\theta}) + \rho c_p S \mathbf{v} \cdot \nabla \boldsymbol{\theta}^T \nabla T) dx, \end{aligned} \quad (2.6.11)$$

$$\begin{aligned} \frac{\partial \mathfrak{J}}{\partial (\mathbf{v}, p)}(\dot{\mathbf{v}}(\boldsymbol{\theta}), \dot{p}(\boldsymbol{\theta})) &= \int_{\Gamma} \mathbf{r} \cdot \sigma_f(\dot{\mathbf{v}}(\boldsymbol{\theta}), \dot{p}(\boldsymbol{\theta})) \cdot \mathbf{n} ds + \int_{\Omega_f} \rho c_p S \nabla T \cdot \dot{\mathbf{v}}(\boldsymbol{\theta}) dx \\ &\quad + \int_{\Omega_f} (\operatorname{div}(\mathbf{f}_f \otimes \boldsymbol{\theta}) \cdot \mathbf{w} - (\rho \mathbf{w} \cdot \nabla \mathbf{v} \cdot \mathbf{v} + \sigma_f(\mathbf{v}, p) : \nabla \mathbf{w}) \operatorname{div}(\boldsymbol{\theta})) dx \\ &\quad + \int_{\Omega_f} (\rho \mathbf{w} \cdot \nabla \mathbf{v} \nabla \boldsymbol{\theta} \cdot \mathbf{v} + \sigma_f(\mathbf{v}, p) : (\nabla \mathbf{w} \nabla \boldsymbol{\theta}) + \sigma_f(\mathbf{w}, q) : (\nabla \mathbf{v} \nabla \boldsymbol{\theta})) dx. \end{aligned} \quad (2.6.12)$$

Formula (2.4.13) follows by summing up the above three equations. If  $H^2$  regularity holds for  $\mathbf{v}$ ,  $\mathbf{u}$ ,  $T$  and  $H^1$  regularity holds for  $p$  on their respective domains of definition, then an integration by parts allows

to rewrite (2.4.13) as;

$$\begin{aligned}
& \frac{d}{d\boldsymbol{\theta}} \left[ J(\Gamma_{\boldsymbol{\theta}}, \mathbf{v}(\Gamma_{\boldsymbol{\theta}}), p(\Gamma_{\boldsymbol{\theta}}), T(\Gamma_{\boldsymbol{\theta}}), \mathbf{u}(\Gamma_{\boldsymbol{\theta}})) \right] (\boldsymbol{\theta}) \\
&= \int_{\Gamma} \mathbf{g}_{\mathfrak{J}} \cdot \boldsymbol{\theta} ds + \int_{\Gamma} (\mathbf{f}_f \cdot \mathbf{w} - \sigma_f(\mathbf{v}, p) : \nabla \mathbf{w} - \rho \mathbf{w} \cdot \nabla \mathbf{v} \cdot \mathbf{v})(\boldsymbol{\theta} \cdot \mathbf{n}) ds \\
&+ \int_{\Gamma} [\mathbf{n} \cdot \sigma_f(\mathbf{v}, p) \nabla \mathbf{w} \cdot \boldsymbol{\theta} + \mathbf{n} \cdot \sigma_f(\mathbf{w}, q) \nabla \mathbf{v} \cdot \boldsymbol{\theta} + \rho(\mathbf{v} \cdot \mathbf{n}) \mathbf{w} \cdot \nabla \mathbf{v} \cdot \boldsymbol{\theta}] ds \\
&+ \int_{\Gamma} (k_s \nabla T_s \cdot \nabla S_s - k_f \nabla T_f \cdot \nabla S_f + Q_f S - Q_s S_s)(\boldsymbol{\theta} \cdot \mathbf{n}) ds \\
&+ \int_{\Gamma} (-k_s (\nabla T_s \cdot \boldsymbol{\theta})(\nabla S_s \cdot \mathbf{n}) - k_s (\nabla S_s \cdot \boldsymbol{\theta})(\nabla T_s \cdot \mathbf{n})) ds \\
&+ \int_{\Gamma} (k_f (\nabla T_f \cdot \boldsymbol{\theta})(\nabla S_f \cdot \mathbf{n}) + k_f (\nabla S_f \cdot \boldsymbol{\theta})(\nabla T_f \cdot \mathbf{n})) ds \\
&+ \int_{\Gamma} [(\sigma_s(\mathbf{u}, T_s) : \nabla \mathbf{r} - \mathbf{f}_s \cdot \mathbf{r})(\boldsymbol{\theta} \cdot \mathbf{n}) - \mathbf{n} \cdot \sigma_s(\mathbf{u}, T_s) \nabla \mathbf{r} \cdot \boldsymbol{\theta} - \mathbf{n} \cdot A\epsilon(\mathbf{r}) \nabla \mathbf{u} \cdot \boldsymbol{\theta}] ds + \int_{\Gamma} \boldsymbol{\Lambda} \cdot \boldsymbol{\theta} dx, \quad (2.6.13)
\end{aligned}$$

where  $\boldsymbol{\Lambda}$  is a  $L^1(D, \mathbb{R}^d)$  function obtained from Green's identity. The Hadamard structure theorem implies that (2.6.13) vanishes for compactly supported fields  $\boldsymbol{\theta}$  or for fields  $\boldsymbol{\theta}$  tangent to  $\Gamma$ . This implies that in fact,  $\boldsymbol{\Lambda} = 0$ , and (2.4.14) follows by removing the terms depending on the tangential component of  $\boldsymbol{\theta}$  on  $\Gamma$ .

## 2.6.2 Calculating the shape derivative of a particular objective functional and its adjoint system with C ea's method

In this appendix, we consider the simplified setting of section 2.3, and we show that the adjoint boundary condition  $p_s = p_f$  corresponding to the equality of normal derivatives for the primal variables in (2.3.2) can be retrieved formally with the classical C ea's Lagrangian method [84] also reviewed in chapter 1, section 1.2.3. We shall as well illustrate once again the calculation of the shape derivative with this method on a particular type of objective functional  $J$ .

We restrict ourselves in this part to objective functionals  $J$  of the form:

$$J(\Gamma, u_s(\Gamma), u_f(\Gamma)) = \int_{\Omega_f} j_f(u_f(\Gamma)) dx + \int_{\Omega_s} j_s(u_s(\Gamma)) dx, \quad (2.6.14)$$

for two  $C^2$  functions  $j_s, j_f : \mathbb{R} \rightarrow \mathbb{R}$  with bounded second-order derivatives:

$$\|j_s''\|_{L^\infty(\mathbb{R})} < \infty, \quad \|j_f''\|_{L^\infty(\mathbb{R})} < \infty.$$

**Proposition 2.5.** *The functional  $\boldsymbol{\theta} \mapsto J(\Gamma_{\boldsymbol{\theta}}, u_s(\Gamma_{\boldsymbol{\theta}}), u_f(\Gamma_{\boldsymbol{\theta}}))$ , from  $W_0^{1,\infty}(D, \mathbb{R}^d)$  into  $\mathbb{R}$ , as defined in (2.6.14), is differentiable at  $\boldsymbol{\theta} = 0$  and, under  $H^2$  regularity of the variables  $u_f, u_s, p_f, p_s$ , the shape derivative reads:*

$$\begin{aligned}
& \frac{d}{d\boldsymbol{\theta}} \left[ J(\Gamma_{\boldsymbol{\theta}}, u_f(\Gamma_{\boldsymbol{\theta}}), u_s(\Gamma_{\boldsymbol{\theta}})) \right] (\boldsymbol{\theta}) = \\
& \int_{\Gamma} \left[ j_f(u_f) - j_s(u_s) + f_f p_f - f_s p_s + \mu \nabla u_s \cdot \nabla p_s + \nu \frac{\partial u_f}{\partial \mathbf{n}} \frac{\partial p_f}{\partial \mathbf{n}} \right] (\boldsymbol{\theta} \cdot \mathbf{n}) ds. \quad (2.6.15)
\end{aligned}$$

**Remark 2.18.** Formula (2.6.15) is equivalent to the following one, which was obtained in (2.3.34) by the Lagrangian method in section 2.3.2,

$$\begin{aligned}
& \frac{d}{d\boldsymbol{\theta}} \left[ J(\Gamma_{\boldsymbol{\theta}}, u_f(\Gamma_{\boldsymbol{\theta}}), u_s(\Gamma_{\boldsymbol{\theta}})) \right] (\boldsymbol{\theta}) = \\
& \int_{\Gamma} \left[ j_f(u_f) - j_s(u_s) + f_f p_f - f_s p_s - \nu \nabla u_f \cdot \nabla p_f + \mu \nabla u_s \cdot \nabla p_s + 2\nu \frac{\partial u_f}{\partial \mathbf{n}} \frac{\partial p_f}{\partial \mathbf{n}} \right] (\boldsymbol{\theta} \cdot \mathbf{n}) ds, \quad (2.6.16)
\end{aligned}$$

because the boundary condition  $u_f = 0$  on  $\Gamma$  implies that  $\nu \frac{\partial u_f}{\partial \mathbf{n}} \frac{\partial p_f}{\partial \mathbf{n}} = \nu \nabla u_f \cdot \nabla p_f$ .



*Proof.* We introduce the following Lagrangian

$$\begin{aligned} \mathcal{L}(\Gamma, \widehat{u}_f, \widehat{u}_s, \widehat{p}_f, \widehat{p}_s, \widehat{\lambda}) &= \int_{\Omega_f} j_f(\widehat{u}_f) dx + \int_{\Omega_s} j_s(\widehat{u}_s) dx \\ &\quad - \int_{\Omega_f} (-\nu \Delta \widehat{u}_f - f_f) \widehat{p}_f dx - \int_{\Omega_s} (\mu \nabla \widehat{u}_s \cdot \nabla \widehat{p}_s - f_s \widehat{p}_s) dx - \int_{\Gamma} \nu \frac{\partial \widehat{u}_f}{\partial \mathbf{n}} \widehat{p}_s ds - \int_{\Gamma} \widehat{\lambda} \widehat{u}_f ds, \end{aligned} \quad (2.6.17)$$

where the above ‘hat’ functions all belong to the Sobolev space  $H_0^1(D)$  which is independent of the position of the interface  $\Gamma$ . Following the methodology described in [17, 84, 253], Lagrange multipliers  $\widehat{p}_f, \widehat{p}_s, \widehat{\lambda}$  are introduced in (2.6.17) to enforce the state equations (2.3.1) and (2.3.2) and the Dirichlet boundary condition  $u_f = 0$  on the moving interface  $\Gamma$ . The main idea of Céa’s method consists in finding the equations for the values  $p_f, p_s$  of the adjoint states  $\widehat{p}_f$  and  $\widehat{p}_s$  and the value  $\lambda$  of the Lagrange multiplier  $\widehat{\lambda}$  by requiring that the partial derivatives  $\frac{\partial \mathcal{L}}{\partial \widehat{u}_f}, \frac{\partial \mathcal{L}}{\partial \widehat{u}_s}$  at  $\widehat{u}_f = u_f, \widehat{u}_s = u_s$  vanish. The latter partial derivatives read, for arbitrary  $v_f, v_s \in H_0^1(D)$ :

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \widehat{u}_f}(v_f) &= \int_{\Omega_f} j'_f(u_f) v_f dx + \int_{\Omega_f} \nu \Delta v_f p_f dx - \int_{\Gamma} \nu \frac{\partial v_f}{\partial \mathbf{n}} p_s ds - \int_{\Gamma} \lambda v_f ds \\ &= \int_{\Omega_f} (j'_f(u_f) + \nu \Delta p_f) v_f dx + \int_{\Gamma} \nu \left( \frac{\partial v_f}{\partial \mathbf{n}} p_f - \frac{\partial p_f}{\partial \mathbf{n}} v_f \right) ds - \int_{\Gamma} \nu \frac{\partial v_f}{\partial \mathbf{n}} p_s ds - \int_{\Gamma} \lambda v_f ds \quad (2.6.18) \\ \frac{\partial \mathcal{L}}{\partial \widehat{u}_s}(v_s) &= \int_{\Omega_s} j'_s(u_s) v_s dx - \int_{\Omega_s} \mu \nabla v_s \cdot \nabla p_s dx \end{aligned}$$

Now requiring  $\frac{\partial \mathcal{L}}{\partial \widehat{u}_s}(v_s)$  vanish for any  $v_s \in H_0^1(D)$  yields the adjoint equation (and the attached boundary conditions) (2.3.23) for  $p_s$ .

Likewise, requiring  $\frac{\partial \mathcal{L}}{\partial \widehat{u}_f}(v_f)$  to vanish for any  $v_f \in H_0^1(D)$  should lead to the other adjoint system (2.3.24) but the derivation is a bit more subtle. First, choosing arbitrary  $v_f$  with compact support in  $\Omega_f$  in (2.6.18) yields:

$$-\nu \Delta p_f = j'_f(u_f).$$

Second, choosing smooth  $v_f$  such that  $v_f = 0$  on  $\Gamma$  and that  $v_f$  has an arbitrary trace  $\frac{\partial v_f}{\partial \mathbf{n}}$  yields the Dirichlet interface condition  $p_f = p_s$  on  $\Gamma$ . Thus, the adjoint system (2.3.24) for  $p_f$  is completely recovered. Finally, choosing  $v_f$  in (2.6.18) with arbitrary trace on  $\Gamma$  leads to the optimal value of the Lagrange multiplier  $\lambda = -\nu \frac{\partial p_f}{\partial \mathbf{n}}$  on  $\Gamma$ .

Assuming that the solutions  $u_s$  and  $u_f$  to (2.3.1) and (2.3.2) are differentiable with respect to  $\Gamma$  (which is where Céa’s method is only formal), and using that the partial derivatives of  $\mathcal{L}$  with respect to  $\widehat{u}_f$  and  $\widehat{u}_s$  vanish at  $(\boldsymbol{\theta}, \widehat{u}_f, \widehat{u}_s, \widehat{p}_f, \widehat{p}_s, \widehat{\lambda}) = (0, u_f, u_s, p_f, p_s, \lambda)$ , a simple use of the chain rule produces:

$$\begin{aligned} \frac{d}{d\boldsymbol{\theta}} \left[ J(\Gamma_{\boldsymbol{\theta}}, u_f(\Gamma_{\boldsymbol{\theta}}), u_s(\Gamma_{\boldsymbol{\theta}})) \right] (\boldsymbol{\theta}) &= \frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} (\boldsymbol{\theta}) \\ &= \int_{\Gamma} \left( j_f(u_f) - j_s(u_s) + \mu \nabla u_s \cdot \nabla p_s - f_s p_s - \nu \operatorname{div}(p_s \cdot \nabla u_f) - \lambda \frac{\partial u_f}{\partial \mathbf{n}} \right) (\boldsymbol{\theta} \cdot \mathbf{n}) ds \\ &= \int_{\Gamma} \left( j_f(u_f) - j_s(u_s) + \mu \nabla u_s \cdot \nabla p_s - f_s p_s - \nu \nabla p_s \cdot \nabla u_f + f_f p_s + \nu \frac{\partial u_f}{\partial \mathbf{n}} \frac{\partial p_f}{\partial \mathbf{n}} \right) (\boldsymbol{\theta} \cdot \mathbf{n}) ds, \end{aligned} \quad (2.6.19)$$

where we have used the technical proposition 1.7 of chapter 1 when differentiating  $\int_{\Gamma} \nu \frac{\partial u_f}{\partial \mathbf{n}} p_s ds$  with the normal vector depending on  $\Gamma$ . Using that  $u_f = 0$  and  $\frac{\partial p_s}{\partial \mathbf{n}} = 0$  on  $\Gamma$ , it follows that  $\nabla u_f \cdot \nabla p_s = 0$  on  $\Gamma$ , and since  $p_f = p_s$  on  $\Gamma$ , we retrieve expression (2.6.16).  $\square$



## CHAPTER 3

# NULL SPACE GRADIENT FLOWS FOR CONSTRAINED SHAPE OPTIMIZATION

### Contents

---

<b>3.1</b>	<b>Introduction</b> . . . . .	<b>107</b>
<b>3.2</b>	<b>Null space gradient flows for equality-constrained optimization in Hilbert spaces</b> . . . . .	<b>111</b>
3.2.1	Notation and first-order optimality conditions . . . . .	111
3.2.2	Definitions and properties of the null space and range space steps $\xi_J$ and $\xi_C$ . . . . .	113
3.2.3	Decrease properties of the equality constrained gradient flow . . . . .	115
<b>3.3</b>	<b>Extension to equality and inequality constraints</b> . . . . .	<b>117</b>
3.3.1	Notation and preliminaries . . . . .	117
3.3.2	Definition of the range space step . . . . .	118
3.3.3	Definition and properties of the null space step . . . . .	118
3.3.4	Decrease properties of the trajectories of the null space ODE . . . . .	123
<b>3.4</b>	<b>Numerical discretization and time-stepping schemes for the null space ODE</b> . . . . .	<b>125</b>
3.4.1	Accounting for discontinuities near the inequality constraint barriers . . . . .	125
3.4.2	Time step adaptation based on a merit function. . . . .	126
3.4.3	Overall algorithm pseudo code . . . . .	127
<b>3.5</b>	<b>Comparisons with other methods and illustrations on academic test cases</b> . . . . .	<b>127</b>
3.5.1	Comparison with the method of slack variables for inequality constraints . . . . .	127
3.5.2	Comparisons with ‘iterative’ optimization algorithms . . . . .	130
3.5.3	Comparative academic test cases in the euclidean space $\mathbb{R}^k$ . . . . .	132
<b>3.6</b>	<b>Optimization within the set of Lipschitz subdomains: applications to shape optimization</b> . . . . .	<b>138</b>
3.6.1	Manifold structures for shape optimization . . . . .	138
3.6.2	Adaptive normalizations for the null space and range space directions . . . . .	140
3.6.3	Illustrations on a multiple load structural shape optimization test case . . . . .	141

---

*Note* : most of the content of this chapter has been submitted in the preprint [155]:

F. FEPPON, G. ALLAIRE, AND C. DAPOGNY, *Null space gradient flows for constrained optimization with applications to shape optimization*, Submitted, (2019).

However, several changes have been made in the introduction and in [section 3.4](#) for a better understanding, as well as a new [section 3.5](#) providing additional discussions and comparisons of our method with other algorithms. Some redundant introductory parts with the previous chapters have also been removed.

### 3.1 INTRODUCTION

The purpose of this chapter is to introduce efficient and reliable algorithmic methodologies for the resolution of general shape optimization problems featuring an arbitrary number of equality and inequality constraints; for multiphysics applications in the context of [chapter 2](#) , these problems can be generically formulated as

$$\begin{aligned} \min_{\Gamma} \quad & J(\Gamma, \mathbf{v}(\Gamma), p(\Gamma), T(\Gamma), \mathbf{u}(\Gamma)) \\ \text{s.t.} \quad & \begin{cases} g_i(\Gamma, \mathbf{v}(\Gamma), p(\Gamma), T(\Gamma), \mathbf{u}(\Gamma)) = 0, & 1 \leq i \leq p, \\ h_j(\Gamma, \mathbf{v}(\Gamma), p(\Gamma), T(\Gamma), \mathbf{u}(\Gamma)) \leq 0, & 1 \leq j \leq q. \end{cases} \end{aligned} \tag{3.1.1}$$

The ability to account for equality or inequality constraints  $g_i$  or  $h_j$  in the resolution of such optimal design problems is a *sine qua non* condition towards the application of shape and topology optimization methods to realistic industrial applications. Indeed, industrial designs are in most cases devised with

respect to a variety of load specifications, such as minimum thickness, maximum curvature radius, upper bound limits on the stress in the case of a mechanical structure or on the pressure drop between the outlet and inlet in the case of a fluid device.

Over the past decades, many iterative algorithms have been proposed to solve generic constrained optimization problems of the form:

$$\begin{aligned} \min_{x \in V} \quad & J(x) \\ \text{s.t.} \quad & \begin{cases} \mathbf{g}(x) = 0 \\ \mathbf{h}(x) \leq 0, \end{cases} \end{aligned} \quad (3.1.2)$$

where  $V$  is the optimization set,  $J : V \rightarrow \mathbb{R}$  is a differentiable objective function,  $\mathbf{g} : V \rightarrow \mathbb{R}^p$  and  $\mathbf{h} : V \rightarrow \mathbb{R}^q$  are differentiable functions accounting for  $p$  equality and  $q$  inequality differentiable constraints, respectively. Classical gradient-based algorithms for the numerical resolution of (3.1.2) include, e.g., Penalty, Lagrangian, Interior Point and Trust Region Methods, Sequential Quadratic or Linear Programming (SQP or SLP) [70, 244, 321, 161, 310], the Method of Moving Asymptotes (MMA) [298], the Method of Feasible Directions [327, 308].

A major difficulty related to the practical use of these algorithms for topology optimization lies in that all the aforementioned techniques require fine tuning of the algorithm parameters in order to actually solve the minimization problem. These parameters are e.g. the penalty coefficients in the Augmented Lagrangian and Interior Point methods, the size of the trust region in SLP algorithms, the strategy for approximating the Hessian matrix in SQP, the bounds on the asymptotes in MMA and the Topkis parameters in MFD. The correct determination of these parameters is strongly case-dependent and often unintuitive: for instance, the penalty coefficients must be neither ‘too large’ nor ‘too small’ in Lagrangian methods, the SLP trust region size—which acts as a step length—cannot be chosen too small (otherwise the involved quadratic subproblems may not have a solution). For shape and topology optimization applications, a fair amount of trials and errors is often required in order to obtain satisfying minimizing sequences of shapes. Since every optimization step depends on the resolution of partial differential equations, such tunings are very tedious, time consuming for 2-d cases, and simply not affordable for realistic 3-d applications.

As a matter of fact, advanced mathematical programming methods are not frequently described for shape optimization based on Hadamard’s method. Rather, for simplicity of implementation, Penalty and Augmented Lagrangian Methods are often used, all the more when only one constraint is considered [32, 113]. Morin et. al. introduced a variant of SQP in [237] but they treated a volume constraint with a Lagrange Multiplier method. For more complex applications, some authors have proposed adapted variants of Sequential Linear Programming [135] or of the Method of Feasible Directions [150]. Let us remark that in shape optimization based on the method of Hadamard, an additional difficulty comes into play due to the classical identification and regularization step of the descent direction. So far, this matter has not been treated explicitly for constrained optimization problems: common approaches rather compute a descent direction *first*, before performing a regularization (see e.g. [135, 150]), which may alter the decrease property of the updating step.

In this chapter, we propose a novel method for constrained optimization which is rather easy to implement and reliable in the sense that it allows to solve (3.1.2) without the need for tuning non physical parameters; this makes it particularly well adapted to the specificities of shape and topology optimization applications. The essence of our method is a modification of the celebrated gradient flow so as to make it able to ‘see the constraints’: optimization trajectories  $x(t)$  are obtained by solving an Ordinary Differential Equation (ODE):

$$\dot{x}(t) = -\alpha_J \boldsymbol{\xi}_J(x(t)) - \alpha_C \boldsymbol{\xi}_C(x(t)). \quad (3.1.3)$$

An admissible local minimizer to (3.1.1) is then reached as the stationary point  $x^*$  of the continuous trajectory  $x(t)$  starting from any (feasible or not) initialization  $x(0)$  and whatever the value of  $\alpha_J, \alpha_C > 0$ . This property is retrieved at the discrete level provided (3.1.3) is discretized with a sufficiently small Euler step size  $\Delta t$ .

The descent direction  $\dot{x}$  is a combination of a so-called ‘null space’ direction  $\boldsymbol{\xi}_J(x)$  and a ‘range space’ direction  $\boldsymbol{\xi}_C(x)$ , lying respectively in the null space of the constraints and in its orthogonal complement (for this reason, we call the ODE (3.1.3) a ‘null space’ gradient flow). The null space direction  $\boldsymbol{\xi}_J(x)$  is the projection of the gradient  $\nabla J(x)$  onto the cone of feasible directions. Our approach relies on a suitable dual program for the resolution of the combinatorial character of the inequality constraints: we

precisely identify those active inequality constraints that the optimization path is allowed to ‘unstick’ from (thus re-entering into the feasible domain) and, conversely, those inequality constraints to which it must remain tangent (see Figure 3.1 below). As a result,  $-\xi_J(x)$  is always the best possible descent direction respecting locally both equality and inequality constraints. The range space direction  $\xi_C(x)$  is a Gauss-Newton direction  $\xi_C(x)$  which is aimed to smoothly lead the optimization path toward the feasible domain. Finally,  $\alpha_J, \alpha_C > 0$  are two (optional) parameters introduced which scale the relative decrease rates of the objective function and of the violation of the constraints; we shall see in particular that the latter quantity decreases along trajectories  $x(t)$  at least as fast as  $e^{-\alpha_C t}$ .

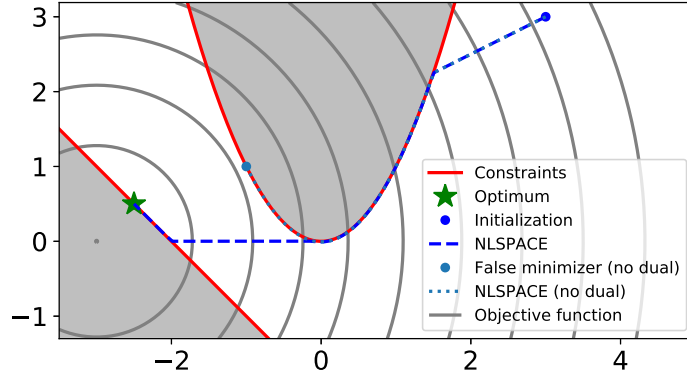


Figure 3.1: An example of optimization trajectory produced by our null space gradient flow (3.1.3). Trajectories travel tangentially to an optimal subset  $\hat{I}(x) \subset \tilde{I}(x)$  of the active constraints  $\tilde{I}(x)$ , which is determined by a dual problem (see section 3.3). A less optimal trajectory is obtained without the identification of the set  $\hat{I}(x)$ , because it is unable to escape the tangent space to the constraints labeled by  $\tilde{I}(x)$ .

More specifically, for a given subset of indices  $I \subset \{1, \dots, q\}$ , denote by  $\mathbf{h}_I(x) := (h_i(x))_{i \in I}$  the collection of inequality constraints and by  $\mathbf{C}_I(x)$  the matrix

$$\mathbf{C}_I(x) := \begin{bmatrix} \mathbf{g}(x) \\ \mathbf{h}_I(x) \end{bmatrix}. \quad (3.1.4)$$

Then, for inequality constrained problems, the directions  $\xi_J(x)$  and  $\xi_C(x)$  in (3.1.3) are defined as follows:

$$\xi_J(x) = (I - \text{DC}_{\tilde{I}(x)}^T (\text{DC}_{\tilde{I}(x)} \text{DC}_{\tilde{I}(x)}^T)^{-1} \text{DC}_{\tilde{I}(x)}) \nabla J(x), \quad (3.1.5)$$

$$\xi_C(x) = \text{DC}_{\tilde{I}(x)}^T (\text{DC}_{\tilde{I}(x)} \text{DC}_{\tilde{I}(x)}^T)^{-1} \mathbf{C}_{\tilde{I}(x)}(x), \quad (3.1.6)$$

where  $I$  is the identity matrix and  $(\text{DC}(x))_{ij} = \partial_j C_i(x)$  denotes the Jacobian matrix of a vector function  $\mathbf{C}(x) = (C_i(x))_i$  (the dependence with respect to  $x$  is omitted when the context is clear). The symbol  $\mathcal{T}$  denotes the transposition operator; it may differ from the usual transpose  $T$  if the optimization set  $V$  is infinite dimensional (see below and section 3.2). Formulas (3.1.5) and (3.1.6) involve two different subsets  $\hat{I}(x) \subset \tilde{I}(x)$  of indices of the inequality constraints  $\{1, \dots, q\}$ : the first one  $\tilde{I}(x)$  is the set of all saturated or violated constraints, defined by

$$\tilde{I}(x) = \{i \in \{1, \dots, q\} \mid h_i(x) \geq 0\}. \quad (3.1.7)$$

The set  $\hat{I}(x) \subset \tilde{I}(x)$  is an optimal subset which is characterized by the following ‘dual’ quadratic optimization subproblem:

$$(\lambda^*(x), \mu^*(x)) := \arg \min_{\substack{\lambda \in \mathbb{R}^p \\ \mu \in \mathbb{R}^{q(x)}, \mu \geq 0}} \|\nabla J(x) + \text{D}\mathbf{g}(x)^T \lambda + \text{D}\mathbf{h}_{\tilde{I}(x)}(x)^T \mu\|_V. \quad (3.1.8)$$

The problem (3.1.8) amounts to compute the projection of  $\nabla J(x)$  onto the cone of feasible directions. It allows to determine the optimal set  $\hat{I}(x)$  from the positive components of the optimal Lagrange multiplier  $\mu^*(x)$ :

$$\hat{I}(x) := \{i \in \tilde{I}(x) \mid \mu_i^*(x) > 0\}. \quad (3.1.9)$$

In [proposition 3.4](#), we show that this definition of  $\widehat{I}(x)$  ensures that  $-\xi_J(x)$  is the ‘best’ descent direction respecting locally both equality and inequality constraints. This key idea to use these two different sets of indices  $\widehat{I}(x)$  and  $\bar{I}(x)$  ensures that  $\xi_J(x)$  is obtained by projection of  $\nabla J(x)$  on the largest possible subspace staying tangential to the set of constraints. This approach turns out to be very efficient efficient in the case of a large number of (possibly violated) inequality constraints ([Figure 3.1](#)). As a result of the flexibility of ODE approaches, our method depends truly only on the discretization step  $\Delta t$  for [\(3.1.3\)](#), and on the physically interpretable, dimensionless parameters  $\alpha_J, \alpha_C$ , which makes them relatively easy to tune for the user.

It turns out that our ODE [\(3.1.3\)](#) is a generalization of rather classical dynamical system approaches to nonlinear constrained optimization, which are maybe less known in the topology optimization community. When the problem [\(3.1.2\)](#) features no constraint, [\(3.1.5\)](#) and [\(3.1.6\)](#) become  $\xi_J(x) = \nabla J(x)$  and  $\xi_C(x) = 0$ , so that the ODE [\(3.1.3\)](#) reduces to the standard gradient flow

$$\dot{x}(t) = -\nabla J(x(t)). \quad (3.1.10)$$

When [\(3.1.2\)](#) also features equality constraints  $\mathbf{g}(x) = 0$ , but no inequality constraint (in that case  $\mathbf{C}_{\widehat{I}(x)}(x) = \mathbf{C}_{\bar{I}(x)}(x) = \mathbf{g}(x)$ ), the same ODE as ours/[\(3.1.3\)](#) was previously derived and studied in the early 1980s by Tanabe [\[300\]](#) (without the Gauss-Newton direction  $\xi_C(x(t))$ ) and by Yamashita [\[317\]](#) (with both  $\xi_J(x(t))$  and  $\xi_C(x(t))$ ). In this particular case, the solution to the dual problem [\(3.1.8\)](#) admits a closed-form expression and [\(3.1.3\)](#) reads with our notation

$$\dot{x} = -\alpha_J(I - \mathbf{Dg}^T(\mathbf{DgDg}^T)^{-1}\mathbf{Dg})\nabla J(x) - \alpha_C\mathbf{Dg}^T(\mathbf{DgDg}^T)^{-1}\mathbf{g}(x). \quad (3.1.11)$$

In the general situation where [\(3.1.2\)](#) features both inequality and equality constraints, variants of the ODE [\(3.1.11\)](#) have been considered by different authors, however with a different method from ours [\[277, 200, 201, 282\]](#). The most common approach in the literature consists in introducing  $q$  slack variables  $\{z_i\}_{1 \leq i \leq q} \in \mathbb{R}^q$  so as to convert the  $q$  inequalities  $h_i(x) \leq 0$  for  $1 \leq i \leq q$  into as many equality constraints  $h_i(x) + z_i^2 = 0$ , before then solving the ODE [\(3.1.3\)](#) in the augmented space  $(x, z) \in V \times \mathbb{R}^q$ . This approach offers convergence guarantees [\[277\]](#) and could also be beneficial for shape optimization, however this is not the strategy we have retained. Indeed, our method does not need to resort to slack variables for handling inequality constraints, and it present small additional advantages described in [section 3.5.1](#).

Another important contribution of this chapter is the exposure of our dynamical system strategy in a setting compatible with the inherently infinite dimensional aspect of shape optimization based on the method of Hadamard. In such context where the optimization variable  $x$  belongs to a Hilbert space  $V$  (if not a more general set, see below), this is achieved by making a clear distinction made between the Fréchet derivative  $\mathbf{D}J(x(t))$  (which is an element of the dual space  $V'$ ) and the gradient  $\nabla J(x(t))$  (which is an element of  $V$ ): the gradient  $\nabla J(x(t))$  is obtained by identification of  $\mathbf{D}J(x(t))$  with an element of  $V$  thanks to the Riesz representation theorem. The same distinction is also needed between the differential of a vector valued function  $\mathbf{DC}(x(t))$  and its transposition  $\mathbf{DC}(x(t))^T$ .

Several works in the field of shape and topology optimization can be related to ours. In fact, our method is very close in spirit to the recent work of Barbarosie et. al. [\[60\]](#), who derived an iterative algorithm for equality constrained optimization which turns out to be a discretization of [\(3.1.3\)](#) with a variable scaling for the parameter  $\alpha_C$ . For inequality constraints, the authors proposed (without convergence results) an active set strategy also based on the extraction of an appropriate subset of the active constraints. However their method relies on a different algorithm from ours, that yields generally a different (suboptimal) set than  $\widehat{I}(x)$ , see [remark 3.7](#) below for more details. Yulin and Xiaoming [\[322\]](#) also suggested to project the gradient of the objective function onto the convex cone of feasible directions; nevertheless, they remained elusive regarding how the projection problem should be solved or how violated constraints should be tackled.

Finally, let us mention as well that the discretization of the flow [\(3.1.3\)](#) can be related to non dynamical system approaches for constrained optimization such as SQP methods (see [\[277\]](#) for a comparison with the ODE approach) and to null space iterative methods [\[69, 70, 244\]](#).

The present chapter is organized as follows. [section 3.2](#) introduces useful notation for distinguishing gradient and differentials for constrained optimization on Hilbert spaces  $V$ . We review as well a few motivations at the origin of the definitions of  $\xi_J(x(t))$  and  $\xi_C(x(t))$  (in [\(3.1.5\)](#) and [\(3.1.6\)](#)) in the case where inequality constraints are absent. In this context, the properties of the flow [\(3.1.3\)](#), classical when

the minimization set  $V$  is a finite dimensional vector space, are reviewed for the more general context where  $V$  is a Hilbert space. We detail then in [section 3.3](#) the necessary adaptations to account for inequality constraints and in particular the introduction of the dual subproblem allowing to determine the null space direction  $\xi_J(x)$ . Under some non restrictive technical assumptions, we prove in [proposition 3.5](#) decreasing properties for the trajectories of the “null space” gradient flow [\(3.1.3\)](#) towards points satisfying the Karush, Kuhn and Tucker (KKT) condition. Numerical implementation and discretization aspects are detailed in [section 3.4](#). [Section 3.5](#) provides pedagogical illustrations of our method on simple academic test cases, and comparisons to a few other optimization methods of the literature including the method of slack variables for inequality constraints. Shape optimization applications are eventually considered in [section 3.6](#). After clarifying the necessary adaptations required to extend the discretization of [\(3.1.3\)](#) to sets featuring a rather generic manifold structure, we explain how our algorithm can be integrated within the level set method for shape optimization [[311](#), [32](#), [24](#)]. The ease of implementation and efficiency of our method is demonstrated numerically for the optimal design of a bridge structure subject to multiple loads, which involves up to ten constraints.

## 3.2 NULL SPACE GRADIENT FLOWS FOR EQUALITY-CONSTRAINED OPTIMIZATION IN HILBERT SPACES

In this section, we consider the case where the optimization takes place on a Hilbert space  $V$  with inner product  $\langle \cdot, \cdot \rangle_V$  and relative norm  $\|\cdot\|_V = \langle \cdot, \cdot \rangle_V^{1/2}$ ; see [section 3.6](#) for the description of the more general situation associated to our shape optimization applications. In this part only, we consider the problem [\(3.1.2\)](#) where only equality constraints are present, namely:

$$\begin{aligned} \min_{x \in V} \quad & J(x) \\ \text{s.t.} \quad & \mathbf{g}(x) = 0, \end{aligned} \tag{3.2.1}$$

where  $J : V \rightarrow \mathbb{R}$  and  $\mathbf{g} : V \rightarrow \mathbb{R}^p$  are Fréchet differentiable functions. The purpose of this section is to motivate the introduction of the ODE [\(3.1.3\)](#) for equality constrained optimization, and to review its properties in the present Hilbertian setting. Let us emphasize that, although this section is elementary and not completely new, it is not easily found as is in the literature. Since it is key in understanding our technique for handling inequality constraints in [section 3.3](#), the present context is thoroughly detailed for the reader’s convenience.

The section is organized as follows. [Section 3.2.1](#) recalls definitions for the differential, the gradient, and the transpose operation in the Hilbertian context. We then sketch briefly in [section 3.2.2](#) how the formulas [\(3.1.5\)](#) and [\(3.1.6\)](#) can be formally obtained, before stating the properties of the null space step  $\xi_J(x)$  and its relation to Lagrange multipliers by means of a dual problem in [lemma 3.1](#). Finally, the decrease properties of the obtained dynamical system are reviewed in [section 3.2.3](#).

### 3.2.1 Notation and first-order optimality conditions

The following definition sets the notation conventions about differentiability and gradients in Hilbert spaces used throughout this chapter. Note that they may differ from those used by other authors because a clear distinction is needed for our shape optimization purposes between gradient and Fréchet derivatives.

**Definition 3.1.** 1. A vector-valued function  $\mathbf{g} : V \rightarrow \mathbb{R}^p$  is *differentiable* at a point  $x \in V$  if there exists a continuous linear mapping  $D\mathbf{g}(x) : V \rightarrow \mathbb{R}^p$  such that

$$\mathbf{g}(x+h) = \mathbf{g}(x) + D\mathbf{g}(x)h + o(h) \text{ with } \frac{o(h)}{\|h\|_V} \xrightarrow{h \rightarrow 0} 0. \tag{3.2.2}$$

$D\mathbf{g}(x)$  is called the Fréchet derivative of  $\mathbf{g}$  at  $x$ .

2. If  $\mathbf{g} : V \rightarrow \mathbb{R}^p$  is differentiable, for any  $\boldsymbol{\mu} \in \mathbb{R}^p$ , the Riesz representation theorem [[76](#)] ensures the existence of a unique vector  $D\mathbf{g}(x)^T \boldsymbol{\mu} \in V$  satisfying

$$\forall \boldsymbol{\mu} \in \mathbb{R}^p, \forall \boldsymbol{\xi} \in V, \langle D\mathbf{g}(x)^T \boldsymbol{\mu}, \boldsymbol{\xi} \rangle_V = \boldsymbol{\mu}^T D\mathbf{g}(x) \boldsymbol{\xi}, \tag{3.2.3}$$

where the superscript  $T$  stands for the usual transpose of a vector in the Euclidean space  $\mathbb{R}^p$ . The linear operator  $D\mathbf{g}(x)^T : \mathbb{R}^p \rightarrow V$  thus defined is called the *transpose* of  $D\mathbf{g}(x)$ .

3. If  $J : V \rightarrow \mathbb{R}$  is a scalar function differentiable at  $x \in V$ , the Riesz representation theorem ensures the existence of a unique vector  $\nabla J(x) \in V$  satisfying

$$\forall \xi \in V, \langle \nabla J(x), \xi \rangle_V = DJ(x)\xi. \quad (3.2.4)$$

This vector  $\nabla J(x)$  is called the *gradient* of  $J$  at  $x$ .

Throughout the chapter, the explicit mention to  $x$  is sometimes omitted in the notation for differentials or gradients when the considered point  $x \in V$  is clear.

**Remark 3.1.** 1. If  $V$  is the (finite-dimensional) Euclidean space  $\mathbb{R}^k$ , equipped with the standard inner product, the Fréchet derivative and the transpose of the differential of a vector valued function  $\mathbf{g} : \mathbb{R}^k \rightarrow \mathbb{R}^p$  are respectively given by the Jacobian matrix  $(D\mathbf{g})_{ij} = \partial_j g_i$  and its transpose  $(D\mathbf{g}^T)_{ij} = (D\mathbf{g})_{ji} = \partial_i g_j$ . In the literature, the differential matrix  $D\mathbf{g}$  is often denoted with the nabla symbol  $\nabla \mathbf{g}$ . For the sake of clarity, we reserve the  $\nabla$  symbol to denote the gradient of scalar functions  $J : V \rightarrow \mathbb{R}$ , and it holds that  $\nabla J(x) = DJ(x)^T \mathbf{1}$ . The calligraphic transpose notation  $\mathcal{T}$  appearing in the objects  $DJ(x)^{\mathcal{T}}$  or  $D\mathbf{g}(x)^{\mathcal{T}}$  encodes at the same time the operator transposition (reversing the input and range spaces) and the Riesz identifications.

2. Still in the case where  $V = \mathbb{R}^k$  is finite-dimensional and  $a$  is given by a symmetric definite positive matrix  $A$  (that is  $\langle \xi, \xi \rangle_V = \xi^T A \xi$ ), the transpose of a  $p \times k$  matrix  $M : \mathbb{R}^k \rightarrow \mathbb{R}^p$  with respect to  $a$  is  $M^{\mathcal{T}} = A^{-1} M^T$ . As we shall see in [section 3.6](#), in shape optimization applications,  $\langle \cdot, \cdot \rangle_V$  often stands for the bilinear form associated to an elliptic operator, hence the calligraphic transpose  $\mathcal{T}$  encompasses the extension and regularization step of the shape derivative outlined in [section 1.4.1](#). If  $V$  is the tangent space to some Riemannian manifold, the inner product  $\langle \cdot, \cdot \rangle_V$  can be interpreted as a metric and  $\nabla J(x)$ , as given by [\(3.2.4\)](#), is the covariant gradient with respect to this metric.
3. When  $V$  is a general Hilbert space, for a vector-valued function  $\mathbf{g} : V \rightarrow \mathbb{R}^p$  with coordinates  $\mathbf{g}(x) = (g_i(x))_{1 \leq i \leq p}$ ,  $D\mathbf{g} : V \rightarrow \mathbb{R}^p$  is the ‘row’ matrix whose entries are the  $p$  linear forms  $Dg_i(x) : V \rightarrow \mathbb{R}$ . The transpose  $D\mathbf{g}(x)^{\mathcal{T}}$  is the ‘column’ matrix gathering the  $p$  vectors  $(\nabla g_i(x))_{1 \leq i \leq p}$  obtained by solving the  $p$  identification problems:

$$\forall \xi \in V, \langle \nabla g_i(x), \xi \rangle_V = Dg_i(x)\xi; \quad (3.2.5)$$

more precisely:

$$\forall \mu \in \mathbb{R}^d, D\mathbf{g}(x)^{\mathcal{T}} \mu = \sum_{i=1}^p \mu_i \nabla g_i(x).$$

In particular, the  $p \times p$  matrix  $D\mathbf{g}D\mathbf{g}^{\mathcal{T}} \in \mathbb{R}^{p \times p}$  has entries

$$(D\mathbf{g}D\mathbf{g}^{\mathcal{T}})_{ij} = \langle \nabla g_i, \nabla g_j \rangle_V = Dg_i(x)(\nabla g_j(x)).$$

Throughout this section, the equality constraints defined by the function  $\mathbf{g}$  are said to be qualified at a point  $x \in V$  if

$$\text{rank}(D\mathbf{g}(x)) = p, \text{ or equivalently } D\mathbf{g}(x)D\mathbf{g}(x)^{\mathcal{T}} \text{ is an invertible } p\text{-by-}p \text{ matrix.} \quad (3.2.6)$$

Note that [\(3.2.6\)](#) makes sense even at points  $x \in V$  where  $\mathbf{g}(x) \neq 0$ , a fact that we shall use extensively in the sequel.

Let us then recall the classical definitions of critical points in terms of the first-order necessary optimality conditions (KKT) for the equality-constrained problem [\(3.2.1\)](#):

**Definition 3.2** (see [\[70, 244\]](#)). A point  $x^*$  is said to satisfy the Karush, Kuhn and Tucker conditions (KKT) for the equality constrained minimization problem [\(3.2.1\)](#) if and only if there exists  $\lambda(x^*) \in \mathbb{R}^p$  such that:

$$\begin{cases} \nabla J(x^*) + D\mathbf{g}(x^*)^{\mathcal{T}} \lambda(x^*) = 0, \\ \mathbf{g}(x^*) = 0. \end{cases} \quad (3.2.7)$$



### 3.2.2 Definitions and properties of the null space and range space steps $\xi_J$ and $\xi_C$

We are now in position to define the null space and range space steps  $\xi_J(x)$  and  $\xi_C(x)$  featured in the ODE (3.1.11) for equality constrained problems in the present Hilbert space setting, owing to the definition (3.2.3) of the transpose:

**Definition 3.3.** For any point  $x \in V$  satisfying the constraint qualification condition (3.2.6), the *null space* and *range space* directions  $\xi_J(x)$  and  $\xi_C(x)$  associated to (3.2.1) are respectively defined by:

$$\xi_J(x) := (I - \text{Dg}^\mathcal{T}(\text{DgDg}^\mathcal{T})^{-1}\text{Dg})\nabla J(x), \quad (3.2.8)$$

$$\xi_C(x) := \text{Dg}^\mathcal{T}(\text{DgDg}^\mathcal{T})^{-1}\mathbf{g}(x). \quad (3.2.9)$$

Let us provide a formal intuition motivating the expressions (3.2.8), (3.2.9) and the dynamical system (3.1.3). A classical practice in Lagrange multiplier methods for optimization which is inspired from the KKT optimality conditions (3.2.7) is to search for an iterative optimization scheme (indexed by the iteration number  $n$ ) of the form (see also [59])

$$x_{n+1} = x_n - \Delta t(\alpha_J \nabla J(x_n) + \text{Dg}(x_n)^\mathcal{T} \lambda_n), \quad (3.2.10)$$

where  $\lambda_n \in \mathbb{R}^p$  is a tentative value for the Lagrange multiplier  $\lambda$  in (3.2.7),  $\alpha_J$  is a user-defined coefficient and  $\Delta t$  is the step increment between successive iterations. We determine the value of  $\lambda_n$  by imposing that the value  $\mathbf{g}(x_{n+1})$  of the constraint decrease by a factor  $(1 - \alpha_C \Delta t)$ , up to some lower order term in  $\Delta t$ . Since

$$\mathbf{g}(x_{n+1}) = \mathbf{g}(x_n) - \Delta t \text{Dg}(x_n)(\alpha_J \nabla J(x_n) + \text{Dg}(x_n)^\mathcal{T} \lambda_n) + o(\Delta t),$$

the requirement that  $\mathbf{g}(x_{n+1}) \simeq (1 - \alpha_C \Delta t)\mathbf{g}(x_n)$  suggests the rule:

$$\lambda_n = (\text{Dg}(x_n)\text{Dg}(x_n)^\mathcal{T})^{-1}(\alpha_C \mathbf{g}(x_n) - \alpha_J \text{Dg}(x_n)\nabla J(x_n)). \quad (3.2.11)$$

We recognize then a time discretization scheme of the ODE (3.1.11) by replacing  $\lambda_n$  with the above value (3.2.11) in (3.2.10).

#### Properties of the null space step $\xi_J$

In the finite-dimensional case where  $V = \mathbb{R}^k$ , it is well-known that the null space step  $\xi_J(x)$  defined by (3.2.8) is the orthogonal projection of the gradient  $\nabla J(x)$  of the objective function onto the null space of the constraints

$$\text{Ker}(\text{Dg}(x)) = \{\xi \in V \mid \text{Dg}(x)\xi = 0\},$$

which is also the tangent space at  $x$  to the manifold  $\{y \in V \mid \mathbf{g}(y) = \mathbf{g}(x)\}$ . Of course, this is still true when  $V$  is a Hilbert space, as recalled in the next lemma.

**Lemma 3.1.** *Let  $x \in V$  be a point satisfying the qualification condition (3.2.6) The following properties hold:*

1. *The space  $V$  has the following orthogonal decomposition:*

$$V = \text{Ker}(\text{Dg}(x)) \oplus \text{Ran}(\text{Dg}(x)^\mathcal{T}),$$

where we have introduced the range  $\text{Ran}(\text{Dg}(x)^\mathcal{T}) := \{\text{Dg}(x)^\mathcal{T} \lambda \mid \lambda \in \mathbb{R}^p\}$  of  $\text{Dg}(x)^\mathcal{T}$ .

Moreover, the operator  $\Pi_{\text{g}(x)} : V \rightarrow V$  defined by

$$\Pi_{\text{g}(x)} = I - \text{Dg}^\mathcal{T}(\text{DgDg}^\mathcal{T})^{-1}\text{Dg}(x) \quad (3.2.12)$$

is the orthogonal projection onto  $\text{Ker}(\text{Dg}(x))$ .

2. *When  $\Pi_{\text{g}(x)}(\nabla J(x)) \neq 0$ ,  $-\xi_J(x) = -\Pi_{\text{g}(x)}(\nabla J(x))$  is the best normalized feasible descent direction for  $J$  in the sense that*

$$-\frac{\xi_J(x)}{\|\xi_J(x)\|_V} = \arg \min_{\xi \in V} \text{DJ}(x)\xi \quad (3.2.13)$$

$$\text{s.t.} \begin{cases} \text{Dg}(x)\xi = 0 \\ \langle \xi, \xi \rangle_V \leq 1. \end{cases}$$

3. The null space direction  $\xi_J(x) = \Pi_{g(x)}(\nabla J(x))$  is the closest least squares approximation to  $\nabla J(x)$  within the space  $\text{Ker}(\mathbf{Dg}(x))$ . It alternatively reads

$$\xi_J(x) = \nabla J(x) + \mathbf{Dg}(x)^T \lambda^*(x), \quad (3.2.14)$$

where the Lagrange multiplier  $\lambda^*(x) := -(\mathbf{DgDg}^T)^{-1} \mathbf{Dg} \nabla J(x)$  is the unique solution to the following least squares problem that is the dual of (3.2.13):

$$\lambda^*(x) = \arg \min_{\lambda \in \mathbb{R}^p} \|\nabla J(x) + \mathbf{Dg}(x)^T \lambda\|_V. \quad (3.2.15)$$

**Remark 3.2.** The qualification condition (3.2.6) is essentially used in the computation of the orthogonal projection onto  $\text{Ker}(\mathbf{Dg}(x))$  from formula (3.2.12). Most of the statements of this chapter could circumvent this hypothesis by relying on a Singular Value Decomposition (SVD) of  $\mathbf{Dg}(x)$ .

*Proof.* 1. Any  $\xi \in V$  may be decomposed as  $\xi = \Pi_{g(x)}(\xi) + (I - \Pi_{g(x)})(\xi)$ , where it is straightforward to verify that  $\Pi_{g(x)}(\xi) \in \text{Ker}(\mathbf{Dg}(x))$ , and  $(I - \Pi_{g(x)})(\xi) \in \text{Ran}(\mathbf{Dg}(x)^T)$ . In addition,  $\text{Ker}(\mathbf{Dg}(x))$  and  $\text{Ran}(\mathbf{Dg}(x)^T)$  are orthogonal for the inner product  $\langle \cdot, \cdot \rangle_V$  since from (3.2.3), one has,

$$\forall \zeta \in \text{Ker}(\mathbf{Dg}(x)), \forall \lambda \in \mathbb{R}^p, \langle \mathbf{Dg}(x)^T \lambda, \zeta \rangle_V = \lambda^T \mathbf{Dg}(x) \zeta = 0.$$

2. It follows from the first point that for any  $\xi \in \text{Ker}(\mathbf{Dg}(x))$  such that  $\|\xi\|_V \leq 1$ ,

$$DJ(x)\xi = \langle \nabla J(x), \xi \rangle_V = \langle \Pi_{g(x)}(\nabla J(x)), \xi \rangle_V \geq -\|\Pi_{g(x)}(\nabla J(x))\|_V,$$

whence we easily infer that  $\xi := -\Pi_{g(x)}(\nabla J(x)) / \|\Pi_{g(x)}(\nabla J(x))\|_V$  is the global minimizer of (3.2.13).

3. The Pythagore identity yields, for any  $\xi \in \text{Ker}(\mathbf{Dg}(x))$ ,

$$\|\nabla J(x) - \xi\|_V^2 = \|(I - \Pi_{g(x)})\nabla J(x)\|_V^2 + \|\Pi_{g(x)}\nabla J(x) - \xi\|_V^2 \geq \|\nabla J(x) - \Pi_{g(x)}\nabla J(x)\|_V^2.$$

Hence the orthogonal projection  $\Pi_{g(x)}(\nabla J(x))$  is the best approximation of  $\nabla J(x)$  within the space  $\text{Ker}(\mathbf{Dg}(x))$ . Recalling from the first point that the range  $\text{Ran}(\mathbf{Dg}(x)^T)$  is the orthogonal complement of  $\text{Ker}(\mathbf{Dg}(x))$ , we obtain also, for any  $\lambda \in \mathbb{R}^p$ ,

$$\|\Pi_{g(x)}(\nabla J(x))\|_V = \|\nabla J(x) - (I - \Pi_{g(x)})(\nabla J(x))\|_V \leq \|\nabla J(x) + \mathbf{Dg}(x)^T \lambda\|_V,$$

whence the expression (3.2.14) and the minimization property (3.2.15) follow. Note that the uniqueness of the solution  $\lambda^*(x)$  to (3.2.15) results from the qualification condition (3.2.6).

Finally, the optimization problem (3.2.13) can be rewritten as

$$\min_{\substack{\xi \in V \\ \langle \xi, \xi \rangle_V \leq 1}} \max_{\lambda \in \mathbb{R}^p} DJ(x)\xi + \lambda^T \mathbf{Dg}(x)\xi.$$

Hence the (formal) dual problem of (3.2.13) reads:

$$\max_{\lambda \in \mathbb{R}^p} \min_{\substack{\xi \in V \\ \langle \xi, \xi \rangle_V \leq 1}} DJ(x)\xi + \lambda^T \mathbf{Dg}(x)\xi.$$

According to the definitions (3.2.3) and (3.2.4) of the gradient and of the Hilbertian transpose, the latter problem rewrites:

$$\max_{\lambda \in \mathbb{R}^p} \min_{\substack{\xi \in V \\ \langle \xi, \xi \rangle_V \leq 1}} \langle \nabla J(x) + \mathbf{Dg}(x)^T \lambda, \xi \rangle_V = - \min_{\lambda \in \mathbb{R}^p} \|\nabla J + \mathbf{Dg}^T \lambda\|_V,$$

where for given  $\lambda \in \mathbb{R}^p$ , the value

$$\xi^* := - \frac{\nabla J(x) + \mathbf{Dg}(x)^T \lambda}{\|\nabla J(x) + \mathbf{Dg}(x)^T \lambda\|_V}$$

is that achieving the minimum in the minimization problem at the left-hand side of the above identity. This shows that (3.2.15) is the dual problem of (3.2.13).  $\square$

### Properties of the range space step $\xi_C$

The next lemma characterizes the range space step  $\xi_C(x)$ , defined by (3.2.9), as the unique Gauss-Newton direction for the minimization of the constraint function  $\mathbf{g}(x)$  which is tangential to the (linearized) set of constraints:

**Lemma 3.2.** *Let  $x \in V$  satisfy the qualification condition (3.2.6); then:*

1. *The range space step  $\xi_C(x) = \mathbf{Dg}^T(\mathbf{DgDg}^T)^{-1}\mathbf{g}(x)$  is orthogonal to  $\text{Ker}(\mathbf{Dg}(x))$ :*

$$\forall \xi \in \text{Ker}(\mathbf{Dg}(x)), \quad \langle \xi_C(x), \xi \rangle_V = 0.$$

2.  *$-\xi_C(x)$  is a descent direction for the violation of the constraints:*

$$\mathbf{Dg}(x)(-\xi_C(x)) = -\mathbf{g}(x). \quad (3.2.16)$$

3. *The set of solutions to the Gauss-Newton program*

$$\min_{\xi \in V} \|\mathbf{g}(x) + \mathbf{Dg}(x)\xi\|^2 \quad (3.2.17)$$

*is the affine subspace  $\{-\xi_C(x) + \zeta \mid \zeta \in \text{Ker}(\mathbf{Dg}(x))\}$  of  $V$ .*

*Proof.* 1. This easily follows from the point 1. of lemma 3.1.

2. This is an immediate consequence of the definition (3.2.9) of  $\xi_C(x)$ . Note that (3.2.16) means that  $-\xi_C(x)$  is a descent direction for the violation of the constraints in the sense that it ensures that any coordinate  $g_i(x)$ ,  $i = 1, \dots, p$ , decreases along  $-\xi_C(x)$  if  $g_i(x) \geq 0$  and increases if  $g_i(x) \leq 0$ .
3. Since (3.2.17) is a convex optimization problem, a necessary and sufficient condition for  $\xi \in V$  to be one solution is given by the usual first-order condition:

$$\forall \zeta \in V, \quad (\mathbf{g}(x) + \mathbf{Dg}(x)\xi)^T(\mathbf{Dg}(x)\zeta) = \langle \mathbf{Dg}(x)^T(\mathbf{g}(x) + \mathbf{Dg}(x)\xi), \zeta \rangle_V = 0,$$

which rewrites:

$$\mathbf{Dg}(x)^T \mathbf{Dg}(x)\xi = -\mathbf{Dg}(x)^T \mathbf{g}(x).$$

Since the matrix  $(\mathbf{DgDg}^T)$  is invertible (as a consequence of the qualification condition (3.2.6)), this is in turn equivalent to:

$$\mathbf{Dg}(x)\xi = -\mathbf{g}(x).$$

Finally, (3.2.16) states that  $-\xi_C(x)$  is one particular solution to the above equation; therefore, any two solutions of this problem differ by some  $\zeta$  such that  $\mathbf{Dg}(x)\zeta = 0$ . □

### 3.2.3 Decrease properties of the equality constrained gradient flow

The main features of the definitions of  $\xi_J(x)$  and  $\xi_C(x)$  are the facts that  $\xi_J$  is orthogonal to the set of constraints, i.e.  $\mathbf{Dg}(x)\xi_J(x) = 0$ , and that  $-\xi_C(x)$  decreases the violation of the constraints while being orthogonal to  $\xi_J(x)$ . These ensure that the values of the constraint functional  $\mathbf{g}(x(t))$  decrease to zero along the trajectories of the ODE (3.1.3), independently of the behavior of  $\xi_J(x)$ . Then, as soon as the violation of the constraint becomes sufficiently small, the objective function  $J$  decreases without affecting the asymptotic vanishing of  $\mathbf{g}(x(t))$ . We review these properties in the next proposition, which was also observed in [317] in the finite-dimensional context.

**Proposition 3.1.** *Assume that the trajectories  $x(t)$  of the flow*

$$\begin{cases} \dot{x} = -\alpha_J(I - \mathbf{Dg}^T(\mathbf{DgDg}^T)^{-1}\mathbf{Dg}(x))\nabla J(x) - \alpha_C \mathbf{Dg}^T(\mathbf{DgDg}^T)^{-1}\mathbf{g}(x) \\ x(0) = x_0 \end{cases} \quad (3.2.18)$$

*exist on some time interval  $[0, T]$  for  $T > 0$ , and that the qualification condition (3.2.6) holds at any point  $x(t)$ ,  $t \in [0, T]$ . Then the following properties hold true:*

1. *The violation of the constraints decreases exponentially:*

$$\forall t \in [0, T], \quad \mathbf{g}(x(t)) = e^{-\alpha_C t} \mathbf{g}(x_0). \quad (3.2.19)$$

2.  $J(x(t))$  decreases ‘as soon as the violation (3.2.19) of the constraints is sufficiently small’ in the following sense: assume that  $\text{rank}(\mathbf{Dg}) = p$  on  $K = \{x \in V \mid \|\mathbf{g}(x)\|_\infty \leq \|\mathbf{g}(x_0)\|_\infty\}$  and that

$$\sup_{x \in K} \|\nabla J(x)\|_V |\sigma_p^{-1}(x)| < +\infty, \quad (3.2.20)$$

where  $\sigma_p(x)$  is the smallest singular value of  $\mathbf{Dg}(x)$ . Then there exists a constant  $C > 0$  such that

$$\forall t \in [0, T], \|\Pi_{g(x)}(\nabla J(x(t)))\|_V^2 > Ce^{-\alpha_C t} \Rightarrow \frac{d}{dt} J(x(t)) < 0. \quad (3.2.21)$$

3. Any stationary point  $x^*$  of (3.2.18) satisfies the first-order KKT conditions (3.2.7) of the optimization program (3.2.1), that is:

$$\begin{cases} \mathbf{g}(x^*) = 0 \\ \exists \boldsymbol{\lambda}^* \in \mathbb{R}^p, \nabla J(x^*) + \mathbf{Dg}^\mathcal{T}(x^*)\boldsymbol{\lambda}^* = \Pi_{g(x^*)}(\nabla J(x^*)) = 0. \end{cases} \quad (3.2.22)$$

*Proof.* 1. Using the definition (3.2.18), the decreasing property (3.2.16) together with the fact that  $\boldsymbol{\xi}_J(x)$  is orthogonal to  $\text{Ker}(\mathbf{Dg}(x))$ , we obtain:

$$\frac{d}{dt} (\mathbf{g}(x(t))) = -\alpha_C \mathbf{g}(x(t)),$$

whence (3.2.19) follows easily.

2. Let us introduce the eigenvalue decomposition

$$\mathbf{Dg}(x)\mathbf{Dg}(x)^\mathcal{T} = \sum_{i=1}^p \sigma_i(x)^2 \mathbf{u}_i(x)\mathbf{u}_i(x)^\mathcal{T}, \text{ where } \sigma_1(x) \geq \dots \geq \sigma_p(x) > 0, \mathbf{u}_i(x)^\mathcal{T} \mathbf{u}_j(x) = \delta_{ij},$$

of the symmetric, positive definite  $p \times p$  matrix  $\mathbf{Dg}(x)\mathbf{Dg}(x)^\mathcal{T}$ . Let then  $\mathbf{v}_i(x)^\dagger : V \rightarrow \mathbb{R}$  be the linear form defined for any  $\boldsymbol{\xi} \in V$  by  $\mathbf{v}_i(x)^\dagger \boldsymbol{\xi} = \sigma_i(x)^{-1} \mathbf{u}_i(x)^\mathcal{T} \mathbf{Dg}(x)\boldsymbol{\xi}$  and let  $\mathbf{v}_i(x)$  be the vector in  $V$  such that  $\forall \boldsymbol{\xi} \in V, \langle \mathbf{v}_i(x), \boldsymbol{\xi} \rangle_V = \mathbf{v}_i(x)^\dagger \boldsymbol{\xi}$ ; more explicitly,  $\mathbf{v}_i(x) = \sigma_i(x)^{-1} \mathbf{Dg}(x)^\mathcal{T} \mathbf{u}_i(x)$ . These definitions allow to write a singular value decomposition for  $\mathbf{Dg}(x)$ ; it is indeed easily verified from the definitions of  $\mathbf{u}_i(x)$  and  $\mathbf{v}_i(x)$  that:

$$\mathbf{Dg}(x) = \sum_{i=1}^p \sigma_i(x) \mathbf{u}_i(x) \mathbf{v}_i(x)^\dagger, \text{ and } \mathbf{Dg}(x)^\mathcal{T} = \sum_{i=1}^p \sigma_i(x) \mathbf{v}_i(x) \mathbf{u}_i(x)^\mathcal{T}$$

with  $\langle \mathbf{v}_i(x), \mathbf{v}_j(x) \rangle_V = \mathbf{v}_i(x)^\dagger \mathbf{v}_j(x) = \delta_{ij}$ . We now calculate:

$$\mathbf{Dg}^\mathcal{T} (\mathbf{Dg}\mathbf{Dg}^\mathcal{T})^{-1} \mathbf{g}(x) = \sum_{i=1}^p \sigma_i^{-1}(x) (\mathbf{u}_i(x)^\mathcal{T} \mathbf{g}(x)) \mathbf{v}_i(x),$$

whence we obtain the following inequality:

$$\forall x \in V, |\mathbf{Dg}^\mathcal{T} (\mathbf{Dg}\mathbf{Dg}^\mathcal{T})^{-1} \mathbf{g}| \leq \sigma_p^{-1}(x) \|\nabla J(x)\|_V \|\mathbf{g}(x)\|. \quad (3.2.23)$$

Since

$$\frac{d}{dt} J(x(t)) = -\alpha_J \mathbf{Dg}^\mathcal{T} (\mathbf{Dg}\mathbf{Dg}^\mathcal{T})^{-1} \mathbf{g}_J(x(t)) - \alpha_C \mathbf{Dg}^\mathcal{T} (\mathbf{Dg}\mathbf{Dg}^\mathcal{T})^{-1} \mathbf{g}_C(x(t)),$$

it follows that  $\frac{d}{dt} J(x(t)) < 0$  as soon as  $\alpha_J |\mathbf{Dg}^\mathcal{T} (\mathbf{Dg}\mathbf{Dg}^\mathcal{T})^{-1} \mathbf{g}_J(x(t))| > \alpha_C |\mathbf{Dg}^\mathcal{T} (\mathbf{Dg}\mathbf{Dg}^\mathcal{T})^{-1} \mathbf{g}_C(x(t))|$ . Thus, from (3.2.19) and (3.2.23), the constant  $C$  in (3.2.21) can be selected as

$$C = p \frac{\alpha_C}{\alpha_J} \|\mathbf{g}(x_0)\| \sup_{x \in K} [\sigma_p^{-1}(x) \|\nabla J(x)\|_V]. \quad (3.2.24)$$

3. Since the vectors  $\boldsymbol{\xi}_J(x)$  and  $\boldsymbol{\xi}_C(x)$  are orthogonal for any point  $x \in V$ , a stationary point  $x^*$  of (3.2.18) must satisfy

$$\Pi_{g(x^*)}(\nabla J(x^*)) = 0, \text{ and } \mathbf{Dg}^\mathcal{T} (\mathbf{Dg}\mathbf{Dg}^\mathcal{T})^{-1} \mathbf{g}(x^*) = 0, \quad (3.2.25)$$

and so the first KKT condition in (3.2.7) is satisfied with the value  $\boldsymbol{\lambda} = -(\mathbf{Dg}\mathbf{Dg}^\mathcal{T})^{-1} \mathbf{Dg}(x^*) \nabla J(x^*)$  of the Lagrange multiplier. Then left multiplication by  $\mathbf{Dg}$  in the second identity in (3.2.25) implies  $\mathbf{g}(x^*) = 0$ , which completes the proof.  $\square$

**Remark 3.3.** The solutions to the dynamical system (3.2.18) are defined for small time if  $\xi_J$  and  $\xi_C$  are locally Lipschitz vector fields, which is the case if e.g.  $J$  and  $\mathbf{g}$  are of class  $\mathcal{C}^2$  [127]. In the case where  $V$  is finite-dimensional, the assumption (3.2.20) is satisfied if the set  $K = \{x \in V \mid \mathbf{g}(x) \leq \mathbf{g}(x_0)\}$  is bounded and the functions  $J$  and  $\mathbf{g}$  are  $\mathcal{C}^1$  functions. It is worth noting that even if the present regularity of  $J$  and  $\mathbf{g}$  is not strong enough to ensure the existence of solutions to (3.2.18), similar properties to those of proposition 3.1 hold for the discretized scheme

$$x_{n+1} = x_n - \Delta t(\alpha_J \xi_J(x_n) + \alpha_C \xi_C(x_n)), \quad (3.2.26)$$

which is sufficient for optimization. One can indeed verify that:

1. At first order, the constraints decrease with a geometric rate:  $\mathbf{g}(x_{n+1}) = (1 - \alpha_C \Delta t) \mathbf{g}(x_n) + o(\Delta t)$ .
2. An accumulation point  $x^*$  of the sequence  $(x_n)_{n \in \mathbb{N}}$  satisfies  $\mathbf{g}(x^*) = 0$  and is a KKT point of the problem (3.2.1), satisfying (3.2.7).

**Remark 3.4.** It is possible to control more accurately the pace at which each of the constraints decreases in (3.2.19): consider a diagonal matrix of positive coefficients  $\mathbf{K} = \text{diag}(\kappa_i)_{1 \leq i \leq p}$  and replace the definition (3.2.9) of  $\xi_C(x)$  by

$$\xi_C(x) := \mathbf{Dg}^T (\mathbf{Dg} \mathbf{Dg}^T)^{-1} \mathbf{K} \mathbf{g}(x).$$

Then each constraint function  $g_i$  decreases at its own rate  $\kappa_i \alpha_C$  along the solution  $x(t)$  of (3.2.18):

$$\forall t \in [0, T], g_i(x(t)) = e^{-\kappa_i \alpha_C t} g_i(x_0).$$

### 3.3 EXTENSION TO EQUALITY AND INEQUALITY CONSTRAINTS

We now proceed to extend the dynamical system (3.1.3) or (3.2.18) so as to handle inequality constraints as well. We return to the full optimization problem (3.1.2), still posed in a Hilbert space  $V$  with inner product  $\langle \cdot, \cdot \rangle_V$ , and where the objective  $J : V \rightarrow \mathbb{R}$ , equality constraints  $\mathbf{g} : V \rightarrow \mathbb{R}^p$  and inequality constraints  $\mathbf{h} : V \rightarrow \mathbb{R}^q$  are differentiable functions.

Inspired by the methodology developed in section 3.2, we still propose to solve the equality and inequality constrained problem (3.1.2) thanks to a dynamical system of the form:

$$\begin{cases} \dot{x}(t) = -\alpha_J \xi_J(x(t)) - \alpha_C \xi_C(x(t)) \\ x(0) = x_0, \end{cases} \quad (3.3.1)$$

whose discretized version reads:

$$x_{n+1} = x_n - \Delta t(\alpha_J \xi_J(x_n) + \alpha_C \xi_C(x_n)). \quad (3.3.2)$$

In what follows, notation conventions related to index sets associated to inequality constraints are introduced in section 3.3.1. The range space step  $\xi_C(x)$  is defined in section 3.3.2 from a formula analogous to (3.1.6). The definition of the null space step  $\xi_J(x)$  is examined in details in section 3.3.3; it involves a procedure discriminating a relevant subset  $\tilde{I}(x) \subset \hat{I}(x)$  of the saturated or violated constraints, which relies on the introduction of the dual problem (3.1.8). Finally, the properties of the flow (3.3.1) are outlined in section 3.3.4.

#### 3.3.1 Notation and preliminaries

For the convenience of the reader, a few notation conventions from the introduction are now made more precise. The set of indices of saturated or violated inequality constraints at  $x \in V$  is denoted by  $\tilde{I}(x)$ :

$$\tilde{I}(x) = \{i \in \{1, \dots, q\} \mid h_i(x) \geq 0\}, \quad (3.3.3)$$

and  $\tilde{q}(x) := \text{Card}(\tilde{I}(x))$  is the number of such constraints. For a subset  $I \subset \{1, \dots, q\}$ , the vector  $\mathbf{h}_I(x) = (h_i(x))_{i \in I}$  collects the inequality constraints indexed by  $I$  and  $\mathbf{C}_I(x)$ , defined by (3.1.4), collects all equality constraints  $\mathbf{g}(x)$  and those selected inequality constraints  $\mathbf{h}_I(x)$ .

In the present context, the constraints are said to be qualified at  $x \in V$  if the linearized saturated or violated constraints are independent, that is,

$$\text{rank}(\mathbf{DC}_{\tilde{I}(x)}(x)) = p + \tilde{q}(x). \quad (3.3.4)$$

If the point  $x$  satisfies the constraints, (3.3.4) is one usual qualification condition (of course, there are other possible qualification conditions, see [70, 244]), but let us insist again that we leave the room for this definition to apply to points  $x$  where constraints are not satisfied. Define  $\Pi_{C_I} : V \rightarrow V$ , the orthogonal projection operator onto  $\text{Ker}(\text{DC}_I(x))$ , by

$$\Pi_{C_I} = I - \text{DC}_I(x)^T (\text{DC}_I(x) \text{DC}_I(x)^T)^{-1} \text{DC}_I(x), \quad (3.3.5)$$

and let  $(\boldsymbol{\lambda}_I(x), \boldsymbol{\mu}_I(x)) \in \mathbb{R}^p \times \mathbb{R}_+^{\text{Card}(I)}$  be the corresponding Lagrange multipliers:

$$\begin{bmatrix} \boldsymbol{\lambda}_I(x) \\ \boldsymbol{\mu}_I(x) \end{bmatrix} := -(\text{DC}_I \text{DC}_I^T)^{-1} \text{DC}_I(x) \nabla J(x). \quad (3.3.6)$$

Last but not least, let us recall the necessary first-order optimality conditions (the KKT conditions) for equality and inequality constrained problems:

**Definition 3.4** (KKT conditions, [70, 244]). A point  $x^* \in V$  is said to satisfy the Karush, Kuhn and Tucker conditions for (3.1.2) if and only if there exist  $\boldsymbol{\lambda}(x^*) \in \mathbb{R}^p$  and  $\boldsymbol{\mu}(x^*) \in \mathbb{R}_+^q$  such that:

$$\begin{cases} \nabla J(x^*) + \text{Dg}(x^*)^T \boldsymbol{\lambda}(x^*) + \text{Dh}(x^*)^T \boldsymbol{\mu}(x^*) = 0, \\ \mathbf{g}(x^*) = 0, \quad \mathbf{h}(x^*) \leq 0, \\ \forall i = 1, \dots, q, \quad \mu_i h_i(x^*) = 0. \end{cases} \quad (3.3.7)$$

### 3.3.2 Definition of the range space step

**Definition 3.5** (range space step). The range step  $\boldsymbol{\xi}_C(x)$  associated with the optimization problem (3.1.2) is defined by

$$\boldsymbol{\xi}_C(x) := \text{DC}_{\tilde{I}(x)}^T (\text{DC}_{\tilde{I}(x)} \text{DC}_{\tilde{I}(x)}^T)^{-1} \mathbf{C}_{\tilde{I}(x)}(x), \quad (3.3.8)$$

where  $\tilde{I}(x)$  is the subset of saturated or violated constraints, defined by (3.3.3).

The purpose of the range space step  $\boldsymbol{\xi}_C(x)$  is to decrease the violation of the constraints as we shall see in proposition 3.5 below. The counterpart of lemma 3.2 holds exactly in this context, in particular:

1.  $\boldsymbol{\xi}_C(x)$  is orthogonal to  $\text{Ker}(\text{DC}_{\tilde{I}(x)})$ .
2.  $-\boldsymbol{\xi}_C(x)$  is a Gauss-Newton direction for the violation of the constraints:

$$\text{DC}_{\tilde{I}(x)}(-\boldsymbol{\xi}_C(x)) = -\mathbf{C}_{\tilde{I}(x)}(x).$$

### 3.3.3 Definition and properties of the null space step

The definition of the null space direction  $\boldsymbol{\xi}_J(x)$  is slightly more involved than in the equality constrained case since it is not obtained by replacing  $\text{Dg}(x)$  by  $\text{DC}_{\tilde{I}(x)}$  in (3.2.8). It requires the introduction of a different subset  $\hat{I}(x) \subset \tilde{I}(x)$ , which is now detailed.

The null space step  $\boldsymbol{\xi}_J(x)$  is sought, up to a change of sign, as a best normalized descent direction diminishing violated or saturated inequality constraints. Following the characterization lemma 3.1 for equality constrained problems,  $-\boldsymbol{\xi}_J(x)$  shall be set positively proportional to the solution  $\boldsymbol{\xi}^*(x)$  of the following minimization problem:

$$\begin{aligned} \boldsymbol{\xi}^*(x) &= \arg \min_{\boldsymbol{\xi} \in V} \text{DJ}(x) \boldsymbol{\xi} \\ \text{s.t.} \quad &\begin{cases} \text{Dg}(x) \boldsymbol{\xi} = 0 \\ \text{Dh}_{\tilde{I}(x)}(x) \boldsymbol{\xi} \leq 0 \\ \|\boldsymbol{\xi}\|_V \leq 1. \end{cases} \end{aligned} \quad (3.3.9)$$

The problem (3.3.9) could be solved directly with standard quadratic programming algorithms. However, it is convenient to characterize explicitly the minimizer  $\boldsymbol{\xi}^*(x)$  of (3.3.9) by examining the dual problem. This will allow us to obtain in definition 3.6 an explicit formula for the null space direction  $\boldsymbol{\xi}_J(x)$ , in the form of (3.1.5).

**Proposition 3.2.** *Let  $x \in V$  satisfy the qualification condition (3.3.4). There exists a unique couple of multipliers  $\boldsymbol{\lambda}^*(x) \in \mathbb{R}^p$  and  $\boldsymbol{\mu}^*(x) \in \mathbb{R}_+^{\tilde{q}(x)}$  solution to the following quadratic optimization problem which is the dual of (3.3.9):*

$$(\boldsymbol{\lambda}^*(x), \boldsymbol{\mu}^*(x)) := \arg \min_{\substack{\boldsymbol{\lambda} \in \mathbb{R}^p \\ \boldsymbol{\mu} \in \mathbb{R}_+^{\tilde{q}(x)}, \boldsymbol{\mu} \geq 0}} \|\nabla J(x) + \mathbf{Dg}(x)^\top \boldsymbol{\lambda} + \mathbf{Dh}_{\tilde{I}(x)}(x)^\top \boldsymbol{\mu}\|_V. \quad (3.3.10)$$

*Proof.* Problem (3.3.9) is equivalent to the following min-max formulation:

$$\min_{\substack{\boldsymbol{\xi} \in V \\ \langle \boldsymbol{\xi}, \boldsymbol{\xi} \rangle_V \leq 1}} \max_{\substack{\boldsymbol{\lambda} \in \mathbb{R}^p \\ \boldsymbol{\mu} \in \mathbb{R}_+^{\tilde{q}(x)}}} \mathbf{D}J(x)\boldsymbol{\xi} + \boldsymbol{\lambda}^\top \mathbf{Dg}(x)\boldsymbol{\xi} + \boldsymbol{\mu}^\top \mathbf{Dh}_{\tilde{I}(x)}(x)\boldsymbol{\xi}.$$

Exchanging formally the min and the max and performing the maximization with respect to  $\boldsymbol{\xi}$  as in the proof of lemma 3.1 yields that (3.3.10) is the dual problem of (3.3.9) up to a change of sign (the duality gap between (3.3.10) and (3.3.9) will be shown to vanish in proposition 3.3). The program (3.3.10) brings into play the closed convex set  $\mathbb{R}^p \times \mathbb{R}_+^{\tilde{q}(x)}$  and the least squares functional

$$(\boldsymbol{\lambda}, \boldsymbol{\mu}) \mapsto \left\| \nabla J(x) + \mathbf{D}\mathbf{C}_{\tilde{I}(x)}(x)^\top \begin{bmatrix} \boldsymbol{\lambda} \\ \boldsymbol{\mu} \end{bmatrix} \right\|_V.$$

The latter is strictly convex over  $\mathbb{R}^p \times \mathbb{R}_+^{\tilde{q}(x)}$  by virtue of (3.3.4). Hence, (3.3.10) has a unique solution.  $\square$

The optimization problem (3.3.10) belongs to the class of non negative least squares problems; it can be solved efficiently with a number of dedicated solvers, such as `cvxopt` [41] or `IPOPT` [310]. One nice feature of (3.3.10) lies in that its dimension is the number  $p + \tilde{q}(x)$  of saturated or violated constraints, which is small for many practical cases, e.g. in all the shape optimization applications considered in this thesis. It is also possible to exploit the ‘sparsity’ of the constraints when  $p + \tilde{q}(x)$  is large, see remark 3.8 below.

The next proposition relates the optimal values and the solutions  $\boldsymbol{\xi}^*(x)$  and  $(\boldsymbol{\lambda}^*(x), \boldsymbol{\mu}^*(x))$  of the primal and dual problems (3.3.9) and (3.3.10). In essence, we show that the optimal feasible descent direction  $\boldsymbol{\xi}^*(x)$  of (3.3.12) is the projection of the gradient  $\nabla J(x)$  onto the cone of feasible directions. The proof follows classical arguments of linear programming duality theory and it is detailed for the convenience of the reader.

**Proposition 3.3.** *Let  $x \in V$  satisfy the qualification condition (3.3.4) and denote*

$$m^*(x) := \|\nabla J(x) + \mathbf{Dg}(x)^\top \boldsymbol{\lambda}^*(x) + \mathbf{Dh}_{\tilde{I}(x)}(x)^\top \boldsymbol{\mu}^*(x)\|_V$$

*the value of the dual problem (3.3.10). Then the value of the primal problem (3.3.9) is  $p^*(x) = -m^*(x)$  and the following alternative holds:*

1.  $m^*(x) = 0$ : the first line of the KKT conditions (3.3.7) for the minimization problem (3.1.2) holds with (necessarily unique) Lagrange multipliers  $(\boldsymbol{\lambda}^*(x), \boldsymbol{\mu}^*(x)) \in \mathbb{R}^p \times \mathbb{R}_+^{\tilde{q}(x)}$ :

$$\nabla J(x) + \mathbf{Dg}(x)^\top \boldsymbol{\lambda}^*(x) + \mathbf{Dh}_{\tilde{I}(x)}(x)^\top \boldsymbol{\mu}^*(x) = 0. \quad (3.3.11)$$

*One particular minimizer of (3.3.9) is  $\boldsymbol{\xi}^*(x) = 0$ .*

2.  $m^*(x) > 0$ : (3.3.11) does not hold and there exists a unique minimizer  $\boldsymbol{\xi}^*(x)$  to (3.3.9), given by

$$\boldsymbol{\xi}^*(x) = -\frac{\nabla J(x) + \mathbf{Dg}(x)^\top \boldsymbol{\lambda}^*(x) + \mathbf{Dh}_{\tilde{I}(x)}(x)^\top \boldsymbol{\mu}^*(x)}{\|\nabla J(x) + \mathbf{Dg}(x)^\top \boldsymbol{\lambda}^*(x) + \mathbf{Dh}_{\tilde{I}(x)}(x)^\top \boldsymbol{\mu}^*(x)\|_V}. \quad (3.3.12)$$

*Proof.* Let  $\boldsymbol{\xi} \in V$  be a feasible direction for the problem (3.3.9), i.e.  $\mathbf{Dg}(x)\boldsymbol{\xi} = 0$ ,  $\mathbf{Dh}_{\tilde{I}(x)}(x)\boldsymbol{\xi} \leq 0$  and  $\|\boldsymbol{\xi}\|_V \leq 1$ . Then for any  $(\boldsymbol{\lambda}, \boldsymbol{\mu}) \in \mathbb{R}^p \times \mathbb{R}_+^{\tilde{q}(x)}$ , it holds

$$\begin{aligned} \mathbf{D}J(x)\boldsymbol{\xi} &\geq \mathbf{D}J(x)\boldsymbol{\xi} + \boldsymbol{\lambda}^\top \mathbf{Dg}(x)\boldsymbol{\xi} + \boldsymbol{\mu}^\top \mathbf{Dh}_{\tilde{I}(x)}(x)\boldsymbol{\xi} \\ &= \langle \nabla J(x) + \mathbf{Dg}(x)^\top \boldsymbol{\lambda} + \mathbf{Dh}_{\tilde{I}(x)}(x)^\top \boldsymbol{\mu}, \boldsymbol{\xi} \rangle_V \\ &\geq -\|\nabla J(x) + \mathbf{Dg}(x)^\top \boldsymbol{\lambda} + \mathbf{Dh}_{\tilde{I}(x)}(x)^\top \boldsymbol{\mu}\|_V \end{aligned} \quad (3.3.13)$$

Since (3.3.13) holds for any feasible direction  $\xi$  for (3.3.9), and for any  $(\lambda, \mu) \in \mathbb{R}^p \times \mathbb{R}_+^{\tilde{q}(x)}$ , it follows:

$$\min_{\substack{\xi \in V \\ \xi \text{ feasible for (3.3.9)}}} DJ(x)\xi \geq - \min_{\substack{\lambda \in \mathbb{R}^p \\ \mu \in \mathbb{R}_+^{\tilde{q}(x)}, \mu \geq 0}} \|\nabla J(x) + Dg(x)^T \lambda + Dh_{\tilde{I}(x)}(x)^T \mu\|_V. \quad (3.3.14)$$

Therefore, we have proven that  $p^*(x) \geq -m^*(x)$ . We now examine the alternative  $m^*(x) = 0$  or  $m^*(x) > 0$ :

1. If  $m^*(x) = 0$ , then (3.3.14) implies  $p^*(x) \geq 0$ . Therefore, the value of (3.3.9) is  $p^*(x) = -m^*(x) = 0$ , attained in particular at  $\xi^* = 0$ , and more generally at any feasible  $\xi^* \in V$  satisfying  $\mu^*(x)^T Dh_{\tilde{I}(x)}(x)\xi^* = 0$ , as follows readily from the KKT conditions for (3.3.9). Furthermore, the KKT equation (3.3.11) is satisfied by definition of  $m^*(x) = 0$ .
2. Assume now  $m^*(x) > 0$ . The KKT condition for (3.3.9) states that for any local optimum  $\xi'$ , there exists  $(\lambda', \mu') \in \mathbb{R}^p \times \mathbb{R}_+^{\tilde{q}(x)}$  and  $\alpha \geq 0$  such that,

$$\forall \xi \in V, (DJ(x) + \lambda'^T Dg(x) + \mu'^T Dh_{\tilde{I}(x)}(x))\xi = -\alpha \langle \xi', \xi \rangle_V. \quad (3.3.15)$$

Using Riesz identifications of the gradient and the differentials, we obtain

$$\alpha \xi' = -(\nabla J(x) + Dg(x)^T \lambda' + Dh_{\tilde{I}(x)}(x)^T \mu'),$$

and since  $m^*(x) > 0$ , it is necessary that  $\alpha > 0$ . The complementarity condition  $\alpha(\langle \xi', \xi' \rangle_V - 1) = 0$  yields then  $\|\xi'\|_V = 1$ , which readily implies:

$$\xi' = -\frac{\nabla J(x) + Dg(x)^T \lambda' + Dh_{\tilde{I}(x)}(x)^T \mu'}{\|\nabla J(x) + Dg(x)^T \lambda' + Dh_{\tilde{I}(x)}(x)^T \mu'\|_V}.$$

Then the complementarity condition for (3.3.9) implies  $\mu'^T Dh_{\tilde{I}(x)}(x)\xi' = 0$ . Therefore it holds that

$$\begin{aligned} DJ(x)\xi' &= DJ(x)\xi' + \lambda'^T Dg(x)\xi' + \mu'^T Dh_{\tilde{I}(x)}(x)\xi' \\ &= \langle \nabla J(x) + Dg(x)^T \lambda' + Dh_{\tilde{I}(x)}(x)^T \mu', \xi' \rangle_V \\ &= -\|\nabla J(x) + Dg(x)^T \lambda' + Dh_{\tilde{I}(x)}(x)^T \mu'\|_V. \end{aligned} \quad (3.3.16)$$

The previous equation together with the inequality (3.3.13) with  $\xi = \xi'$  then implies that  $(\lambda', \mu')$  achieves the minimum of (3.3.10). By uniqueness, this implies  $\lambda' = \lambda^*(x)$  and  $\mu' = \mu^*(x)$ , hence  $\xi' = \xi^*(x)$ . Furthermore,  $p^*(x) = DJ(x)\xi^*(x) = DJ(x)\xi' = -m^*(x)$ .

□

Finally, the next proposition characterizes explicitly the expression of the optimal descent direction  $\xi^*(x)$  from the signs of the multiplier  $\mu^*(x)$ , and highlights in which sense the problem (3.3.9) is combinatorial. Let us recall the definitions (3.3.5) and (3.3.6) for the projection operator  $\Pi_{C_I}$  and the multipliers  $(\lambda_I(x), \mu_I(x))$  which are used in the result.

**Proposition 3.4.** *In the context of point (2) in proposition 3.3, let  $\xi^*(x)$  and  $(\lambda^*(x), \mu^*(x))$  be the minimizers of the primal and dual problems (3.3.9) and (3.3.10). Define the subset  $\hat{I}(x) \subset \tilde{I}(x)$  by*

$$\hat{I}(x) := \{i \in \tilde{I}(x) \mid \mu_i^*(x) > 0\}. \quad (3.3.17)$$

1.  $(\lambda^*(x), \mu^*(x))$  and  $\xi^*(x)$  are explicitly given in terms of  $\hat{I}(x)$  by:

$$\begin{bmatrix} \lambda^*(x) \\ \hat{\mu}^*(x) \end{bmatrix} = \begin{bmatrix} \lambda_{\hat{I}(x)}(x) \\ \mu_{\hat{I}(x)}(x) \end{bmatrix} = -(DC_{\hat{I}(x)} DC_{\hat{I}(x)}^T)^{-1} DC_{\hat{I}(x)} \nabla J(x), \quad (3.3.18)$$

$$\xi^*(x) = -\frac{\Pi_{C_{\hat{I}(x)}}(\nabla J(x))}{\|\Pi_{C_{\hat{I}(x)}}(\nabla J(x))\|_V}, \quad (3.3.19)$$

where  $\hat{\mu}^*(x) := (\mu_i^*(x))_{i \in \hat{I}(x)}$  is the vector collecting all positive components of  $\mu^*(x)$ .



2.  $\widehat{I}(x)$  is equivalently the unique solution to either of the following discrete optimization problems:

$$\begin{aligned} \widehat{I}(x) = \arg \max_{I \subset \widetilde{I}(x)} \quad & \|\Pi_{C_I}(\nabla J(x))\|_V \\ \text{s.t.} \quad & D\mathbf{h}_{\widetilde{I}(x)}(x)\Pi_{C_I}(\nabla J(x)) \geq 0, \end{aligned} \quad (3.3.20)$$

$$\begin{aligned} \widehat{I}(x) = \arg \min_{I \subset \widetilde{I}(x)} \quad & \|\Pi_{C_I}(\nabla J(x))\|_V \\ \text{s.t.} \quad & \boldsymbol{\mu}_I(x) \geq 0. \end{aligned} \quad (3.3.21)$$

In particular,  $\widehat{I}(x)$  is the unique subset  $I \subset \widetilde{I}(x)$  satisfying simultaneously both feasibility conditions

$$D\mathbf{h}_{\widehat{I}(x)}(x)\Pi_{C_I}(\nabla J(x)) \geq 0 \text{ and } \boldsymbol{\mu}_I(x) \geq 0.$$

*Proof.* 1. The complementary condition for the primal and dual problems (3.3.9) and (3.3.10) reads

$$\forall i \in \widetilde{I}(x), \quad \mu_i^*(x) Dh_i(x) \boldsymbol{\xi}^*(x) = 0. \quad (3.3.22)$$

Therefore,  $Dh_i(x) \boldsymbol{\xi}^*(x) = 0$  for all indices  $i \in \widehat{I}(x)$ , which implies that  $DC_{\widehat{I}(x)}(x) \boldsymbol{\xi}^*(x) = 0$ . Then, after left multiplication of (3.3.12) by  $(DC_{\widehat{I}(x)} DC_{\widetilde{I}(x)}^T)^{-1} DC_{\widehat{I}(x)}$ , we obtain (3.3.18), whence (3.3.19) follows.

2. Let  $I \subset \widetilde{I}(x)$  a subset satisfying  $D\mathbf{h}_{\widetilde{I}(x)}(x)\Pi_{C_I}(\nabla J(x)) \geq 0$ . This implies that

$$\boldsymbol{\xi} = -\Pi_{C_I}(\nabla J(x)) / \|\Pi_{C_I}(\nabla J(x))\|_V$$

is feasible for the primal problem (3.3.9), and we obtain by definition of  $\boldsymbol{\xi}^*(x)$  that

$$-\|\Pi_{C_{\widehat{I}(x)}}(\nabla J(x))\|_V = DJ(x) \boldsymbol{\xi}^*(x) \leq DJ(x) \boldsymbol{\xi} = -\|\Pi_{C_I}(\nabla J(x))\|_V, \quad (3.3.23)$$

whence the maximization property (3.3.20).

For  $I \subset \widetilde{I}(x)$  satisfying  $\boldsymbol{\mu}_I(x) \geq 0$ , we obtain feasible multipliers  $(\boldsymbol{\lambda}, \boldsymbol{\mu})$  for the dual problem (3.3.10) by taking  $\boldsymbol{\mu}$  to be equal to  $\boldsymbol{\mu}_I$  on the indices of  $I$  and extended by 0 in the complementary subset  $\widetilde{I}(x) \setminus I$ . Then the optimality of  $(\boldsymbol{\lambda}^*(x), \boldsymbol{\mu}^*(x))$  for this dual problem reads:

$$\begin{aligned} \|\Pi_{C_{\widehat{I}(x)}}(\nabla J(x))\|_V &= \|\nabla J + D\mathbf{g}(x)^T \boldsymbol{\lambda}^*(x) + D\mathbf{h}_{\widetilde{I}(x)}(x)^T \boldsymbol{\mu}^*(x)\|_V \\ &\leq \|\nabla J(x) + D\mathbf{g}(x)^T \boldsymbol{\lambda} + D\mathbf{h}_{\widetilde{I}(x)}^T \boldsymbol{\mu}\|_V = \|\Pi_{C_I}(\nabla J(x))\|_V, \end{aligned} \quad (3.3.24)$$

whence the minimization property (3.3.21). □

**Remark 3.5.** In view of (3.3.17), the optimal multiplier  $\boldsymbol{\mu}^*(x)$  can be interpreted as an indicator variable specifying which constraints of  $\widetilde{I}(x)$  are ‘not aligned’ with the gradient  $\nabla J(x)$  and should be kept in the subset  $\widehat{I}(x)$ . The best descent direction (in the sense of (3.3.9)) is obtained by projecting the gradient  $\nabla J(x)$  onto the tangent space of the constraint subset  $\widehat{I}(x)$  rather than onto the full set of violated or saturated constraints  $\widetilde{I}(x)$ . Indeed, the descent direction  $\boldsymbol{\xi} = -\Pi_{C_{\widetilde{I}(x)}} \nabla J(x)$  that would be obtained by projecting  $\nabla J(x)$  on the whole set  $\widetilde{I}(x)$  would only keep them all constant at first order, i.e.  $Dh_i(x) \boldsymbol{\xi} = 0$ , (see remark 3.9 for more details). It is therefore more efficient to project  $\nabla J(x)$  only on those constraints associated to the indices  $i \in \widehat{I}(x)$ , thus allowing the remaining ones (associated to  $i \in \widetilde{I}(x) \setminus \widehat{I}(x)$ , indicating vanishing multipliers  $\mu_i^*(x) = 0$ ) to decrease since the calculated descent direction ensures that  $Dh_i(x) \boldsymbol{\xi}^*(x) \leq 0$  holds for all  $i = 1, \dots, q$ .

**Remark 3.6.** The use of a dual problem such as (3.3.10) in order to obtain information about which constraints should remain active is classical in active sets methods, see e.g. [77, 187, 244]. In principle, the optimal subset  $\widehat{I}(x)$  could be found by solving the discrete problems (3.3.20) or (3.3.21). However, we expect that in practice, it is more efficient to rely on iterative solvers relying on gradient descent for solving the dual problem (3.3.10), e.g. a cone programming solver or a non negative least squares algorithm such as [77]. This is what we do in the sequel.

Having introduced the subset  $\widehat{I}(x)$  (defined in (3.3.17)), we are now able to define the null space direction  $\xi_J(x)$  in the present context:  $-\xi_J(x)$  is set to be a positive multiple of the optimal descent direction  $\xi^*(x)$  supplied by (3.3.19).

**Definition 3.6.** For any point  $x \in V$  satisfying the constraint qualification (3.3.4), the null space direction  $\xi_J(x)$  at  $x$  for the optimization problem (3.1.2) is defined by:

$$\xi_J(x) := \Pi_{\mathcal{C}_{\widehat{I}(x)}}(\nabla J(x)) = (I - \text{DC}_{\widehat{I}(x)}(x)^\top (\text{DC}_{\widehat{I}(x)} \text{DC}_{\widehat{I}(x)}^\top)^{-1} \text{DC}_{\widehat{I}(x)}) \nabla J(x), \quad (3.3.25)$$

where  $\widehat{I}(x)$  is the optimal set defined by (3.3.17).

The main point in definition 3.6 is that, while all violated and saturated constraints are taken into account in the Gauss-Newton direction  $\xi_C(x)$  defined by (3.3.8), only those constraints in  $\widehat{I}(x)$ , not aligned with the gradient  $\nabla J(x)$ , occur in the definition of  $\xi_J(x)$ .

**Remark 3.7.** With our notation, the optimization scheme proposed by Barbarosie et. al. [59, 60] reads

$$\begin{cases} x_{n+1} = x_n - \Delta t \nabla J(x_n) - \text{DC}_{I(x_n)}^\top \nu_n \\ \nu_n = -\Delta t (\text{DC}_{I(x_n)} \text{DC}_{I(x_n)}^\top)^{-1} \text{DC}_{I(x_n)} \nabla J(x_n) + \text{DC}_{I(x_n)}^\top (\text{DC}_{I(x_n)} \text{DC}_{I(x_n)}^\top)^{-1} \mathcal{C}_{I(x_n)}, \end{cases} \quad (3.3.26)$$

where the set  $I(x_n)$  is obtained by removing indices from  $\widehat{I}(x_n)$  one by one, starting from the index  $i_0$  associated with the most negative multiplier  $\nu_{n,i_0} < 0$ , until all of them become non negative. Therefore, the set  $I(x_n)$  used in this strategy and that  $\widehat{I}(x_n)$  featured in our strategy, given by (3.3.17), do not coincide in general; one could think of configurations where the procedure of [60] would fail to find the optimal set  $\widehat{I}(x_n)$  (for example if  $i_0 \in \widehat{I}(x_n)$ ) and would project the gradient on a less optimal subset of constraints. We note that no convergence result is given by the authors about this procedure.

**Remark 3.8.** Let us discuss two extreme cases related to the involved computational effort in the numerical implementation of (3.3.25). Upon discretization, we may assume that  $V = \mathbb{R}^k$  is a finite-dimensional space.

1. If the total number  $p + \tilde{q}$  of saturated or violated constraints is small compared to the dimension  $k$  of  $V$ , it is best, for numerical efficiency, to assemble the small square matrix  $(\text{DC}_{\widehat{I}(x)} \text{DC}_{\widehat{I}(x)}^\top)$  and to invert it by a direct method.
2. If  $V = \mathbb{R}^k$  is equipped with an inner product encoded by a matrix  $A$ , and if  $p + \tilde{q}$  is of the order of  $k$  or larger, the computation of the inverse of  $(\text{DC}_{\widehat{I}(x)} \text{DC}_{\widehat{I}(x)}^\top)$  can be expensive. However, it can be still tractable if both  $\text{DC}$  and  $A$  are sparse matrices. For instance, this occurs in the case of bound constraints on the optimization variable  $x = (x_1, \dots, x_k)$ , e.g. constraints of the form  $\alpha_i \leq x_i \leq \beta_i$ ,  $i = 1, \dots, k$ . Recalling from remark 3.1 that in this setting,  $\text{DC}_{\widehat{I}(x)}^\top = A^{-1} \text{DC}_{\widehat{I}(x)}^\top$ , it can be verified that the vector

$$X := A^{-1} \text{DC}_{\widehat{I}(x)}^\top (\text{DC}_{\widehat{I}(x)} A^{-1} \text{DC}_{\widehat{I}(x)}^\top)^{-1} \text{DC}_{\widehat{I}(x)} \nabla J(x)$$

can be computed as the solution to the sparse linear system

$$\begin{bmatrix} A & -\text{DC}_{\widehat{I}(x)}^\top \\ \text{DC}_{\widehat{I}(x)} & 0 \end{bmatrix} \begin{bmatrix} X \\ \Lambda \end{bmatrix} = \begin{bmatrix} 0 \\ \text{DC}_{\widehat{I}(x)} \nabla J(x) \end{bmatrix},$$

where  $\Lambda \in \mathbb{R}^{p + \text{Card}(\widehat{I}(x))}$  is an extra slack variable, which yields the null space directions  $\xi_J(x) = \nabla J(x) - X$ . A similar strategy based on the sparsity of  $A$  and  $\text{DC}_{\widehat{I}(x)}^\top$  can be used to compute the range space direction  $\xi_C(x)$  of (3.3.8), or to solve the dual quadratic subproblem (3.3.10).

**Remark 3.9.** As we have already mentioned, the Lagrange multiplier  $\mu^*(x)$  given by (3.3.18) may be understood as an indicator of which inequality constraints are aligned with the gradient of  $J$ . To further highlight this, it is instructive to consider the particular situation where the gradients of the constraint functions are orthogonal, i.e.:

$$\begin{aligned} \langle \nabla g_i(x), \nabla g_j(x) \rangle_V &= 0, \text{ for } i, j = 1, \dots, p, \quad i \neq j, \\ \langle \nabla h_i(x), \nabla h_j(x) \rangle_V &= 0, \text{ for } i, j = 1, \dots, q, \quad i \neq j, \\ \langle \nabla g_i(x), \nabla h_j(x) \rangle_V &= 0, \text{ for } i = 1, \dots, p, \quad j = 1, \dots, q. \end{aligned}$$

In this case, it easily follows from the Pythagore theorem that for any  $(\boldsymbol{\lambda}, \boldsymbol{\mu}) \in \mathbb{R}^p \times \mathbb{R}_+^{\tilde{q}(x)}$ ,

$$\begin{aligned} \|\nabla J(x) + \mathbf{D}\mathbf{g}(x)^T \boldsymbol{\lambda} + \mathbf{D}\mathbf{h}_{\tilde{I}(x)}(x)^T \boldsymbol{\mu}\|_V^2 &= \left\| \nabla J(x) + \sum_{i=1}^p \lambda_i \nabla g_i(x) + \sum_{j \in \tilde{I}(x)} \mu_j \nabla h_j(x) \right\|_V^2 \\ &= \|\nabla J(x)\|_V^2 + \sum_{i=1}^p (\lambda_i^2 \|\nabla g_i(x)\|_V^2 + 2\lambda_i \langle \nabla J(x), \nabla g_i(x) \rangle_V) \\ &\quad + \sum_{j \in \tilde{I}(x)} (\mu_j^2 \|\nabla h_j(x)\|_V^2 + 2\mu_j \langle \nabla J(x), \nabla h_j(x) \rangle_V). \end{aligned}$$

Therefore the minimization problem (3.3.10) is separable with respect to the variables  $(\boldsymbol{\lambda}, \boldsymbol{\mu}) \in \mathbb{R}^p \times \mathbb{R}_+^{\tilde{q}(x)}$ :  $(\lambda_i^*(x))_{1 \leq i \leq p}$  and  $(\mu_i^*(x))_{i \in \tilde{I}(x)}$  are the respective solutions to the minimization problems:

$$\forall i \in 1 \dots p, \quad \lambda_i^*(x) = \arg \min_{t \in \mathbb{R}} (t^2 \|\nabla g_i(x)\|_V^2 + 2t \langle \nabla J(x), \nabla g_i(x) \rangle_V),$$

$$\forall i \in \tilde{I}(x), \quad \mu_i^*(x) = \arg \min_{\substack{t \in \mathbb{R} \\ t \geq 0}} (t^2 \|\nabla h_i(x)\|_V^2 + 2t \langle \nabla J(x), \nabla h_i(x) \rangle_V),$$

which yields eventually:

$$\lambda_i^*(x) = -\frac{\langle \nabla J(x), \nabla g_i(x) \rangle_V}{\|\nabla g_i(x)\|_V^2}, \quad \mu_i^*(x) = \begin{cases} 0 & \text{if } \langle \nabla J(x), \nabla h_i(x) \rangle_V \geq 0, \\ -\frac{\langle \nabla J(x), \nabla h_i(x) \rangle_V}{\|\nabla h_i(x)\|_V^2} & \text{otherwise.} \end{cases}$$

Hence,  $\mu_i^*(x)$  is positive if and only if the direction  $-\nabla J(x)$  leads to an increase (i.e. a violation) of the  $i^{\text{th}}$  inequality constraint.

In the general case where all the constraint gradients are not mutually orthogonal, the interpretation of  $\boldsymbol{\mu}^*(x)$  is similar, up to the additional complication that (3.3.10) accounts for the combinatorics behind the possible alignments between different constraint gradients. In the following, with a slight abuse of language, we shall nevertheless refer to the indices  $i \in \tilde{I}(x) \setminus \hat{I}(x)$  as those associated to constraints which are ‘aligned’ with  $\nabla J(x)$ .

### 3.3.4 Decrease properties of the trajectories of the null space ODE

The final result of this section is the counterpart of proposition 3.1 in the case of the equality and inequality constrained optimization problem (3.1.2).

**Proposition 3.5.** *Assume that the trajectories  $x(t)$  of the flow*

$$\begin{cases} \dot{x}(t) = -\alpha_J \boldsymbol{\xi}_J(x(t)) - \alpha_C \boldsymbol{\xi}_C(x(t)) \\ x(0) = x_0, \end{cases} \quad (3.3.27)$$

with  $\boldsymbol{\xi}_J$  and  $\boldsymbol{\xi}_C$  given by (3.3.8) and (3.3.25) exist on some interval  $[0, T]$  for  $T > 0$  and are such that:

(a) the set  $\tilde{I}(x(t))$  defined in (3.3.3) is constant over  $[0, T]$ :

$$\forall t \in [0, T], \quad \tilde{I}(x(t)) = \tilde{I}(x_0);$$

(b) the constraints remain qualified along the flow  $x(t)$ , in the sense of (3.3.4).

Then the following properties hold true:

1. The violation of the constraints decreases exponentially:

$$\forall t \in [0, T], \quad \mathbf{g}(x(t)) = e^{-\alpha_C t} \mathbf{g}(x_0) \quad \text{and} \quad \mathbf{h}_{\tilde{I}(x(t))}(x(t)) \leq e^{-\alpha_C t} \mathbf{h}_{\tilde{I}(x_0)}(x_0). \quad (3.3.28)$$

2.  $J(x(t))$  decreases ‘as soon as the violation (3.3.28) of the constraints is sufficiently small’ in the following sense. Assume that  $\text{rank}(\mathbf{D}\mathbf{C}_{\tilde{I}(x_0)}(x))$  is maximal for all  $x$  in  $K = \{x \in V \mid \|\mathbf{C}_{\tilde{I}(x_0)}(x)\|_\infty \leq \|\mathbf{C}_{\tilde{I}(x_0)}(x_0)\|_\infty\}$  and

$$\sup_{x \in K} \|\nabla J(x)\|_V |\sigma_p^{-1}(x)| < +\infty. \quad (3.3.29)$$

where  $\sigma_p(x)$  is the smallest singular value of  $\mathrm{DC}_{\tilde{I}(x)}(x)$ . Then there exists a constant  $C > 0$  such that

$$\forall t \in [0, T], \|\Pi_{C_{\tilde{I}(x(t))}}(\nabla J(x(t)))\|_{\tilde{V}}^2 > Ce^{-\alpha_C t} \Rightarrow \frac{d}{dt} J(x(t)) < 0. \quad (3.3.30)$$

3. Any stationary point  $x^*$  of the flow (3.3.27) satisfies the KKT optimality conditions (3.3.7) which equivalently rewrite:

$$\begin{cases} \nabla J(x^*) + \mathrm{Dg}(x^*)^\top \boldsymbol{\lambda}^*(x^*) + \mathrm{Dh}_{\tilde{I}(x^*)}(x^*)^\top \boldsymbol{\mu}^*(x^*) = 0, \\ (\mathbf{g}(x^*) = 0 \text{ and } \mathbf{h}_{\tilde{I}(x^*)}(x^*) = 0) \Leftrightarrow \mathbf{C}_{\tilde{I}(x^*)}(x^*) = 0, \end{cases} \quad (3.3.31)$$

where  $(\boldsymbol{\lambda}^*(x^*), \boldsymbol{\mu}^*(x^*)) \in \mathbb{R}^p \times \mathbb{R}_+^{\tilde{q}(x^*)}$  are defined in (3.3.10) or (3.3.18).

*Proof.* 1. The definition (3.3.8) of  $\boldsymbol{\xi}_C(x(t))$  implies  $\mathrm{DC}_{\tilde{I}(x(t))}\boldsymbol{\xi}_C(x(t)) = \mathbf{C}_{\tilde{I}(x(t))}(x(t))$ , and since  $-\boldsymbol{\xi}_J(x(t))$  is positively proportional to  $\boldsymbol{\xi}^*(x(t))$  (proposition 3.3), it holds

$$\mathrm{DC}_{\tilde{I}(x(t))}\boldsymbol{\xi}_J(x(t)) = 0, \quad -\mathrm{Dh}_{\tilde{I}(x(t)) \setminus \hat{I}(x(t))}(x(t))\boldsymbol{\xi}_J(x(t)) \leq 0.$$

Therefore we obtain

$$\frac{d}{dt} \mathbf{C}_{\tilde{I}(x(t))}(x(t)) = -\alpha_C \mathbf{C}_{\tilde{I}(x(t))}(x(t)) \text{ and } \frac{d}{dt} \mathbf{h}_{\tilde{I}(x(t)) \setminus \hat{I}(x(t))}(x(t)) \leq -\alpha_C \mathbf{h}_{\tilde{I}(x(t)) \setminus \hat{I}(x(t))}(x(t)) \quad (3.3.32)$$

from which (3.3.28) follows by application of Gronwall's lemma.

2. The proof is identical to that of proposition 3.1.

3. A stationary point  $x^*$  of (3.3.27) satisfies by definition  $-\alpha_J \boldsymbol{\xi}_J(x^*) - \alpha_C \boldsymbol{\xi}_C(x^*) = 0$ . Left multiplication of this identity by  $\mathrm{DC}_{\tilde{I}(x^*)}(x^*)$  yields:

$$-\alpha_J \mathrm{DC}_{\tilde{I}(x^*)}(x^*)\boldsymbol{\xi}_J(x^*) - \alpha_C \mathbf{C}_{\tilde{I}(x^*)}(x^*) = 0. \quad (3.3.33)$$

Remembering now that from definition (3.3.9),

$$-\mathrm{DC}_{\tilde{I}(x^*)}\boldsymbol{\xi}_J(x^*) \leq 0 \text{ and } \mathbf{C}_{\tilde{I}(x^*)}(x^*) \geq 0,$$

equality in (3.3.33) can hold only if both terms vanish. In particular, we infer that  $\mathbf{C}_{\tilde{I}(x^*)}(x^*) = 0$ , a fact which implies  $\boldsymbol{\xi}_C(x^*) = 0$  and which encompasses the last two lines of the KKT conditions (3.3.7). Returning to the fact that  $-\alpha_J \boldsymbol{\xi}_J(x^*) - \alpha_C \boldsymbol{\xi}_C(x^*) = 0$ , we obtain that  $\boldsymbol{\xi}_J(x^*) = 0$ , which corresponds to the first line in (3.3.7). This completes the proof.  $\square$

**Remark 3.10.** Since the sets  $\tilde{I}(x)$  or  $\hat{I}(x)$  are subject to change as soon as inequality constraints become active or inactive, or if not enough regularity holds, the ODE (3.3.27) has in general a discontinuous right-hand side and is defined only formally (note that a rigorous mathematical meaning could still be provided with the theory of non smooth ODEs, see [126, 158]). However and as discussed further on in the next remarks, its discretization makes sense and exhibits the same decrease properties as its continuous counterpart for sufficiently small steps  $\Delta t$ .

**Remark 3.11.** The assumption (a) in proposition 3.5, whereby the index set  $\tilde{I}(x(t))$  remains constant is essentially made to ensure that the right-hand side of the flow (3.3.27) is continuous. Indeed, in such a case, the range space direction  $\boldsymbol{\xi}_C(x(t))$  is continuous by its definition (3.3.8), while the null space step  $\boldsymbol{\xi}_J(x(t))$  is continuous because

$$\boldsymbol{\xi}_J(x(t)) = \nabla J(x(t)) + \mathrm{DC}_{\tilde{I}(x(t))} \begin{bmatrix} \boldsymbol{\lambda}^*(x(t)) \\ \boldsymbol{\mu}^*(x(t)) \end{bmatrix}$$

and it can be shown that the multipliers  $(\boldsymbol{\lambda}^*(x(t)), \boldsymbol{\mu}^*(x(t)))$  defined by (3.3.10) are continuous functions. At a time  $T$  corresponding to a sudden change of the index set  $\tilde{I}(x(t))$ , we assume that the solution  $x(t)$  can be extended by restarting the ODE (3.3.27) with the new index set  $\tilde{I}(x(T))$ . From (1) in proposition 3.5, the bound  $\mathbf{h}(x(t)) \leq e^{-\alpha_C t} \mathbf{h}(x(0))$  still holds after the time  $T$  for all constraints  $i \in \{1, \dots, q\}$ :

constraints are asymptotically satisfied. Properties (2) and (3) remain true, up to an adjustment of the constant  $C$  in (3.3.30) (which can be taken global since there are finitely many possible sets  $\tilde{I}(x(t))$ ). There may exist situations where the set of asymptotically saturated constraints  $\tilde{I}(x(t))$  could oscillate indefinitely. However (2) states that  $x(t)$  always keeps improving (in the sense of (3.3.30)), and (3) states that if  $x(t)$  eventually converges, it is necessarily towards a KKT point.

**Remark 3.12.** In practice, the analysis of proposition 3.5 is sufficient because, similarly to the conclusions of remark 3.3, analogous properties hold for the discrete scheme

$$x_{n+1} = x_n - \Delta t(\alpha_J \xi_J(x_n) + \alpha_C \xi_C(x_n)). \quad (3.3.34)$$

Indeed, one can easily check that:

1. Up to first order, the violation of the constraints decreases at a geometric rate:

$$C(x_{n+1}) = (1 - \alpha_C \Delta t)C(x_n) + o(\Delta t). \quad (3.3.35)$$

This suggests that in order to obtain a stable scheme, one must a priori select  $\alpha_C$  and  $\Delta t$  such that  $0 < \alpha_C \Delta t < 2$ .

2. An accumulation point  $x^*$  of the sequence  $(x_n)_{n \in \mathbb{N}}$  is feasible, i.e.  $C_{\tilde{I}(x^*)}(x^*) = 0$  and is a KKT point.

Finally, note that one flexibility of this ODE approach is that at the continuous level, the results of proposition 3.5 do not depend on the values of the parameters  $\alpha_J > 0$  and  $\alpha_C > 0$ . Therefore the convergence of the discrete scheme towards the continuous trajectory should hold as soon as the discretization step size  $\Delta t > 0$  is sufficiently small.

### 3.4 NUMERICAL DISCRETIZATION AND TIME-STEPPING SCHEMES FOR THE NULL SPACE ODE

This short section describes practical implementation details for the discretization of the ODE (3.1.3) by an explicit Euler scheme. Two important issues are discussed respectively in sections 3.4.1 and 3.4.2. First, we propose small adaptations in the computation of  $\xi_J(x)$  and  $\xi_C(x)$  in order to account for the discontinuous changes of the right-hand side  $-(\alpha_J \xi_J + \alpha_C \xi_C)$ . Then, a merit function is proposed for adapting the time step  $\Delta t$ . The complete implementation of the algorithm is summarized in section 3.4.3 below.

#### 3.4.1 Accounting for discontinuities near the inequality constraint barriers

A potential issue when implementing directly the Euler time-stepping scheme (3.3.2) comes from the fact that the vector fields  $\xi_J$  and  $\xi_C$  given by (3.3.8) and (3.3.25) are characterized by the same discontinuities as the discrete index mapping  $x \mapsto \tilde{I}(x)$ . As a result, abrupt oscillations of the discrete optimization path  $(x_n)$  may occur near the boundary of the feasible set: if  $h_i(x_n) = 0$  and  $i \in \tilde{I}(x_n)$  for some index  $i \in \{1, \dots, q\}$ , then in the definition (3.3.25) of  $\xi_J(x_n)$ , the gradient  $\nabla J(x_n)$  is projected tangentially to the constraint  $h_i$ , but it is not projected after any slight deviation (e.g. due to the discretization) making this constraint inactive ( $h_i(x_{n+1}) < 0$ ). This kind of issue is very classical in the discretization of ODEs with discontinuous vector fields and can be tackled by various methods, see e.g. [126] for a review.

In this section, we suggest a simple alternative that works well in practice to stabilize trajectories near these boundaries: inequality constraints are felt from a short distance by replacing the set  $\tilde{I}(x_n)$  in (3.1.7) with the set  $\tilde{I}_\varepsilon(x_n)$  of inequality constraints violated “up to  $\varepsilon_i$ ”:

$$\tilde{I}_\varepsilon(x_n) = \{i \in \{1, \dots, q\} \mid h_i(x_n) \geq -\varepsilon_i\}. \quad (3.4.1)$$

The tolerances  $\varepsilon_i > 0$  can be estimated in an automatic fashion, independent of an arbitrary rescaling of the constraints, thanks to an a posteriori bound we now detail. Let  $\mathbf{h}$  be a user-defined parameter accounting for the distance from the optimization path at which the constraints should be felt.

Assume now that the current point  $x_n$  satisfies the constraint  $h_i$  up to the uncertainty  $\mathbf{h}$  on its location: by this we mean that there exists some unknown point  $x_n^*$  such that  $\|x_n^* - x_n\| \leq \mathbf{h}$ ,  $h_i(x_n^*) > 0$  and  $h_i(x_n^*) = 0$ . Then the error  $\mathbf{h}$  for the location of  $x_n$  propagates to the constraint values  $h_i(x_n)$  according to the following inequality:

$$h_i(x_n) = |h_i(x_n) - h_i(x_n^*)| \simeq |Dh_i(x_n)(x_n^* - x_n)| \leq \|\nabla h_i(x_n)\| \mathbf{h}. \quad (3.4.2)$$

It is therefore natural to set

$$\varepsilon_i := \|\nabla h_i(x_n)\|_V \mathbf{h} \quad (3.4.3)$$

for the value of  $\varepsilon_i$  in (3.4.1). In our implementation relying on the discretization of (3.1.3), the parameter  $\mathbf{h}$  is set proportional to the time step  $\Delta t$ :  $\mathbf{h} = \mathbf{K}\Delta t$  for a constant  $\mathbf{K}$  to be defined by the user. The parameter  $\mathbf{h}$  should be defined in accordance with the typical distance  $\|\Delta x\|_V = \|x_{n+1} - x_n\|_V$  between two successive iterations; in the academic examples in section 3.5 below considering optimization in  $\mathbb{R}^k$ , we may set e.g.  $\mathbf{h} = 0.01$  for a typical increment size  $\|\Delta x\|_V \simeq 0.1$ . For our shape optimization applications in section 3.6,  $\mathbf{h}$  is typically of the order of the discretization mesh size, see section 3.6.2 below.

Note that more generally, the a posteriori bound (3.4.2) allows to assert whether a constraint  $C_i(x_n)$  can be considered as satisfied or not with respect to the numerical discretization.

The dual problem (3.3.10) is then solved with  $\tilde{I}_\varepsilon(x_n)$  instead of  $\tilde{I}(x_n)$  in order to obtain a new subset of indices  $\hat{I}_\varepsilon(x_n)$  which indicates which constraints are likely to be not aligned with the gradient  $\nabla J(x_n)$  when crossing the barrier  $\{\mathbf{h} = 0\}$ . The null space and range space steps  $\xi_J(x_n)$  and  $\xi_C(x_n)$  in step 4 of algorithm 3.1 are finally replaced with  $\xi_{J,\varepsilon}(x_n)$  and  $\xi_{C,\varepsilon}(x_n)$  computed as follows:

$$\xi_{J,\varepsilon}(x_n) := (I - \text{DC}_{\hat{I}_\varepsilon(x_n)}^T)(\text{DC}_{\hat{I}_\varepsilon(x_n)} \text{DC}_{\hat{I}_\varepsilon(x_n)}^T)^{-1} \text{DC}_{\hat{I}_\varepsilon(x_n)} \nabla J(x_n), \quad (3.4.4)$$

$$\xi_{C,\varepsilon}(x_n) := \text{DC}_{I_\varepsilon^*(x_n)}^T (\text{DC}_{I_\varepsilon^*(x_n)} \text{DC}_{I_\varepsilon^*(x_n)}^T)^{-1} C_{I_\varepsilon^*(x_n)}(x_n), \quad (3.4.5)$$

where  $I_\varepsilon^*(x_n) = \tilde{I}(x_n) \cup \hat{I}_\varepsilon(x_n)$  is the set of constraints that are either violated, saturated or not aligned with the gradient  $\nabla J(x_n)$  at  $\mathbf{h} = -(\varepsilon_1, \dots, \varepsilon_q)^T$ . The use of  $\hat{I}_\varepsilon(x_n)$  in the definition of  $\xi_{J,\varepsilon}(x_n)$  ensures that the gradient  $\nabla J(x_n)$  is being projected tangentially to the constraint on a small layer near the boundary of the feasible set. As a result, no abrupt discontinuity occurs anymore for  $\xi_{J,\varepsilon}$  and  $\xi_{C,\varepsilon}$  when crossing the boundary of the feasible domain while remaining in this layer. Including constraints  $i \in \hat{I}_\varepsilon(x_n)$  in the Gauss-Newton direction  $\xi_{C,\varepsilon}(x_n)$  even if they are satisfied (i.e. if  $-\varepsilon_i \leq h_i(x_n) \leq 0$ ) further allows to stabilize the values of these constraints closer to zero.

### 3.4.2 Time step adaptation based on a merit function.

The ODE (3.1.3) is discretized by an explicit scheme of the form:

$$x_{n+1} = x_n - \Delta t_n (\alpha_J \xi_J(x_n) + \alpha_C \xi_C(x_n)), \quad (3.4.6)$$

with a variable time step  $\Delta t_n > 0$ . The practical implementation of such a strategy is often guided by a merit function, i.e. an indicator allowing to detect that a step has been too large, a situation where a choice has to be made regarding whether the step should be reduced or accepted. For our null space algorithm, a merit function which resembles very much that of the Augmented Lagrangian Method is readily available, however with a specific choice of multipliers:

**Lemma 3.3.** *For a given  $x_n \in V$ , let  $\text{merit}_{x_n} : V \rightarrow \mathbb{R}$  be the function defined by*

$$\text{merit}_{x_n}(x) := \alpha_J \left( J(x) + \mathbf{\Lambda}(x_n)^T C_{\tilde{I}(x_n)}(x) \right) + \frac{\alpha_C}{2} C_{\tilde{I}(x_n)}(x)^T \mathbf{S}(x_n) C_{\tilde{I}(x_n)}(x) \quad (3.4.7)$$

where  $\mathbf{\Lambda}(x_n) = \left[ \boldsymbol{\lambda}^*(x_n)^T \quad \boldsymbol{\mu}^*(x_n)^T \right]^T$  is the vector of multipliers defined as the solution to the dual problem (3.3.10) (see (3.3.18)) and  $\mathbf{S}(x_n) = (\text{DC}_{\tilde{I}(x_n)}(x_n) \text{DC}_{\tilde{I}(x_n)}(x_n)^T)^{-1}$  is symmetric positive definite. Then (3.4.6) is a gradient step for decreasing the function  $\text{merit}_{x_n}$ , namely:

$$\nabla \text{merit}_{x_n}(x_n) = \alpha_J \xi_J(x_n) + \alpha_C \xi_C(x_n).$$

*Proof.* It is a straightforward computation of the gradient of (3.4.7).  $\square$

A possible implementation of an optimization strategy of the form (3.4.6) based on this merit function is summarized in algorithm 3.1, which requires the introduction of a few extra parameters:

- **time\_step**: choose a fixed time step  $\Delta t > 0$ .
- **maxtrials**: the optimization time step is decreased up to **maxtrials** times until the value of the merit function has decreased. If the merit function has not decreased after all **maxtrials** steps, the smallest step is accepted.

- **tolLag**: a small threshold for the values of the Lagrange multipliers  $\mu_i^*$  under which these are considered to be 0 (we took **tolLag**= $1e-8$ ). This value is used for the identification of the set  $\widehat{I}(x_n)$  according to the definition (3.3.17); it should be set in accordance with the machine precision and that of the quadratic programming solver for the dual problem (3.3.10).

Let us emphasize that these parameters have a quite intuitive and physical meaning, so that the task of assigning their values does not involve fine tunings in practice.

Importantly, the rescaling induced by the inverse of the correlation matrix  $(DC_{\widehat{I}(x_n)}DC_{\widehat{I}(x_n)}^T)^{-1}$  normalizes all the constraints; in particular, the whole algorithm 3.1 is invariant under multiplication of the constraints by arbitrary positive constants (up to the machine precision for the step 3); a preliminary rescaling of the constraints is therefore not required from the user.

### 3.4.3 Overall algorithm pseudo code

The resulting algorithmic implementation of the null space gradient flow taking into account both adaptations of section 3.4.1 and (3.4.2) is summarized in algorithm 3.1 below.

## 3.5 COMPARISONS WITH OTHER METHODS AND ILLUSTRATIONS ON ACADEMIC TEST CASES

This section is purely pedagogical and serves to illustrate and compare our null space gradient flow to other classical first order methods for constrained optimization problems. The method of slack variables for treating inequality constraints with the equality constrained gradient flow (3.1.11) is reviewed and compared with our method in section 3.5.1. The comparison of our gradient flow algorithm with the more classical SLP and Augmented Lagrangian methods are presented in section 3.5.2.

### 3.5.1 Comparison with the method of slack variables for inequality constraints

It is classical to introduce slack variables so as to turn inequality constraints in the problem (3.1.2) into equality constraints of an augmented problem including these additional variables. This section briefly reviews the method investigated by [277] in the context of dynamical system approaches to constrained optimization and compares it with our method based on the dual problem (3.3.10).

The method of slack variables consists in replacing the problem (3.1.2) with the following equivalent one, involving  $q$  extra variables  $(z_1, \dots, z_q) \in \mathbb{R}^q$ :

$$\begin{aligned} \min_{\substack{x \in V \\ z \in \mathbb{R}^q}} J(x) \\ \text{s.t. } \mathbf{C}(x, z) = 0, \end{aligned} \quad (3.5.1)$$

where the augmented vector of constraints  $\mathbf{C}(x, z)$  reads:

$$\mathbf{C}(x, z) := \begin{bmatrix} \mathbf{g}(x) \\ h_1(x) + \frac{1}{2}z_1^2 \\ \vdots \\ h_q(x) + \frac{1}{2}z_q^2 \end{bmatrix} \in \mathbb{R}^{p+q}.$$

Problem (3.5.1) is an equality constrained optimization problem of the form (3.2.1), set over the Hilbert space

$$\widetilde{V} := V \times \mathbb{R}^q, \text{ with inner product } \langle (x, z), (x', z') \rangle_{\widetilde{V}} := \langle x, x' \rangle_V + z^T z'.$$

It can be solved thanks to the proposed algorithm in section 3.2 for equality-constrained problems. The

---

**Algorithm 3.1** Discretization of the null space gradient flow (3.3.27), based on a merit function.

---

**for**  $n = 1 \dots \text{maxiter}$  **do**

1. Compute the gradients  $\nabla J(x_n)$ ,  $\nabla g_i(x_n)$  and  $\nabla h_j(x_n)$  for  $1 \leq i \leq p$ ,  $1 \leq j \leq q$  by solving, if necessary, the identification problem (3.2.3) and (3.2.4).
2. For all inequality constraints  $1 \leq i \leq q$ , compute the tolerance

$$\varepsilon_i := \|\nabla h_i(x_n)\|_V \mathbf{h}.$$

3. Determine the set  $\tilde{I}(x_n)$  of active or violated constraints and the set  $\tilde{I}_\varepsilon(x_n)$  of constraints violated “up to  $\varepsilon_i$ ”:

$$\begin{aligned}\tilde{I}(x_n) &= \{i \in \{1, \dots, q\} \mid h_i(x_n) \geq 0\} \\ \tilde{I}_\varepsilon(x_n) &= \{i \in \{1, \dots, q\} \mid h_i(x_n) \geq -\varepsilon_i\}.\end{aligned}$$

4. Denote by  $\tilde{q}_\varepsilon := \text{Card}(\tilde{I}_\varepsilon)$ . Solve the dual problem

$$(\boldsymbol{\lambda}_\varepsilon^*(x_n), \boldsymbol{\mu}_\varepsilon^*(x_n)) := \arg \min_{\substack{\boldsymbol{\lambda} \in \mathbb{R}^p \\ \boldsymbol{\mu} \in \mathbb{R}^{\tilde{q}_\varepsilon(x)}, \boldsymbol{\mu} \geq 0}} \|\nabla J(x) + \text{Dg}(x)^\top \boldsymbol{\lambda} + \text{Dh}_{\tilde{I}_\varepsilon(x)}(x)^\top \boldsymbol{\mu}\|_V$$

to obtain the optimal Lagrange multiplier  $\boldsymbol{\mu}^*(x_n)$ . Infer the subset  $\hat{I}_\varepsilon(x_n) \subset \tilde{I}_\varepsilon(x_n)$  indicating which constraints must remain active (proposition 3.4) :

$$\hat{I}_\varepsilon(x_n) = \{i \in \tilde{I}_\varepsilon(x_n) \mid \mu_{\varepsilon,i}^*(x_n) > \text{tolLag}\}. \quad (3.4.8)$$

5. Let  $I_\varepsilon^*(x_n) := \tilde{I}(x_n) \cup \hat{I}_\varepsilon(x_n)$ . Extract the vectors  $\mathbf{C}_{\hat{I}_\varepsilon(x_n)}(x_n)$  and  $\mathbf{C}_{I_\varepsilon^*(x_n)}(x_n)$  (defined by (3.1.4)) and compute

$$\begin{aligned}\boldsymbol{\xi}_J(x_n) &= (I - \text{DC}_{\hat{I}_\varepsilon(x_n)}^\top (\text{DC}_{\hat{I}_\varepsilon(x_n)} \text{DC}_{\hat{I}_\varepsilon(x_n)}^\top)^{-1} \text{DC}_{\hat{I}_\varepsilon(x_n)}) \nabla J(x_n), \\ \boldsymbol{\xi}_C(x_n) &= \text{DC}_{I_\varepsilon^*(x_n)}^\top (\text{DC}_{I_\varepsilon^*(x_n)} \text{DC}_{I_\varepsilon^*(x_n)}^\top)^{-1} \mathbf{C}_{I_\varepsilon^*(x_n)}.\end{aligned} \quad (3.4.9)$$

**for**  $k = 1 \dots \text{maxtrials}$  **do**

  Compute the step

$$x_{n+1} = x_n - \frac{\Delta t}{2^{k-1}} (\alpha_J \boldsymbol{\xi}_J(x_n) + \alpha_C \boldsymbol{\xi}_C(x_n)).$$

**if**  $\text{merit}_{x_n}(x_{n+1}) < \text{merit}_{x_n}(x_n)$  **then**

**break**

**end if**

**end for**

**end for**

---



associated gradient flow for (3.5.1) reads:

$$\begin{cases} \begin{bmatrix} \dot{x}(t) \\ \dot{z}(t) \end{bmatrix} = -\alpha_J(I - \mathbf{DC}^\mathcal{T}(\mathbf{DCDC}^\mathcal{T})^{-1}\mathbf{DC}) \begin{bmatrix} \nabla J(x(t)) \\ 0 \end{bmatrix} - \alpha_C \mathbf{DC}^\mathcal{T}(\mathbf{DCDC}^\mathcal{T})^{-1}\mathbf{C}(x(t), z(t)), \\ \begin{bmatrix} x(0) \\ z(0) \end{bmatrix} = \begin{bmatrix} x_0 \\ z_0 \end{bmatrix}, \end{cases} \quad (3.5.2)$$

where  $(x_0, z_0) \in \tilde{V}$  is a suitable initialization. In our implementation, for comparison purposes, the variable  $z$  is initialized with a value  $z_0 \in \mathbb{R}^q$  in such a way that the inequality constraints of (3.1.2) which are inactive for  $x_0$  (i.e.  $h_i(x_0) < 0$ ) are associated with satisfied equality constraints  $C_{p+i}(x_0, z_0) = 0$  in (3.5.1):

$$\forall i \in \{1, \dots, q\}, z_{0,i} = \sqrt{2|h_i(x_0)|}.$$

In the finite-dimensional setting  $V = \mathbb{R}^k$  and when  $J$ ,  $\mathbf{g}$  and  $\mathbf{h}$  are  $\mathcal{C}^2$  functions, Schropp and Singer proved in [277] that:

- (i) stationary points of the extended flow (3.5.2) are exactly critical points of (3.1.2), that is points  $x^*$  satisfying (3.3.7) but with  $\boldsymbol{\mu}(x^*) \in \mathbb{R}^q$  possibly negative;
- (ii) among all possible critical points, only KKT points (fulfilling all three conditions (3.3.7) with  $\boldsymbol{\mu}(x^*) \in \mathbb{R}_+^q$ ) are asymptotically stable equilibria.

As a consequence, the solution vector  $x(t)$  to (3.5.2) converges in practice to a KKT point for problem (3.1.2).

The main differences between the slack variable approach and our proposed flow (3.3.27) for dealing with equality and inequality constrained problems can be summarized as follows.

1. Any point  $x^{crit}$  satisfying the constraints ( $\mathbf{C}_{\tilde{I}(x^{crit})}(x^{crit}) = 0$ ) and  $\Pi_{\mathbf{C}_{\tilde{I}(x^{crit})}}(\nabla J(x^{crit})) = 0$  is a stationary point of the extended dynamical system (3.5.2), although it might violate the full KKT condition (because (3.3.6) may yield possible negative values of the multiplier  $\boldsymbol{\mu}_{\tilde{I}(x^{crit})}(x^{crit})$ ). In contrast,  $x^*$  is a stationary points of the flow (3.3.27) if and only if it is a true feasible KKT point; see proposition 3.5.
2. The computation of  $\boldsymbol{\xi}_J(x)$  and  $\boldsymbol{\xi}_C(x)$  in our flow (3.3.27) requires to invert a matrix of size at most  $(p + \tilde{q}(x))$ -by- $(p + \tilde{q}(x))$  with  $\tilde{q}(x)$  the number of active or violated constraints at  $x$ . The process of equalizing inequality constraints as in [277, 282] rather requires to invert the full  $(p + q)$ -by- $(p + q)$  matrix  $\mathbf{DC}(x, z)\mathbf{DC}(x, z)^\mathcal{T}$ . Our method is therefore more efficient if  $\tilde{q}(x) \ll q$ , that is if a lot of inequality constraints are inactive.
3. At feasible points, our ODE (3.3.27) follows the best locally admissible descent direction (with respect to the norm of  $V$ ). This is not the case for the extended ODE (3.5.2). Therefore, from a common feasible point  $x$ , the ODE (3.3.27) always decreases the objective function with a steeper slope  $dJ/ds$  with respect to the parameterization induced by the path length  $s$ , defined as a function of the time  $t$  by

$$s(t) = \int_0^t \|\dot{x}(\alpha)\|_V d\alpha. \quad (3.5.3)$$

This property is illustrated in the academic examples of section 3.5, and in particular on Figure 3.9a below.

All in all, our observations based on the simple numerical examples of the next section 3.5.3 tend to illustrate that both flows (3.3.27) and (3.5.2) may have equivalent performances for solving the non linear optimization problem (3.1.2), this performance being measured in term of the total length

$$S = \int_0^{+\infty} \|\dot{x}\|_V dt$$

covered by the optimization path to reach the optimum. However, the two ODEs (3.3.27) and (3.5.2) yield optimization paths of essentially different natures. Our null space flow (3.3.27) ignores inactive constraints and those aligned with the gradient of the objective function. As a result, it produces non smooth paths that are more likely to reach quickly the saturation of the constraint. The extended flow (3.5.2) yields smoother trajectories that more likely stay away from the constraints, at the cost of inverting at every step the full matrix  $\mathbf{DC}(x, z)\mathbf{DC}(x, z)^T$  of the size of the total number of constraints (active and inactive).

### 3.5.2 Comparisons with ‘iterative’ optimization algorithms

This part compares the null space gradient flow with a few classical first order methods for constrained optimization which are not directly interpretable as the discretization of some dynamical system. We consider in the next paragraphs the Augmented Lagrangian Method and Sequential Linear Programming (SLP). Many other algorithms exist and could have been examined, such as the Method of Feasible Directions (MFD) [327, 308] or the Method of Moving Asymptotes (MMA) [298]. These methods can be considered as more or less sophisticated variant of SLP: they both rely on the resolution of linearized subproblems over a trust region that must be determined by the user. These methods present therefore similar characteristics to the SLP method. The reader is referred to the textbooks [244, 70] for further material on iterative methods for constrained optimization.

#### The Augmented Lagrangian method

The Augmented Lagrangian Method is an optimization algorithm which is very popular partly because of its implementation simplicity: it consists in replacing (3.1.2) with a sequence of unconstrained minimization problems

$$\min_{x \in V} \mathcal{L}_k(x) = J(x) + \boldsymbol{\lambda}_k^T \mathbf{g}(x) + \boldsymbol{\mu}_k^T \mathbf{h}(x) + \frac{\alpha_k}{2} \|\mathbf{g}(x)\|^2 + \frac{\alpha_k}{2} \|\mathbf{h}_{\tilde{I}(x)}(x)\|^2. \quad (3.5.4)$$

Here, the tentative Lagrange multipliers have been denoted  $\boldsymbol{\lambda}_k$  and  $\boldsymbol{\mu}_k$  in order to be consistent with our previous notation. The parameter  $\alpha_k$  serves to penalize the violation of the constraints.

A sequence  $x_n$  which *hopefully* converges to a solution to the original problem (3.1.2) is typically computed by alternating gradient steps for the unconstrained minimization of (3.5.4),

$$x_{n+1} = x_n - \Delta t \left( \nabla J(x_n) + \mathbf{Dg}(x_n)^T \boldsymbol{\lambda}_k + \mathbf{Dh}(x_n)^T \boldsymbol{\mu}_k + \alpha_k \mathbf{Dg}(x_n)^T \mathbf{g}(x_n) + \alpha_k \mathbf{Dh}_{\tilde{I}(x_n)}(x_n)^T \mathbf{h}_{\tilde{I}(x_n)} \right) \quad (3.5.5)$$

with updates of the Lagrange multipliers  $\boldsymbol{\lambda}_k$ ,  $\boldsymbol{\mu}_k$ , and possibly (but this is not compulsory) of the penalization parameter  $\alpha_k$ . Note that the iteration indices  $n$  and  $k$  may be independent: one major difficulty of the method lies in the correct initialization and updates of the parameters  $\boldsymbol{\lambda}_k$ ,  $\boldsymbol{\mu}_k$ ,  $\alpha_k$ . Generally, the parameter  $\alpha_k$  is assumed to be constant, and an update rule is commonly assumed for updating  $\boldsymbol{\lambda}_k$  and  $\boldsymbol{\mu}_k$  (see [244]):

$$\lambda_{k+1,i} := \lambda_{k,i} + \alpha_k \mathbf{g}_i(x_k^*), 1 \leq i \leq p, \quad \mu_{k+1,j} := \begin{cases} \mu_{k,j} + \alpha_k h_j(x_k^*) & \text{if } h_j(x_k^*) > -\frac{\mu_{k,j}}{\alpha_k}, 1 \leq j \leq q \\ 0 & \text{otherwise,} \end{cases} \quad (3.5.6)$$

where  $x_k^*$  is a minimizer for (3.5.4). Of course, in practice, (3.5.6) is not computed with the true minimizer  $x_k^*$  but with a current iterate  $x_n$  obtained after one or more iterations of (3.5.5).

Besides being very difficult to use as soon as the problem involves more than one constraint, this algorithm is rather slow, because it uses *tentative* values for the Lagrange multipliers  $\boldsymbol{\lambda}_k$  and  $\boldsymbol{\mu}_k$  rather than the optimal multipliers  $\boldsymbol{\lambda}^*(x_n)$ ,  $\boldsymbol{\mu}^*(x_n)$  of (3.3.10). As a result, it can take a lot of iterations for  $\boldsymbol{\lambda}_k$  and  $\boldsymbol{\mu}_k$  to converge to the correct multipliers, often resulting in oscillations around the constraints at convergence. Recall that actually, lemma 3.3 highlights that our null space algorithm can be seen as an Augmented Lagrangian Method with a ‘clever’ choice of  $\boldsymbol{\lambda}_k$ ,  $\boldsymbol{\mu}_k$ , and  $\alpha_k$  (in fact,  $\alpha_k$  is replaced by the symmetric positive definite matrix  $\mathbf{S}(x_n)$ ).

The examples of section 3.5.3 rely on the implementation of this algorithm available in [150]. The step  $\alpha_k$  is kept constant equal to 1.1, and (3.5.6) is performed after every single iteration of (3.5.5), which means  $n = k$ .

### Sequential Linear Programming

The SLP algorithm is a very popular first order method for equality and inequality constrained problems. The descent direction given by the SLP algorithm is of the form (3.3.34). The main idea of the method is to obtain the next iterated  $x_{n+1}$  by solving a linearized version of the original problem (3.1.2) around the current guess  $x_n$ :

$$x_{n+1} := \arg \min \quad J(x_n) + DJ(x_n) \cdot (x_{n+1} - x_n) \\ \text{s.t.} \quad \begin{cases} \mathbf{g}(x_n) + D\mathbf{g}(x_n) \cdot (x_{n+1} - x_n) = 0 \\ \mathbf{h}(x_n) + D\mathbf{h}(x_n) \cdot (x_{n+1} - x_n) \leq 0 \\ \|x_{n+1} - x_n\|_V \leq \delta. \end{cases} \quad (3.5.7)$$

The constraint  $\|x_{n+1} - x_n\|_V \leq \delta$  requires that  $x_{n+1}$  should be found in a trust region of size  $\delta > 0$  around  $x_n$  (otherwise the value of the subproblem (3.5.7) could be  $-\infty$  due to unbounded minimizing sequences). Note that when the original problem (3.1.2) is set in  $\mathbb{R}^k$ , this constraint is more often expressed as a bound constraint

$$-\delta \leq x_{n+1,i} - x_{n,i} \leq \delta, \quad \forall 1 \leq i \leq k,$$

because this allows to use gradient projection methods [244].

Actually, there is a connection between null space iterates (3.3.34) and SLP iterates. Indeed, assume that the constraint  $\|x_{n+1} - x_n\|_V \leq \delta$  is saturated and that the inequality constraints

$$h_i(x_n) + Dh_i(x_n) \cdot (x_{n+1} - x_n) \leq 0$$

get saturated for  $i$  belonging to some index set  $I_n$ . Then it is readily verified (in the line of the point 2. of lemma 3.1, or in the line of proposition 3.4) that the solution of (3.5.7) is explicitly given by

$$x_{n+1} = x_n - \alpha_n (I - DC_{I_n}^T (DC_{I_n} DC_{I_n}^T)^{-1} DC_{I_n}) \nabla J(x_n) - DC_{I_n}^T (DC_{I_n} DC_{I_n}^T)^{-1} \mathbf{C}_{I_n}(x_n) \quad (3.5.8)$$

where  $\mathbf{C}_I = \begin{bmatrix} \mathbf{g}(x_n) & \mathbf{h}_I \end{bmatrix}^T$  is the vector of saturated linearized constraints, and  $\alpha_n > 0$  is a scaling coefficient given by

$$\alpha_n = \frac{\sqrt{\delta^2 - \|DC_I^T (DC_I DC_I^T)^{-1} \mathbf{C}_I(x_n)\|_V^2}}{\|(I - DC_I^T (DC_I DC_I^T)^{-1} DC_I) \nabla J(x_n)\|_V}.$$

Note that if the trust region size  $\delta$  is too small, i.e. if

$$\delta < \|DC_I^T (DC_I DC_I^T)^{-1} \mathbf{C}_I(x_n)\|_V$$

then the problem (3.5.7) does not have a solution.

The scheme (3.5.8) is a variant of the discretization (3.3.34) of the null space gradient flow, up to the redefinition of  $\alpha_C$  and  $\alpha_J$ . The index set  $I_n$  may also be different from  $\hat{I}(x_n)$ , because the subproblem (3.5.7) sees all the constraints while the null space method sees only the active or violated ones.

As is underlined by the expression of  $\alpha_n$ , the major problem of the SLP method is the fact that the SLP subproblem (3.5.7) may not have a solution if the trust region parameter  $\delta$  is too small. Some additional tuning operations are usually done to face this issue, see e.g. [174]. This issue is partly due to the fact that the update scheme (3.5.8) is not the discretization of some ODE, because the range space step  $DC_{I_n}^T (DC_{I_n} DC_{I_n}^T)^{-1} \mathbf{C}_{I_n}(x_n)$  has not been scaled by some small parameter (e.g. by  $\delta$  in our case). The SLP method directly imposes constraints to be satisfied at first order, whereas our null space method requires rather constraints to decrease by a factor  $\alpha_C \Delta t$  at the next iteration (see the discussion at the beginning of section 3.2.2 and remark 3.3). Note that in contrast with the SLP method, the primal and dual subproblems (3.3.9) and (3.3.10) always admit a solution.

In the following examples of section 3.5.3, the SLP method is implemented with the code available in [150]. Let us note that this implementation of the SLP method does not solve the subproblem (3.5.7) in the unfeasible domain, but a variant in order to deal with violated constraints.

### 3.5.3 Comparative academic test cases in the euclidean space $\mathbb{R}^k$

In this section, we consider simple and illustrative academic examples in order to compare qualitatively the above strategies for dealing with inequality constraints in optimization problems; we consider:

- the method of [section 3.5.1](#) for equalizing inequality constraints by means of slack variables; see [\(3.3.27\)](#). This strategy is hereafter labeled as ‘SLACK’;
- the proposed null space flow [\(3.3.27\)](#) in [section 3.3.2](#), based on the dual problem [\(3.3.10\)](#) for solving the combinatorial character of the constraints. This method is labeled as ‘NLSPACE’;
- an alternative, naive version of [\(3.3.27\)](#) which does not take advantage of the use of a dual problem, and simply projects  $\nabla J(x(t))$  on all the violated constraints:

$$\begin{cases} \dot{x} = -\alpha_J \tilde{\xi}_J(x(t)) - \alpha_C \xi_C(x(t)) \\ \tilde{\xi}_J(x) := (I - \text{DC}_{\tilde{I}(x)}^T (\text{DC}_{\tilde{I}(x)} \text{DC}_{\tilde{I}(x)}^T)^{-1} \text{DC}_{\tilde{I}(x)}) \nabla J(x) \\ \xi_C(x) := \text{DC}_{\tilde{I}(x)}^T (\text{DC}_{\tilde{I}(x)} \text{DC}_{\tilde{I}(x)}^T)^{-1} \text{C}_{\tilde{I}(x)}(x). \end{cases} \quad (3.5.9)$$

In other words, all the violated or saturated constraints are taken into account in the computation of both the null space and range space directions  $\xi_J(x)$  and  $\xi_C(x)$ . This strategy is labeled as ‘NLSPACE (no dual)’;

- the Augmented Lagrangian Method of [section 3.5.2](#), labeled as ‘AULG’;
- the Sequential Linear Programming algorithm of [section 3.5.2](#), labeled as ‘SLP’.

To achieve our comparison purpose, [algorithm 3.1](#) is implemented for the discretization of [\(3.3.27\)](#), [\(3.5.2\)](#) and [\(3.5.9\)](#), with straightforward adaptations for equalizing slack variables or disabling the resolution of the dual problem. In all the considered cases, we have set the values of  $\alpha_J$  and  $\alpha_C$  such that  $\alpha_J/\alpha_C = 5/3$ . The step size  $\Delta t$  for the discretization of the ODE [\(3.1.3\)](#) was chosen sufficiently small to compute continuous paths with satisfying accuracy. Our discussion is exclusively focused on the continuous trajectories of the considered ODEs. In particular, we do neither discuss the issue of the selection of the time step, nor the efficiency of these methods in terms of the needed number of iterations required to achieve convergence.

In order to compare the three methods without bias, we consider the arc length  $s(t)$  (defined in [\(3.5.3\)](#)) as the common reference time for the three ODEs [\(3.3.27\)](#), [\(3.5.2\)](#) and [\(3.5.9\)](#); recall indeed that this quantity is invariant under any monotone parameterization change of the time  $t$ . In the convergence figures below, optimized quantities are then plotted with respect to the pseudo time  $s(t)$  in abscissa; for example we plot the graph  $t \mapsto (s(t), J(x(s(t))))$  in order to account for the evolution of the objective function  $J$ . The SLP and Augmented Lagrangian methods are not dynamical system approaches: we decided not to include them when plotting convergence histories for the objective function and constraints because the arc length  $s(t)$  would not make sense for these.

We shall also plot the evolution of the Lagrange multipliers  $s \mapsto \mu(x(s))$  associated with  $\xi_J(x(s))$  or  $\tilde{\xi}_J(x(s))$  for the ODEs [\(3.3.27\)](#) and [\(3.5.9\)](#). For that purpose, these Lagrange multipliers are defined on the violated indices  $i \in \tilde{I}(x(s))$  by [\(3.3.18\)](#) for the null space flow [\(3.3.27\)](#), and by

$$\mu(x(s)) := -(\text{DC}_{\tilde{I}(x(s))} \text{DC}_{\tilde{I}(x(s))}^T)^{-1} \text{DC}_{\tilde{I}(x(s))} \nabla J(x(s)) \quad (3.5.10)$$

for the flow [\(3.5.9\)](#) that does not use the dual problem [\(3.3.10\)](#). For the indices  $i \in \{1, \dots, q\} \setminus \tilde{I}(x(s))$ , the value of the Lagrange multiplier is set to  $\mu_i(x(s)) := 0$  by convention. We do not plot Lagrange multipliers for the ODE [\(3.5.2\)](#) using slack variables because these are defined with respect to the extended variables  $(x(s), z(s))$ .

The examples of this section take place in the optimization set  $V = \mathbb{R}^2$ , which is equipped with the usual euclidean inner product; the Hilbert transposition  $\mathcal{T} = T$  coincides with the usual transposition operator (see [definition 3.1](#)). For simplicity, these examples only involve inequality constraints; we consider the following three scenarios:

**Test case 1:** the initial point is unfeasible and the gradient of the objective function  $\nabla J(x)$  is always aligned with the directions of the constraints;

**Test case 2:** the initial point is unfeasible, but the gradient  $\nabla J(x)$  may not be aligned with the direction of constraints;

**Test case 3:** one of the constraints becomes inactive in the course of the optimization path.

**Test case 1 : unfeasible initialization with initial gradient aligned with the constraints.**

Our first example features the following problem, reproduced from [150]:

$$\begin{aligned} \min_{(x_1, x_2) \in \mathbb{R}^2} \quad & J(x_1, x_2) := x_2 + 0.3x_1 \\ \text{s.t.} \quad & \begin{cases} h_1(x_1, x_2) := -x_2 + \frac{1}{x_1} \leq 0, \\ h_2(x_1, x_2) := x_1 + x_2 - 3 \leq 0. \end{cases} \end{aligned} \quad (3.5.11)$$

This test case is designed so that for the chosen initial point  $x_0 = (1.5, 2.25)$ , the gradient of the objective function  $\nabla J(x)$  is ‘aligned’ with the linear constraint  $h_2$ , in the sense that

$$-\nabla h_2(x_0) \cdot \nabla J(x_0) < 0.$$

Hence at least for small times (in fact during the whole optimization path), the constraint  $h_2$  can be ignored since the minimization of  $J$  is naturally concurrent with a decrease of the value of  $h_2$ .

The optimization paths taken by the solutions of the three ODEs (3.3.27), (3.5.2) and (3.5.9) are plotted on Figure 3.2. The associated convergence histories for the values of the objective and constraint functions are displayed on Figure 3.3. Note the oscillations characterizing the Augmented Lagrangian Method to travel tangentially to the constraint boundary. The SLP algorithm is able to compute the tangential projection of  $\nabla J(x(t))$  onto the constraint set; it finds a path that is similar to the one of the null space method.

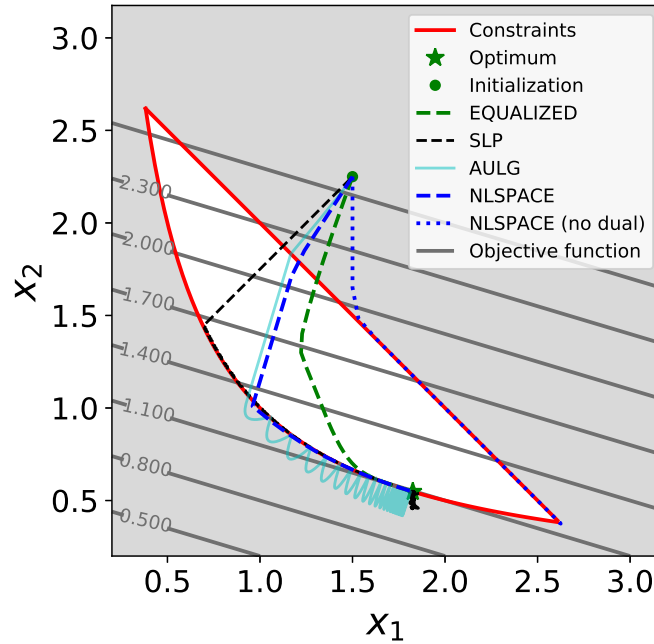


Figure 3.2: Optimization test case 1: optimization paths for an unfeasible initialization  $x_0$  with  $\nabla J(x_0)$  aligned in the direction of the constraints.

Let us comment the trajectory followed by the null space flow (3.3.27) in details. The gradient of the objective function remains aligned with the constraint  $h_2$ , which is associated with a zero Lagrange multiplier  $\mu_2(x(s))$  (see Figure 3.4). During the first part of the optimization, the first constraint  $h_1$  is not violated, hence the multiplier  $\mu_1(x(s))$  is also set to zero. As a result, both constraints are ignored when computing the null space direction, which is set equal to the gradient:  $\xi_J(x(s)) = \nabla J(x(s))$ . The optimization path  $x(s)$  follows then almost the direction of the gradient  $\nabla J(x(s))$  (without projection), up to a small deviation induced by the non zero Gauss-Newton direction  $\xi_C(x(s))$  in the unfeasible domain. When the hyperbolic constraint represented by  $h_1$  becomes violated, the gradient is not aligned anymore with this constraint and the dual problem (3.3.10) yields a non-zero Lagrange multiplier  $\mu_1(x(s))$  (near  $s = 1.4$ ). From this point, the gradient  $\nabla J(x(s))$  is then projected tangentially to the constraint  $h_1$  till the optimum is attained.

In contrast, the path of the ODE (3.5.9) (which does not involve the dual problem (3.3.10)) fails to find the optimum as it is unable to unstick from the first saturated constraint. Notably, the gradient  $\nabla J(x)$  is kept being projected tangentially to the violated constraint  $h_2$  while it should not, which could have been detected from the negativity of the computed Lagrange multiplier  $\mu_2$  (see Figure 3.4).

Finally, the extended ODE (3.5.2) making use of slack variables feels the constraint  $h_1$  from some distance, inducing a deviation of the trajectory in the direction of the optimum before reaching the saturation of the constraint. This allows the trajectory  $x(s)$  to find globally a slightly shorter path than our null space flow (3.3.27).

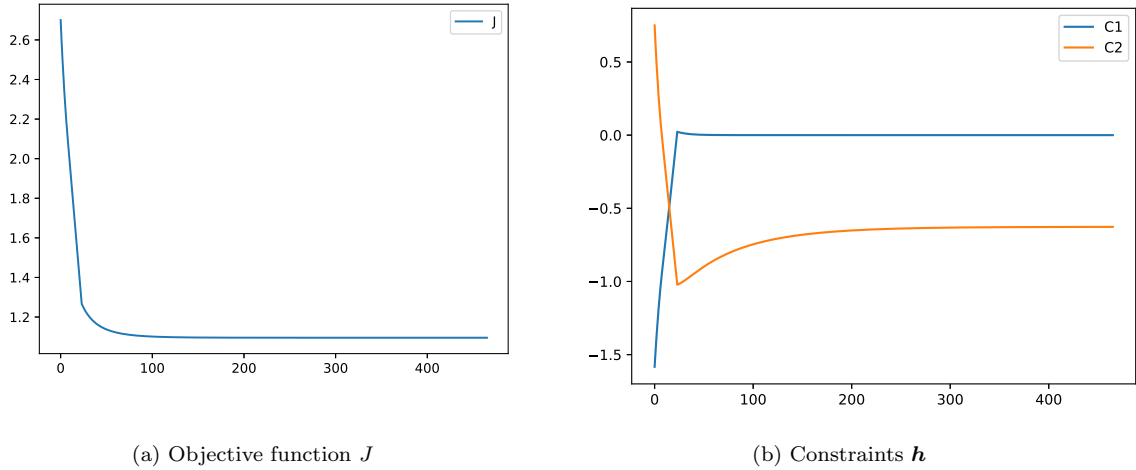


Figure 3.3: History curves for the optimization test case 1.

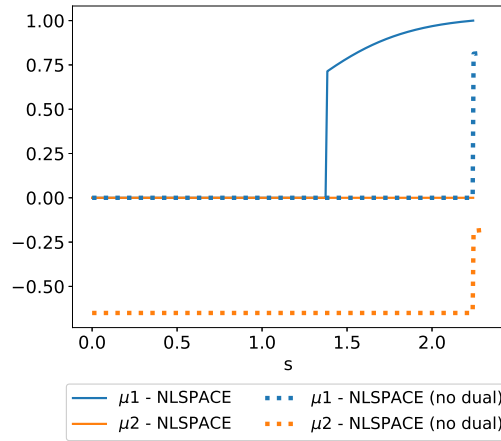


Figure 3.4: Evolution of the Lagrange multipliers  $\mu_1(x(s)), \mu_2(x(s))$  for the optimization test case 1.

### Test case 2 : unfeasible initialization with initial gradient not aligned with the constraints.

We now devise a test case where the gradient of the initialization is not aligned with the constraints. The feasible domain is the same as in the previous test case but the objective function is different:

$$\begin{aligned} \min_{(x_1, x_2) \in \mathbb{R}^2} \quad & J(x_1, x_2) := (x_1 - 2)^2 + (x_2 - 2)^2 \\ \text{s.t.} \quad & \begin{cases} h_1(x_1, x_2) := -x_2 + \frac{1}{x_1} \leq 0 \\ h_2(x_1, x_2) := x_1 + x_2 - 3 \leq 0. \end{cases} \end{aligned} \quad (3.5.12)$$

We keep the same initialization  $x_0 = (1.5, 2.25)$ . Corresponding optimization paths and convergence curves are displayed on Figs. 3.5 and 3.6.

For this example, the linear constraint  $h_2$  is not aligned with the gradient along the optimization path of the null space gradient flow (3.3.27). This is associated with a non-zero Lagrange multiplier  $\mu_2(x(s)) > 0$  (see Figure 3.7): the gradient  $\nabla J(x(s))$  is kept being projected tangentially to the constraint  $h_2$  when computing  $\xi_J(x(s))$ . Here, the combination with the Gauss-Newton direction  $\xi_C(x(s))$  allows to decrease simultaneously the objective function and the violation of the constraints, which enables the optimization path to reach directly the optimum when hitting the feasible set. Note that the convergence curve Figure 3.6a depicts a monotonically increasing objective function  $J$  after  $s \geq 0.5$  (although we are minimizing  $J$ ); this is of course due to the fact that constraints are never satisfied.

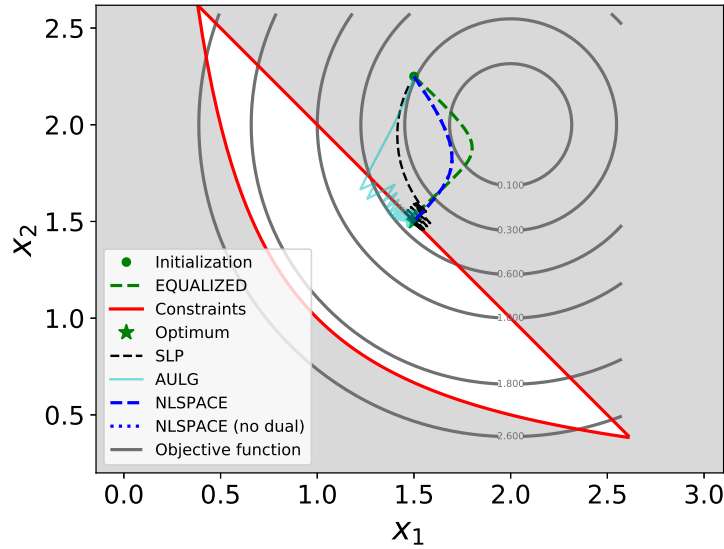


Figure 3.5: Optimization problem of section 3.5.3: unfeasible initialization  $x_0$  with  $\nabla J(x_0)$  not aligned in the direction of the constraints.

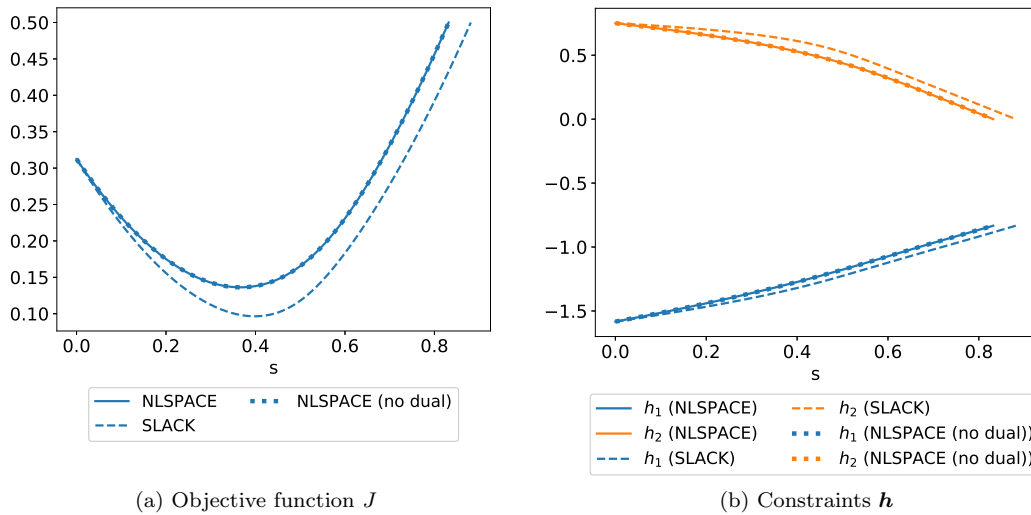


Figure 3.6: History curves for the optimization test case 2.

Since the constraint  $h_2$  remains active from initialization to the optimum, the optimization path is unchanged when disabling the dual problem (ODE (3.5.9)). Finally, the path selected by the extended flow (3.5.2) to reach the optimum is very similar to the one of our method, although slightly longer.

### Test case 3: a saturated inequality constraint becoming inactive along the optimization path

This last optimization test case is designed to illustrate the relevance of the dual problem for detecting when a saturated inequality constraint ceases to be saturated. We consider a disconnected unfeasible

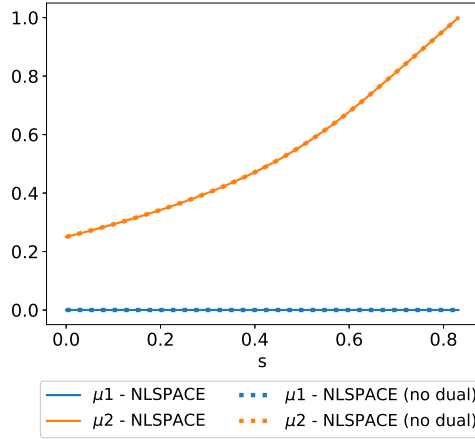


Figure 3.7: Evolution of the Lagrange multipliers  $\mu_1(x(s)), \mu_2(x(s))$  for the optimization test case 2.

domain made from the reunion of a half-space and the interior region of a parabola:

$$\begin{aligned} \min_{(x_1, x_2) \in \mathbb{R}^2} \quad & J(x_1, x_2) = x_1^2 + (x_2 + 3)^2 \\ \text{s.t.} \quad & \begin{cases} h_1(x_1, x_2) = -x_1^2 + x_2 \leq 0 \\ h_2(x_1, x_2) = -x_1 - x_2 - 2 \leq 0. \end{cases} \end{aligned} \quad (3.5.13)$$

The feasible domain and optimization paths starting from the initialization  $x_0 = (3, 3)$  are displayed on [Figure 3.8](#). Associated convergence curves are reported on [Figure 3.9](#). The optimization paths obtained with the Augmented Lagrangian Method and the SLP method are also depicted. We observed oscillations of the SLP method in the region close to both constraints. We note as well, again, the tendency of the Augmented Lagrangian Method iterates to oscillate around the constraints.

For the null space flow [\(3.3.27\)](#), four different stages occur as is visible on the evolution of the Lagrange multipliers  $\mu_1(x(s)), \mu_2(x(s))$  reported on [Figure 3.10](#):

1. From  $s = 0$  to  $s = 1.73$ , the trajectory  $x(s)$  remains in the feasible domain. Lagrange multipliers  $\mu_1(x(s)) = \mu_2(x(s)) = 0$  are set to 0 and the null space direction  $\xi_J(x(s)) = \nabla J(x(s))$  coincides with the gradient of the objective function, until  $x(s)$  hits the parabolic barrier, which corresponds to the saturation of the first constraint  $h_1$ .
2. From  $s = 1.73$  to  $s = 4.4$ , the resolution of the dual problem yields a non zero multiplier  $\mu_1(x(s)) > 0$ . The optimization trajectory  $x(s)$  remains tangent to the first constraint, until reaching a limit point such that  $\nabla J(x(s)) \cdot \nabla h_1(x(s)) = 0$ . At this moment, it is not necessary to project the gradient tangentially to this constraint any more, and the values of both Lagrange multipliers  $\mu_1(x(s)) = \mu_2(x(s))$  are equal to 0.
3. From  $s = 4.4$  to  $s = 6.5$ , both constraints  $h_1$  and  $h_2$  are ignored and the trajectory  $x(s)$  follows the gradient  $-\nabla J(x(s))$ , till the saturation of  $h_2$ .
4. From  $s = 6.5$  to  $s = 7.1$ , the second Lagrange multiplier  $\mu_2(x(s)) > 0$  has a positive value;  $x(s)$  evolves then tangentially to this constraint till the optimum is attained.

As illustrated on [Figure 3.8](#), the use of the dual problem is key in the detection of the moment when the optimization trajectory  $x(s)$  needs to be released from active inequality constraints. Because of the discrete nature of the time stepping, the path followed by the ODE [\(3.5.9\)](#) necessarily enters slightly the violated parabolic domain. Since it does not use the information provided by the dual problem [\(3.3.10\)](#), the gradient  $\nabla J(x(s))$  is kept being projected tangentially to the constraint  $h_1$  till  $x(s)$  converges to some stationary point (which is not a KKT point). As can be seen on [Figure 3.10](#), the optimization trajectory followed by this ODE coincides with the one of the flow [\(3.3.27\)](#), till the instant  $s = 4.4$  at which the Lagrange multiplier  $\mu_1(x(s))$  becomes negative (which violates the feasibility condition of the dual problem [\(3.3.10\)](#)). Note that using larger steps could have allowed the trajectory  $x(s)$  to exit “by chance” the unfeasible domain, and in that case convergence to the optimum would have been obtained. However this would not reflect the actual behavior of the continuous solutions of [\(3.5.9\)](#).



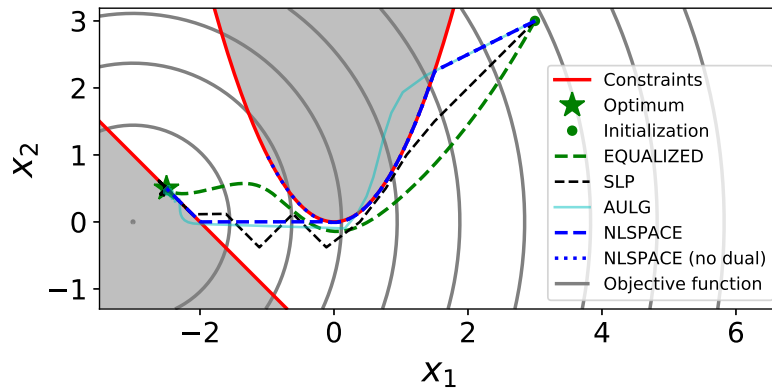


Figure 3.8: Optimization problem of section 3.5.3 : feasible initialization  $x_0$  but the optimization has to find a path across the parabolic domain.

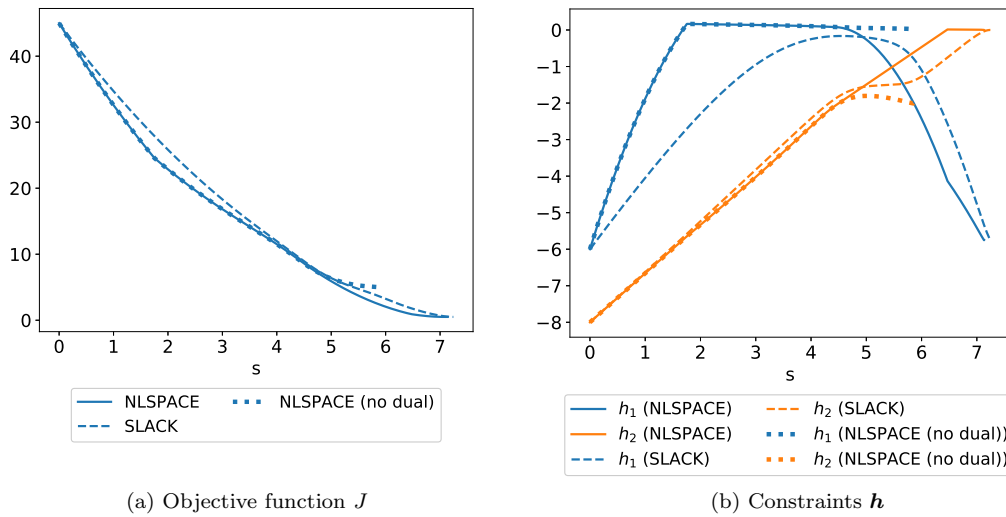


Figure 3.9: History curves of the null space algorithm for the optimization problem of section 3.5.3.

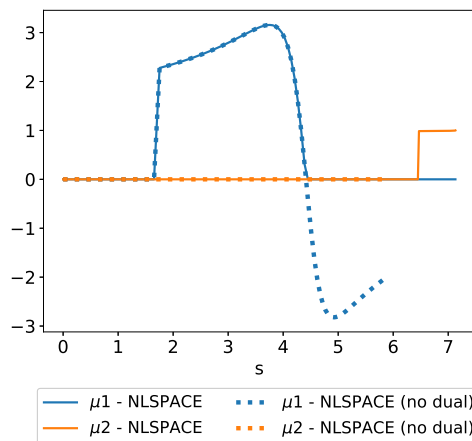


Figure 3.10: Evolution of Lagrange multipliers  $\mu_1(x(s)), \mu_2(x(s))$  for the optimization problem of section 3.5.3.

Finally, the extended ODE (3.5.2) using slack variables finds a smooth path to the optimum. Since inactive constraints are felt from distance, the trajectory  $x(s)$  is able to remain more strictly in the feasible domain for all times. The total length of the optimization path is almost the same than the one of the null space flow (3.3.27) (note the steeper descending slopes  $dJ/ds$  for the latter at intersection points of the two trajectories, see the point (3) in the discussion of section 3.5.1).

### 3.6 OPTIMIZATION WITHIN THE SET OF LIPSCHITZ SUBDOMAINS: APPLICATIONS TO SHAPE OPTIMIZATION

Our ultimate goal is to apply our optimization strategy to shape and topology optimization problems. In such applications, the optimization set is not a vector space  $V$  as in (3.1.2) but a more general set  $\mathcal{X}$  of shapes in  $\mathbb{R}^d$  ( $d = 2$  or  $3$  in standard applications):

$$\mathcal{X} = \{\Omega \subset D \mid \Omega \text{ Lipschitz}\}, \quad (3.6.1)$$

where  $D \subset \mathbb{R}^d$  is an enclosing ‘hold-all’ domain. Since  $\mathcal{X}$  is not a Hilbert space, the present context does not fall into the optimization framework described in sections 3.2 and 3.3. However,  $\mathcal{X}$  may be endowed with a manifold structure, which makes it possible to extend our dynamical system (3.3.1) to this context, up to small adaptations that we now describe.

In the whole section (and in the next chapters of the thesis), we consider a generic shape optimization problem

$$\begin{aligned} \min_{\Omega \in \mathcal{X}} \quad & J(\Omega) \\ \text{s.t.} \quad & \begin{cases} \mathbf{g}(\Omega) = 0 \\ \mathbf{h}(\Omega) \leq 0, \end{cases} \end{aligned} \quad (3.6.2)$$

for shape differentiable objective and constraint functions  $J : \mathcal{X} \rightarrow \mathbb{R}$ ,  $\mathbf{g} : \mathcal{X} \rightarrow \mathbb{R}^p$  and  $\mathbf{h} : \mathcal{X} \rightarrow \mathbb{R}^q$  (in the sense of chapter 1, definition 1.1). In section 3.6.1, we outline the analogy between the resolution of (3.6.2) with the method of Hadamard and classical methods for the optimization on smooth manifolds. Then, we explain how to extend the null space algorithm 3.1 to this setting, and in particular how the classical extension and regularization step of shape derivatives (outlined in chapter 1, section 1.4.1) is naturally included in our method when using the definition (3.2.3) of the Hilbertian transposition  $\mathcal{T}$ . A few implementation details regarding the normalization of the descent direction with respect to the mesh size are discussed in section 3.6.2. Finally, section 3.6.3 concludes this chapter with some further numerical illustrations of our constrained optimization algorithm in the context of shape optimization: we consider the model example of shape optimization of a bridge structure subjected to multiple load cases and featuring 10 constraints.

#### 3.6.1 Manifold structures for shape optimization

Our extension of the previous material to the shape optimization context is inspired by ‘classical’ optimization strategies on a smooth embedded manifold  $\mathcal{M} \subset \mathbb{R}^k$ . In this context, a descent direction at a point  $x_n \in \mathcal{M}$  for some objective functional is typically sought as an element  $\boldsymbol{\xi}_n \in T_{x_n}\mathcal{M}$  of the tangent space  $T_{x_n}\mathcal{M}$  to  $\mathcal{M}$  at  $x_n$ ; see e.g. [141, 5]. Then one relies on a *retraction*  $\rho_{x_n}$ , that is a mapping

$$\rho_{x_n} : T_{x_n}\mathcal{M} \rightarrow \mathcal{M}$$

satisfying the following two consistency conditions:

$$\begin{cases} \rho_{x_n}(0) = x_n \\ \forall \boldsymbol{\xi} \in T_{x_n}\mathcal{M}, \left. \frac{d}{dt} \right|_{t=0} \rho_{x_n}(t\boldsymbol{\xi}) = \boldsymbol{\xi}. \end{cases}$$

The mapping  $\rho_{x_n}$  then allows to convert  $\boldsymbol{\xi}_n$  into a practical update of the actual point  $x_n$  on  $\mathcal{M}$ :

$$x_{n+1} := \rho_{x_n}(\Delta t \boldsymbol{\xi}_n), \quad (3.6.3)$$

where  $\Delta t > 0$  is the descent step; see [6] and Figure 3.11. Since the new point  $x_{n+1}$  belongs to  $\mathcal{M}$ , this procedure can be repeated iteratively.

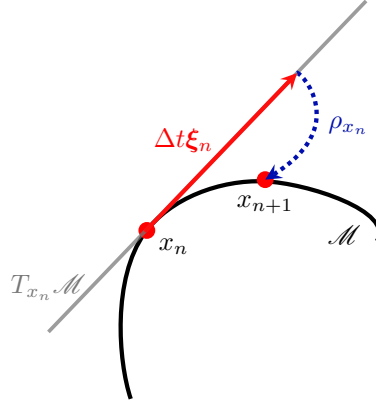


Figure 3.11: Optimization on a manifold  $\mathcal{M}$ : a retraction map  $\rho_{x_n}$  is used to project a tangential motion  $\Delta t \xi_n \in T_{x_n} \mathcal{M}$  from  $x_n \in \mathcal{M}$  back onto the optimization domain  $\mathcal{M}$ .

The same idea can be used to apply the methods of [sections 3.2](#) and [3.3](#) to the optimization problem [\(3.6.2\)](#), posed over the set of shapes  $\mathcal{X}$ . To this end, we rely on Hadamard's method (reviewed in [chapter 1](#)), which considers variations of a shape  $\Omega \in \mathcal{X}$  of the form

$$\rho_\Omega(\boldsymbol{\theta}) := (I + \boldsymbol{\theta})(\Omega), \text{ for } \boldsymbol{\theta} \in W^{1,\infty}(D, \mathbb{R}^d) \text{ with } \|\boldsymbol{\theta}\|_{W^{1,\infty}(D, \mathbb{R}^d)} < 1, \quad (3.6.4)$$

Formally, the set  $W^{1,\infty}(D, \mathbb{R}^d)$  may be interpreted as the tangent space to  $\mathcal{X}$  at  $\Omega$  and the mapping  $\rho_\Omega$ , which is defined by [\(3.6.4\)](#) on a neighborhood of  $\boldsymbol{\theta} = 0$  in  $W^{1,\infty}(D, \mathbb{R}^d)$ , plays the role of a retraction. Other definitions are possible for such a transformation dictating how a shape should evolve according to a vector field  $\boldsymbol{\theta}$ , see the discussion of [chapter 1, section 1.4.2](#). Note also that more rigorous manifold structures on shape spaces can be formulated, see e.g. [\[44, 278\]](#).

Usually, in the context of a general embedded manifold  $\mathcal{M} \subset \mathbb{R}^k$ , a differential structure on  $\mathcal{M}$  is defined first (inducing a notion of derivative on  $\mathcal{M}$ ), and the definition of a suitable retraction is inferred accordingly. In the framework of Hadamard's method however, it is the retraction  $\rho_\Omega$  itself, that is the parametrization [\(3.6.4\)](#) by deformation fields  $\boldsymbol{\theta}$ , that is used to define the notion of derivative. Indeed, the definition of shape derivative of [chapter 1, definition 1.1](#) is equivalent to saying that a functional  $J(\Omega)$  is shape differentiable if and only if  $J \circ \rho_\Omega$  is differentiable.

The key ingredient for applying the null space algorithm [algorithm 3.1](#) to this context is the computation of the gradient  $\nabla J(\Omega)$  and transposes  $D\mathbf{g}(\Omega)^\top$  and  $D\mathbf{h}(\Omega)^\top$  of the shape derivatives of the constraints. Following [definition 3.1](#), these are computed by introducing a Hilbert space  $V \subset W^{1,\infty}(D, \mathbb{R}^d)$  with scalar product  $\langle \cdot, \cdot \rangle_V$  and by solving identification problems [\(3.2.3\)](#) and [\(3.2.4\)](#). Our practical implementation follows the strategy described in [chapter 1, section 1.4.1](#): we set either

$$V = H^1(D, \mathbb{R}^d) \text{ with } \forall \boldsymbol{\theta}, \boldsymbol{\theta}' \in V, \langle \boldsymbol{\theta}, \boldsymbol{\theta}' \rangle_V := \int_D (\gamma^2 \nabla \boldsymbol{\theta} : \nabla \boldsymbol{\theta}' + \boldsymbol{\theta} \cdot \boldsymbol{\theta}') dx \quad (3.6.5)$$

or

$$V = \{v \nabla d_\Omega \mid v \in H^1(D)\} \text{ with } \forall v, w \in V, \langle v \nabla d_\Omega, w \nabla d_\Omega \rangle_V := \int_D (\gamma^2 \nabla v \cdot \nabla w + vw) dx. \quad (3.6.6)$$

The implicit dependence of  $V$  with respect to  $\Omega$  is omitted. As mentioned in [chapter 1](#), the parameter  $\gamma$  is set proportional to the minimum mesh element size  $\mathbf{hmin}$ .

In light of the previous discussion, the proposed dynamical system [\(3.3.1\)](#) for tackling shape optimization problems of the form [\(3.6.2\)](#) is extended and discretized as follows.

1. The null space and range space directions  $\xi_J(\Omega)$  and  $\xi_C(\Omega)$  are computed as elements of  $V$  thanks to the formulas [\(3.3.8\)](#) and [\(3.3.25\)](#). This requires the computation of the gradient  $\nabla J(\Omega)$  and of the transposes  $D\mathbf{g}^\top(\Omega)$ ,  $D\mathbf{h}^\top(\Omega)$  via the resolution of identification problems such as [\(3.2.3\)](#) and [\(3.2.4\)](#). In particular, steps 1 to 5 of [algorithm 3.1](#) including the resolution of the dual problem [\(3.3.10\)](#) are achieved from the knowledge of the Fréchet derivatives and of their transposes.
2. The update [\(3.3.2\)](#) of the design from one iteration of the process to the next is performed by using the retraction map  $\rho_\Omega$  as in [\(3.6.3\)](#):

$$\Omega_{n+1} := \rho_{\Omega_n}(-\Delta t(\alpha_J \xi_J(\Omega_n) + \alpha_C \xi_C(\Omega_n))); \quad (3.6.7)$$

in practice, the step 6 of [algorithm 3.1](#) is adapted accordingly.

In our applications, let us recall that [\(3.6.7\)](#) is carried out by relying on the mesh evolution technique reviewed in [chapter 1, section 1.4.2](#): at every iteration  $n$ , the current shape  $\Omega_n$  is explicitly discretized as a submesh of a triangulated mesh  $\mathcal{T}_n$  of  $D$  (see e.g. [Figure 3.15](#) below). This also means that one does not use the retraction map  $\rho_\Omega$  of [\(3.6.4\)](#) naturally considered by the method of Hadamard but rather the one  $\tilde{\rho}_\Omega$  associated with the use of the level set method described in eqn. [\(1.4.12\)](#) of [chapter 1](#).

### 3.6.2 Adaptive normalizations for the null space and range space directions

We rely on [algorithm 3.1](#) for our implementation of the null space flow [\(3.3.1\)](#) for numerical shape optimization. A few comments are in order regarding the appropriate scaling of the null and range space steps in relation with the size of the mesh discretization; we define accordingly variable coefficients  $\alpha_{J,n}$  and  $\alpha_{C,n}$  for the descent direction  $\boldsymbol{\theta}_n$  in [\(3.6.8\)](#).

For stability reasons, the vertices of the current mesh  $\mathcal{T}_n$  discretizing  $\Omega_n$  should move by a distance which equals at most a few mesh elements in order to produce the subsequent shape  $\Omega_{n+1}$ . Hence, the minimum edge length  $\mathbf{hmin}$  of the computational mesh is a natural candidate for the limiting step size value  $\mathbf{h}$  of the discussion in [section 3.4.1](#). In our practical implementation, we set  $\Delta t = 1$  and a descent direction  $\boldsymbol{\theta}_n(x)$  is computed by estimating

$$\boldsymbol{\theta}_n := -(\alpha_{J,n}\boldsymbol{\xi}_J(\Omega_n) + \alpha_{C,n}\boldsymbol{\xi}_C(\Omega_n)), \quad (3.6.8)$$

where  $\alpha_J$  and  $\alpha_C$  of the update [\(3.6.7\)](#) have been replaced by dynamic coefficients  $\alpha_{J,n}$  and  $\alpha_{C,n}$ .

The parameters  $\alpha_{J,n}$  and  $\alpha_{C,n}$  scaling the null space and range space steps  $\boldsymbol{\xi}_J(\Omega_n)$  and  $\boldsymbol{\xi}_C(\Omega_n)$  are updated dynamically in order to control the step size  $\|\boldsymbol{\theta}_n\|_{L^\infty(D,\mathbb{R}^d)}$ . Note that the infinity norm is considered because *all* values of the displacement  $\boldsymbol{\theta}_n$  should be of the order of the mesh size. We consider  $A_J$  and  $A_C$  two user-defined parameters, which are expressed in terms of the minimum edge length  $\mathbf{hmin}$  for a clearer intuitive meaning. The coefficients  $\alpha_{J,n}$  and  $\alpha_{C,n}$  are updated at every iteration according to the following rules:

$$\alpha_{J,n} := \begin{cases} \frac{A_J \mathbf{hmin}}{\|\boldsymbol{\xi}_J(\Omega_n)\|_{L^\infty(D,\mathbb{R}^d)}} & \text{if } n < \mathbf{n}_0 \\ \frac{A_J \mathbf{hmin}}{\max(\|\boldsymbol{\xi}_J(\Omega_n)\|_{L^\infty(D,\mathbb{R}^d)}, \|\boldsymbol{\xi}_J(\Omega_{\mathbf{n}_0})\|_{L^\infty(D,\mathbb{R}^d)})} & \text{if } n \geq \mathbf{n}_0 \end{cases} \quad (3.6.9)$$

$$\alpha_{C,n} := \min \left( 0.9, \frac{A_C \mathbf{hmin}}{\max(1\mathbf{e}-9, \|\boldsymbol{\xi}_C(\Omega_n)\|_{L^\infty(D,\mathbb{R}^d)})} \right). \quad (3.6.10)$$

These normalizations ensure that the null space and range space steps always remain smaller than  $A_J$  and  $A_C$  times the mesh size:

$$\forall n \geq 0, \|\alpha_{J,n}\boldsymbol{\xi}_J(\Omega_n)\|_{L^\infty(D,\mathbb{R}^d)} \leq A_J \mathbf{hmin} \text{ and } \|\alpha_{C,n}\boldsymbol{\xi}_C(\Omega_n)\|_{L^\infty(D,\mathbb{R}^d)} \leq \min(0.9, A_C \mathbf{hmin}).$$

Actually, the null space component  $\alpha_{J,n}\boldsymbol{\xi}_J(\Omega_n)$  of  $\boldsymbol{\theta}_n$  is scaled to be exactly of the size  $A_J \mathbf{hmin}$  for the first  $\mathbf{n}_0$  iterations:

$$\forall 1 \leq n \leq \mathbf{n}_0, \|\alpha_{J,n}\boldsymbol{\xi}_J(\Omega_n)\|_{L^\infty(D,\mathbb{R}^d)} = A_J \mathbf{hmin}.$$

Then,  $\boldsymbol{\xi}_J(\Omega_n)$  is allowed to converge to 0 as  $n \rightarrow \infty$ .

The range step  $\alpha_{C,n}\boldsymbol{\xi}_C(\Omega_n)$  is also set to remain smaller than the constant 0.9, in view of the stability condition  $0 < \alpha_C \Delta t < 2$  (see [remark 3.12](#)). The role of the constant  $1\mathbf{e}-9$  is only to avoid division by 0 when no constraint is active.

**Remark 3.13.** Since we measure step sizes with the infinity norm  $\|\boldsymbol{\theta}_n\|_{L^\infty(D)}$  rather than with the Hilbertian norm  $\|\boldsymbol{\theta}_n\|_V = \|\boldsymbol{\theta}_n\|_{H^1(D,\mathbb{R}^d)}$ , the tolerance bounds [\(3.4.3\)](#) need to be updated with respect to this norm as follows:

$$\varepsilon_i := \mathbf{hmin} \int_{\partial\Omega} |v_{C_i}(\Omega_n)| ds,$$

where it is assumed that the shape derivative of each constraint functional  $C_i(\Omega_n)$  can be written as a boundary integral featuring the scalar field  $v_{C_i}(\Omega_n)$ :

$$DC_i(\Omega_n)(\boldsymbol{\theta}) := \int_{\partial\Omega} v_{C_i}(\Omega_n) \boldsymbol{\theta} \cdot \mathbf{n} ds.$$

### 3.6.3 Illustrations on a multiple load structural shape optimization test case

In this final section, we illustrate the efficiency of our optimization strategy on a practical structural design problem. Two possible configurations are investigated for the shape optimization of a bridge structure subjected to multiple loads, featuring either multiple objective criteria or multiple constraint functions.

#### Shape optimization setting

We consider the shape optimization of a bridge-like structure  $\Omega$  contained in a two-dimensional rectangular hold-all domain  $D \subset \mathbb{R}^2$  with size  $10 \times 2$ . The boundary of  $\partial\Omega$  is divided into disjoint regions as:

$$\partial\Omega = \Gamma \cup \Gamma_D \cup \bigcup_{i=0}^8 \Gamma_i,$$

where

- $\Gamma_D$  is a non-optimizable part of the boundary on which the structure  $\Omega$  is clamped, made of two segments with unit length at the lower extremities of  $D$ .
- For  $i = 0, \dots, 8$ ,  $\Gamma_i$  is a non-optimizable subset of the upper side of  $D$  given by

$$\Gamma_i := \left[ i\frac{10}{9}, (i+1)\frac{10}{9} \right] \times \{2\}, \quad \forall 0 \leq i \leq 8;$$

$\Gamma_i$  is subjected to a unit, vertical downward traction load  $\mathbf{g}_i = (0, -1)$ .

- The remaining region  $\Gamma$  is traction-free and it is the only region of  $\partial\Omega$  which is subject to optimization.

Non-optimizable material layers of width 0.1 are additionally imposed on the upper part of the domain  $D$  and above each component of  $\Gamma_D$ ; see [Figure 3.12](#). We consider nine different load cases, that are

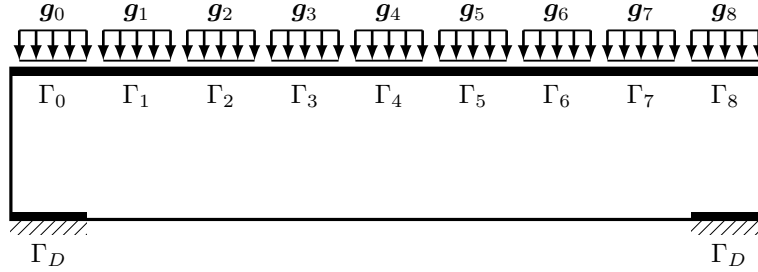


Figure 3.12: Geometric setting for the multiple load case test case

obtained by applying successively and exclusively each of the loads  $\mathbf{g}_i$  on the region  $\Gamma_i$ . In each situation, the corresponding elastic displacement  $\mathbf{u}_i$  is the unique solution in  $H^1(\Omega, \mathbb{R}^d)$  to the linearized elasticity system:

$$\begin{cases} -\operatorname{div}(Ae(\mathbf{u}_i)) = 0 & \text{in } \Omega \\ Ae(\mathbf{u}_i)\mathbf{n} = 0 & \text{on } \Gamma \\ Ae(\mathbf{u}_i)\mathbf{n} = \mathbf{g}_i & \text{on } \Gamma_i \\ Ae(\mathbf{u}_i)\mathbf{n} = 0 & \text{on } \Gamma_j \text{ for } j \neq i \\ \mathbf{u}_i = 0 & \text{on } \Gamma_D, \end{cases} \quad (3.6.11)$$

where we recall that  $e(\mathbf{u}) = (\nabla\mathbf{u} + \nabla\mathbf{u}^T)/2$  is the strain tensor associated to the displacement  $\mathbf{u}$  and  $Ae(\mathbf{u}) = 2\mu e(\mathbf{u}) + \lambda\operatorname{Tr}(e(\mathbf{u}))I$  is the corresponding stress tensor, involving the Hooke's law  $A$ . The Young modulus and the Poisson ratio are set to  $E = 15$  and  $\nu = 0.35$ , which corresponds to  $\lambda = 12.96$  and  $\mu = 5.56$ .

Starting from the initial structure  $\Omega_0$  depicted in [Figure 3.13](#), we are interested in the simultaneous minimization of the volume  $\operatorname{Vol}(\Omega)$  of  $\Omega$  and in the maximization of all the compliances  $C_i(\Omega)$  associated

with each load case  $\mathbf{g}_i$ . These quantities are defined by:

$$\text{Vol}(\Omega) := \int_{\Omega} dx, \quad C_i(\Omega) := \int_{\Omega} A e(\mathbf{u}_i) : e(\mathbf{u}_i) dx, \quad (3.6.12)$$

and their shape derivatives read (see [chapter 2](#)):

$$\text{DVol}(\Omega)(\boldsymbol{\theta}) = \int_{\Gamma} \boldsymbol{\theta} \cdot \mathbf{n} ds, \quad \text{DC}_i(\Omega)(\boldsymbol{\theta}) = - \int_{\Gamma} A e(\mathbf{u}_i) : e(\mathbf{u}_i) \boldsymbol{\theta} \cdot \mathbf{n} ds. \quad (3.6.13)$$

### Volume minimization with maximum compliance constraint

A first possible way to address this case featuring multiple concurrent objectives is to minimize the volume  $\text{Vol}(\Omega)$  subject to an upper bound constraint on the individual compliances  $C_i(\Omega)$ :

$$\begin{aligned} \min_{\Omega \in \mathcal{X}} \quad & \text{Vol}(\Omega) \\ \text{s.t.} \quad & C_i(\Omega) \leq C \quad \text{for all } i \in I \end{aligned} \quad (3.6.14)$$

where  $I \subset \{0, 1, \dots, 8\}$  is a set of indices for the considered load cases. The value of the threshold  $C$  in [\(3.6.14\)](#) is set to a fraction of the maximum of the compliances  $C_i(\Omega_0)$  of the initial design  $\Omega_0$ :

$$C = 0.7 \max_{i=0, \dots, 8} \int_{\Omega_0} A e(\mathbf{u}_i) : e(\mathbf{u}_i) dx \quad (3.6.15)$$

We solve [\(3.6.14\)](#) in the following three configurations:

1. *Case 1: single load case:*  $I = \{4\}$  (only the central load  $\mathbf{g}_4$  is applied)
2. *Case 2: three load cases:*  $I = \{0, 4, 8\}$  (only the central load  $\mathbf{g}_4$  and the two extreme loads  $\mathbf{g}_0$  and  $\mathbf{g}_8$  are applied).
3. *Case 3: all load cases:*  $I = \{0, 1, \dots, 8\}$  (all nine loads are considered).

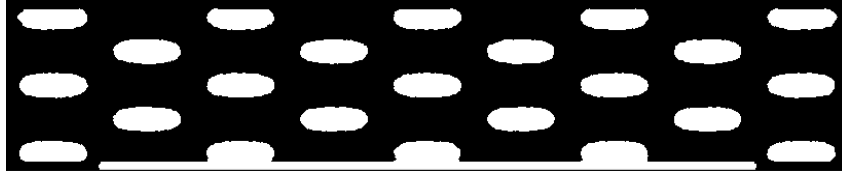


Figure 3.13: Initialisation  $\Omega_0$  (solid in *black*) for the shape optimization examples of [section 3.6](#). The thin white layer at the bottom is a non optimizable part of the domain.

Let us emphasize that for this example (and the next ones), no fine tuning of the algorithm parameters  $A_J$  and  $A_C$  (determining the update of the values of  $\alpha_{J,n}$  and  $\alpha_{C,n}$  in [\(3.6.8\)](#)) of [section 3.6.2](#) is required. The only intuition guiding our choice for this particular test case is that the value of  $A_J$  should be set lower than  $A_C$ . Indeed, a too high value of  $A_J$  might entail a too quick decrease of the volume, which would incur dramatic topological changes violating the rigidity constraints. Therefore these parameters were set to  $A_J = 1$  and  $A_C = 2$  for this test case. The minimum mesh size is  $\text{hmin} = 0.03$ .

The optimized shapes obtained in the three aforementioned situations are shown on [Figure 3.14](#). The meshes of the initial and final designs, as well as several intermediate shapes corresponding to the nine load test-case are shown on [Figs. 3.15](#) and [3.16](#). The convergence histories in the three situations are reported on [Figs. 3.17](#) to [3.19](#). They allow to verify the decrease of the objective function even after the saturation of the constraints. Note that for this example and the one to follow, we observed that  $\hat{I}(\Omega_n)$  coincides with  $\tilde{I}(\Omega_n)$  at every iteration, however this situation is very specific to this test case and does not reproduce in generality<sup>1</sup>. As expected, the optimal value found for the volume of the solid distribution increases with the number of imposed constraints. The major structural change between the different situations is the addition of extra vertical bars of material near the extremities of the structure.

<sup>1</sup>for instance the thermoelasticity test cases of [chapter 2, section 2.5.5](#) for which the volume constraint does not saturate featured  $\hat{I}(\Omega_n) \neq \tilde{I}(\Omega_n)$

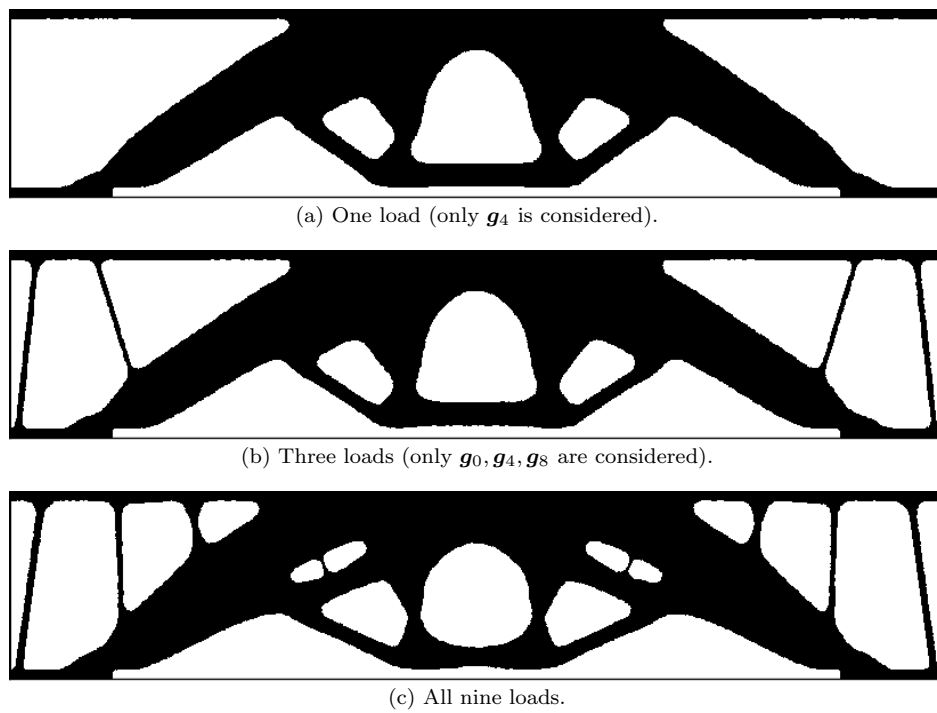


Figure 3.14: Optimized shapes for three possible configurations of the volume minimization problem (3.6.14) subject to maximum compliance constraint.

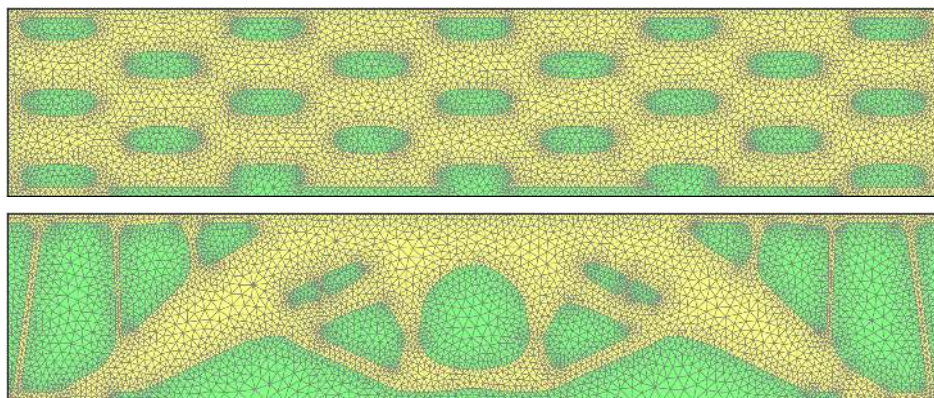


Figure 3.15: Meshes of the initialization and final shapes for the nine load case of Figure 3.14c ((3.6.14)).

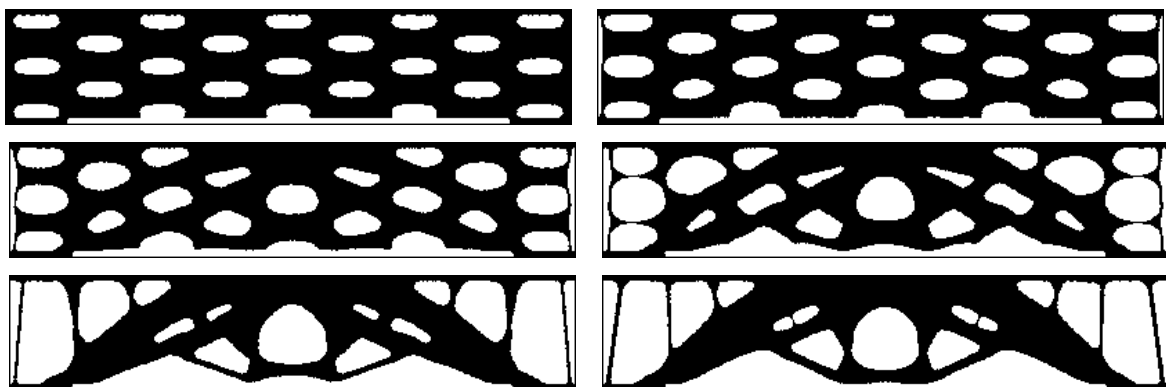


Figure 3.16: Intermediate minimizing shapes for the nine load case of the volume minimization problem of (3.6.14) (iterations 0, 5, 10, 20, 80, and 300).

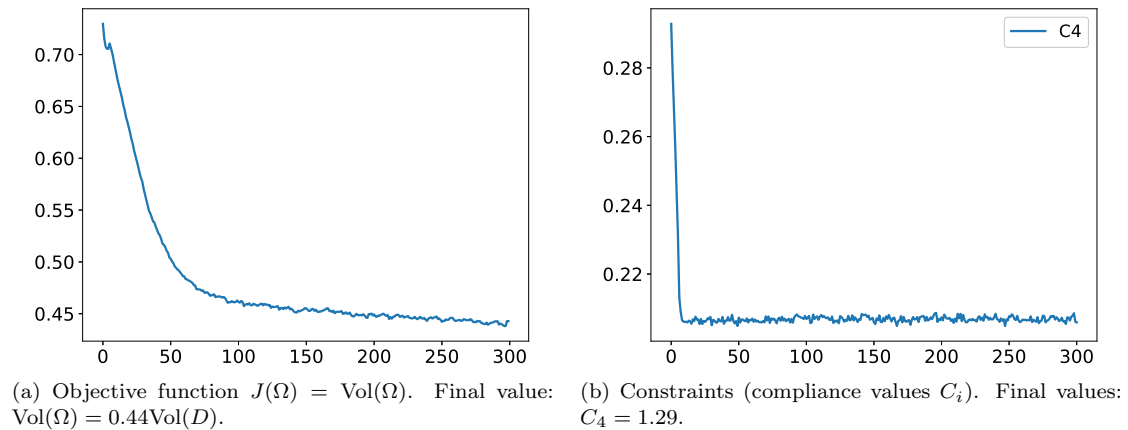


Figure 3.17: Convergence history curves for the single load case of (3.6.14).

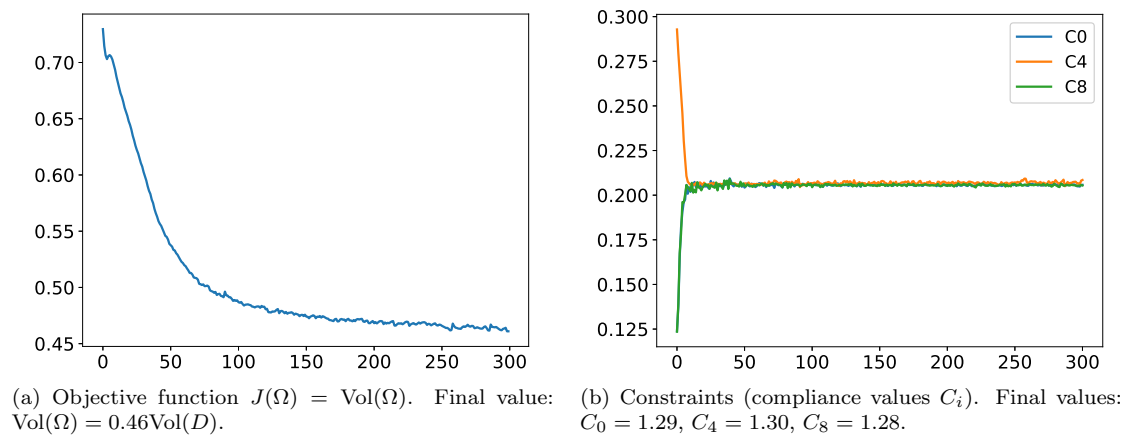


Figure 3.18: Convergence history curves for the three load case of (3.6.14).

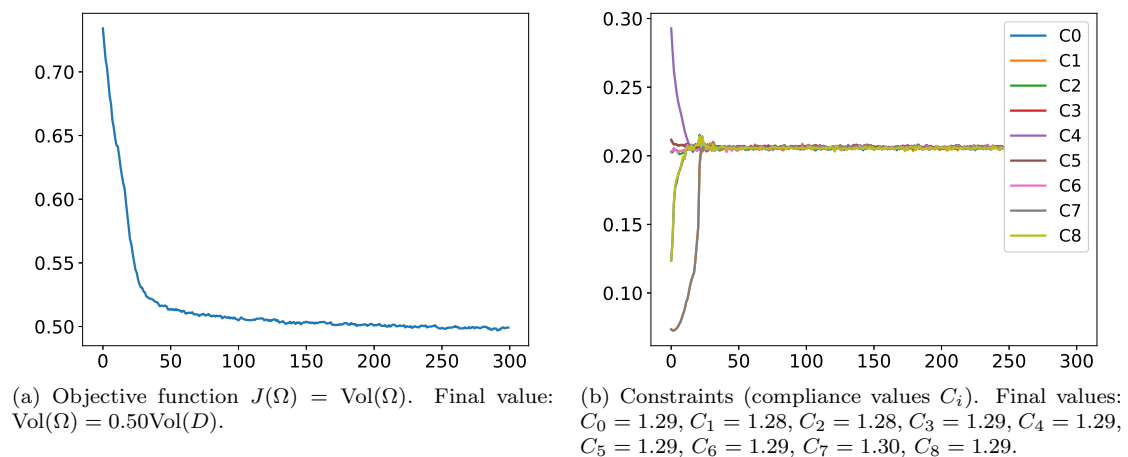


Figure 3.19: Convergence history curves for the nine load case of (3.6.14).



**Min/Max compliance optimization with a volume constraint**

We now consider a different approach where the maximum value of the compliances  $C_i(\Omega)$  is minimized with an equality volume constraint:

$$\begin{aligned} \min_{\Omega \in \mathcal{X}} \quad & \max_{i \in I} C_i(\Omega) \\ \text{s.t.} \quad & \text{Vol}(\Omega) = \rho_0 \text{Vol}(D) \end{aligned} \tag{3.6.16}$$

for a target volume fraction  $\rho_0 = 0.5$  of elastic material and for the three load sets  $I$  introduced in the previous subsection. This problem may be given the form (3.1.2) after introducing a slack variable  $m$ :

$$\begin{aligned} \min_{(\Omega, m) \in \mathcal{X} \times \mathbb{R}} \quad & m \\ \text{s.t.} \quad & \begin{cases} \text{Vol}(\Omega) = \rho_0 \text{Vol}(D) \\ C_i(\Omega) \leq m \quad \text{for all } i \in I. \end{cases} \end{aligned} \tag{3.6.17}$$

The optimization must now be performed with respect to both the slack variable  $m$  and the domain geometry  $\Omega$ , which demands minor adaptations of our optimization algorithm (similar e.g. to those in section 3.5.1): the optimization set  $\mathcal{X} \times \mathbb{R}$  is equipped with the tensorized tangent space  $\tilde{V} = V \times \mathbb{R}$  and differentials are identified to gradients thanks to the inner product  $\langle \cdot, \cdot \rangle_{\tilde{V}}$  defined by

$$\forall (v, w) \in H^1(D, \mathbb{R}) \times H^1(D, \mathbb{R}), (l, m) \in \mathbb{R} \times \mathbb{R}, \quad \langle (v, l), (w, m) \rangle_{\tilde{V}} := \langle v, w \rangle_V + lm, \tag{3.6.18}$$

where  $\langle \cdot, \cdot \rangle_V$  is the scalar product of (3.6.5) or (3.6.6). The slack variable  $m$  is initialized with the maximum value of the compliance of the initial structure  $\Omega_0$  over all the considered loads:

$$m_0 := \max_{i \in I} C_i(\Omega_0), \tag{3.6.19}$$

and its values  $m_n$  are then updated along with the shape  $\Omega_n$  according to algorithm 3.1.

The resulting optimized structures are shown on Figure 3.20 for each of the three considered configurations and the associated convergence histories are displayed on Figs. 3.21 to 3.23 for the single, triple and nine load cases respectively. Note that sudden, abrupt peaks on the constraint curves correspond to topological changes (e.g. at iteration 38 for the nine load case) for which the displacements corresponding to the extremal loads  $\mathbf{g}_0$  and  $\mathbf{g}_8$  are especially sensitive. The decrease of all compliance functions is observed even after all the inequality constraints are saturated, which occurs as soon as where all compliances achieve a common value. As expected, the optimal design found for the nine load minimum compliance case (Figure 3.14c) is similar (up to a few bars) to the corresponding one found for the volume minimization (Figure 3.20c): indeed, both cases reach at convergence a volume fraction  $\text{Vol}(\Omega) = 0.5 \text{Vol}(D)$  and a maximum compliance  $\max C_i(\Omega) \simeq 1.30$ .

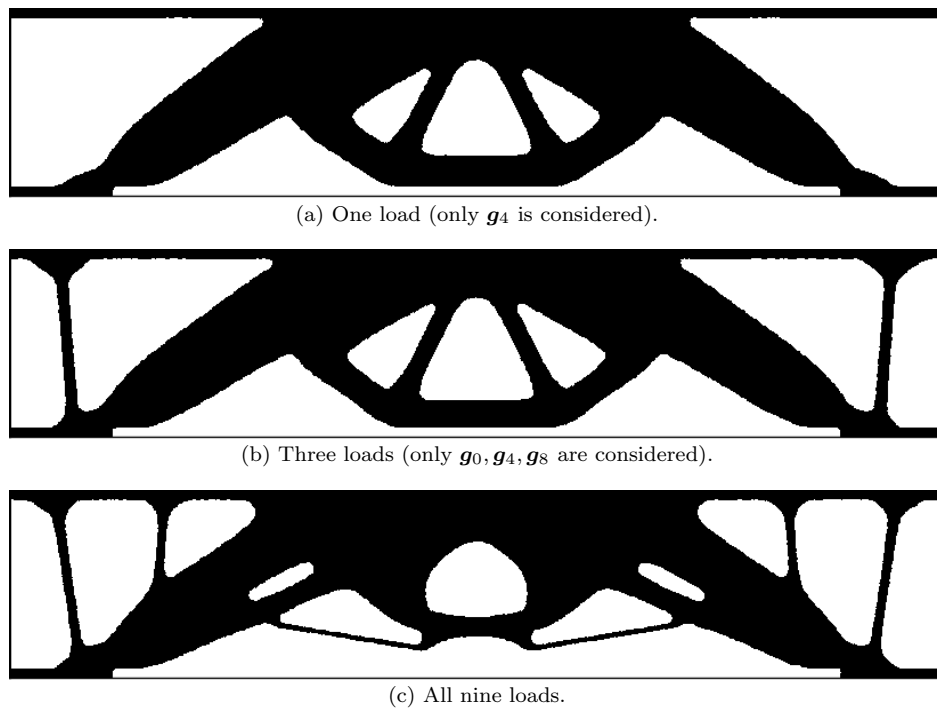


Figure 3.20: Optimized shapes for three possible configurations of the min/max optimization problem (3.6.17).

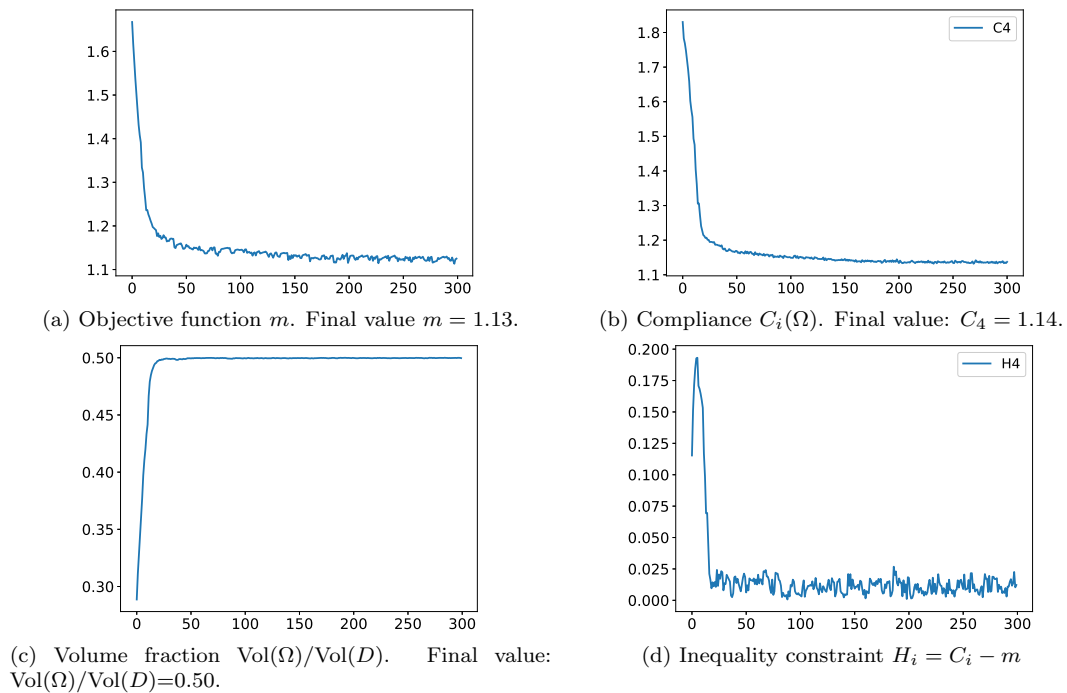


Figure 3.21: Convergence history curves for one load case of (3.6.17).

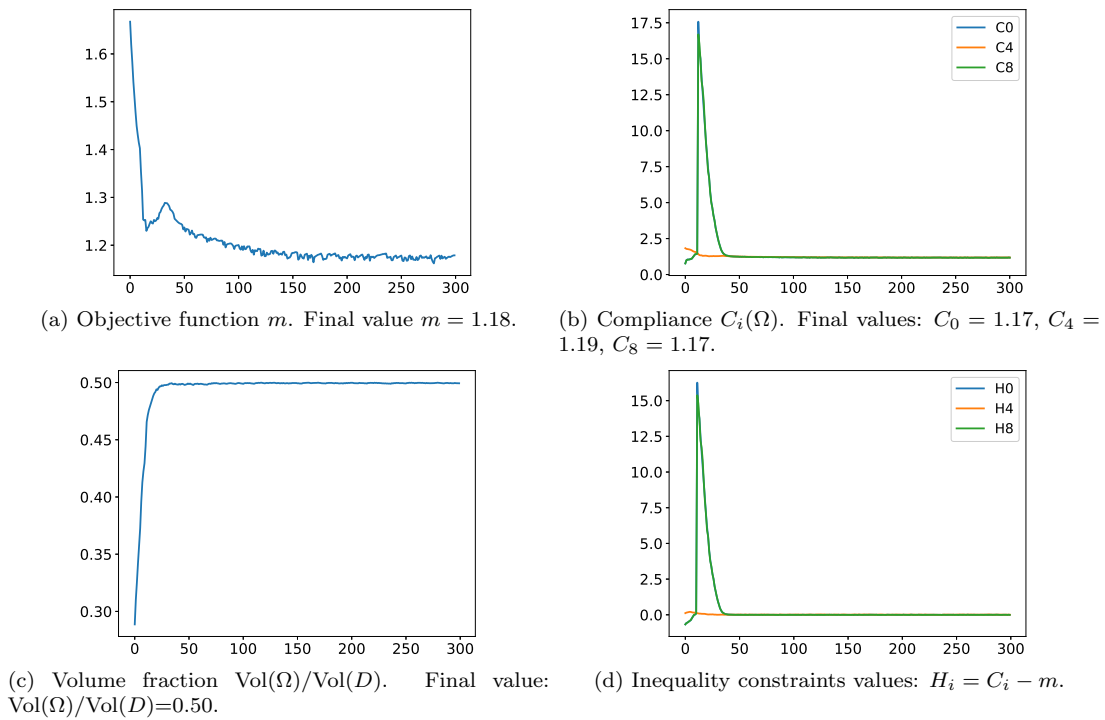


Figure 3.22: Convergence history curves for three load case of (3.6.17).

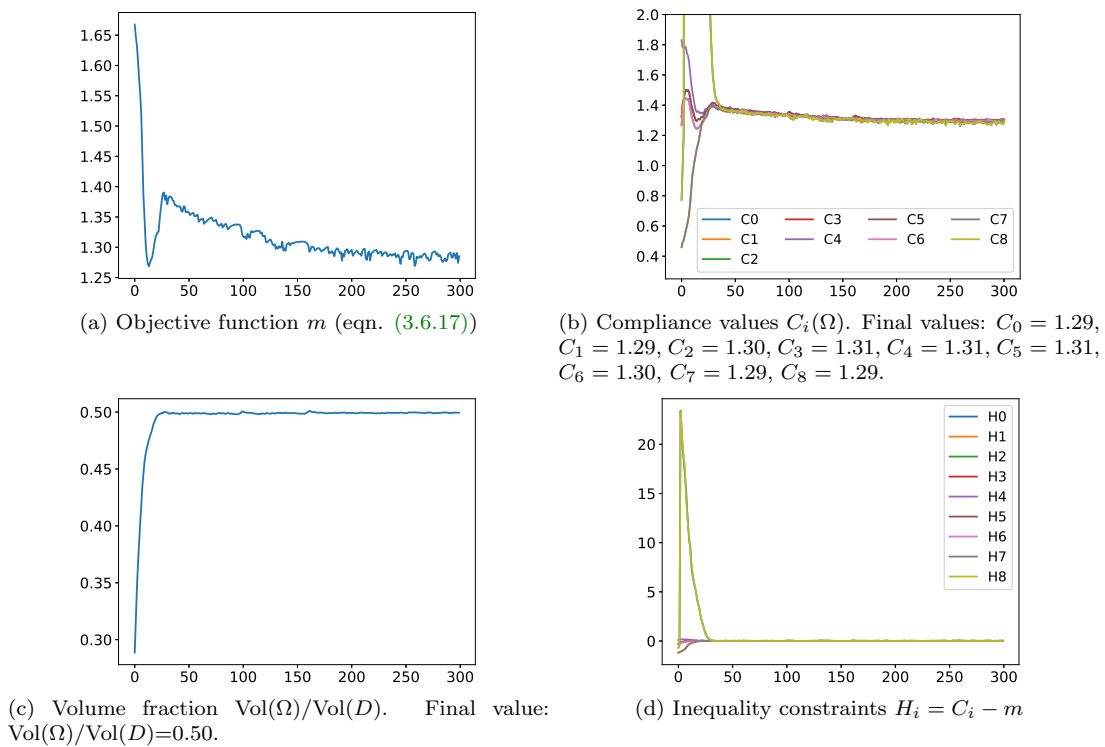


Figure 3.23: Convergence history curves for nine load case of (3.6.17).



## CHAPTER 4

# A VARIATIONAL METHOD FOR COMPUTING SHAPE DERIVATIVES OF GEOMETRIC CONSTRAINTS ALONG RAYS

### Contents

---

<b>4.1 Introduction</b>	<b>149</b>
<b>4.2 Weighted graph space of the advection operator <math>\beta \cdot \nabla</math> for velocity fields of class <math>\mathcal{C}^1</math></b>	<b>152</b>
4.2.1 Preliminaries, notation and assumptions	153
4.2.2 Definition of the graph space $V_\omega$ of the advection operator $\beta \cdot \nabla$	156
4.2.3 Density of functions of class $\mathcal{C}^1$ in the weighted space $V_\omega$	157
4.2.4 Trace theorem and Poincaré inequality in $V_\omega$	160
4.2.5 Well-posedness of the variational problem (4.2.5)	163
<b>4.3 Numerical methods for integration along normal rays</b>	<b>163</b>
4.3.1 Shape optimization context: normal rays and flow map of the signed distance function gradient	164
4.3.2 Computing curvatures and detecting the skeleton for direct integration along the rays	164
4.3.3 Admissible numerical weights built upon the signed distance function	166
4.3.4 Numerical comparisons between the variational method and direct integration along rays	170
<b>4.4 Applications to maximum and minimum thickness constraints in shape optimization</b>	<b>178</b>
4.4.1 Shape optimization setting for linearly elastic structures	178
4.4.2 Shape optimization under a maximum thickness constraint	179
4.4.3 Shape optimization examples under a minimum thickness constraint	181

---

*Note* : most of the content of this chapter has been published in [154]. F. FEPPON, G. ALLAIRE, AND C. DAPOGNY, *A variational formulation for computing shape derivatives of geometric constraints along rays*, To appear in M2AN, (2019).

The introduction has been partly modified in regards with this new section and the original motivations of the thesis. Some background material redundant with the previous chapters has also been removed.

### 4.1 INTRODUCTION

One of the key challenges regarding the systematic integration of topology optimization methods in industrial applications is the need for taking into account geometric constraints. These come very often from the final manufacturing process which generally brings into play constraints on the geometry of shapes; for instance minimum thickness constraints, maximum thickness constraints, minimum distance between members, casting constraints; see [320, 88, 178, 234] or the recent survey [220].

For our multiphysics shape optimization purposes, the starting motivation for the present chapter originates from the need for tackling non penetration constraints in liquid liquid heat exchangers design. In the next [chapter 6 \(section 5.1\)](#), we consider the heat exchanger design problem of [Figure 4.1](#) featuring two fluid phases  $\Omega_{f,hot} \subset D$  and  $\Omega_{f,cold} \subset D$ . A very natural specification of the problem is that these two phases should not interpenetrate during the optimization process of the union  $\Omega = \Omega_{f,cold} \cup \Omega_{f,hot}$  of the two phases. Very few works have investigated the enforcement of this non-mixing condition; we are actually only aware of the thesis of Papazoglou [255] in the context of density based topology optimization.

Among the variety of methods proposed in the literature to enforce geometric constraints, several works [29, 30, 109, 110, 234] have proposed to formulate them by means of integral functionals involving the signed distance function  $d_\Omega$  (see [chapter 1, section 1.3](#)) to the optimized domain  $\Omega$ . More precisely, the shape  $\Omega$  is sought among all possible subsets of a fixed ‘hold-all’ domain  $D \subset \mathbb{R}^d$  (with still  $d = 2, 3$  in applications) as the solution of a constrained minimization problem of the form:

$$\min_{\Omega \subset D} J(\Omega), \text{ s.t. } P(\Omega) \leq 0, \text{ where } P(\Omega) := \int_D j(d_\Omega(x)) dx, \quad (4.1.1)$$

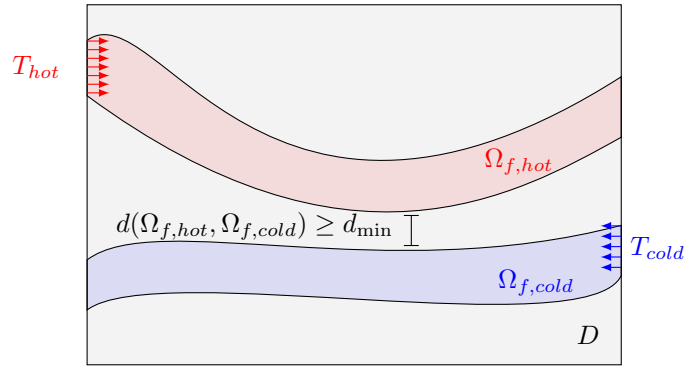


Figure 4.1: Setting for a 2-d fluid fluid heat exchanger design problem (treated in [chapter 6, section 5.1](#)). A hot fluid phase  $\Omega_{f,1} \subset D$  is entering at the top left with a temperature  $T_{hot}$ , and a hot fluid phase  $\Omega_{f,2}$  is entering in the reverse direction at the bottom right inlet. Both phases should remain at a distance  $d(\Omega_{f,1}, \Omega_{f,2}) > d_{\min}$  from one another.

where  $J(\Omega)$  is the objective function,  $j : \mathbb{R} \rightarrow \mathbb{R}$  is a given, smooth function, and  $P(\Omega)$  is a geometric constraint. The above framework is quite appealing insofar as the signed distance function  $d_\Omega$  naturally lends itself to clear mathematical formulations of maximum, minimum thickness constraint functionals (and the other aforementioned geometric criteria). For instance, the non penetration constraint of the 2-d heat exchanger problem of [Figure 4.1](#) can be conveniently formulated with the signed distance function of either of the phases  $\Omega_{f,hot}$  or  $\Omega_{f,cold}$ , because it is sufficient to require the two phases to remain at a distance  $d_{\min} > 0$  from one another:

$$\forall x \in \Omega_{f,hot}, d_{\Omega_{f,cold}}(x) \geq d_{\min}.$$

This pointwise constraint can then be captured by an integral criterion  $P(\Omega_{f,hot})$  as in [\(4.1.1\)](#), see [section 5.1](#).

Allaire, Jouve and Michailidis [\[30\]](#) have shown that the problem [\(4.1.1\)](#) is amenable to numerical treatments by means of standard (e.g. steepest-descent) optimization algorithms. However, several technical stages in its numerical implementation pose difficulties. The leading motivation of the present chapter is to introduce a new variational method that make geometric constraints substantially simpler to implement in the framework of the method of Hadamard.

Let us provide a little more details about the main numerical issues raised by the implementation of [\(4.1.1\)](#), while staying at an explanatory level; see [sections 4.3.1](#) and [4.4](#) below for full details. The use of traditional optimization methods for the program [\(4.1.1\)](#) relies (in particular) on the knowledge of the *shape derivative* of the considered constraint functional  $\Omega \mapsto P(\Omega)$  ([chapter 1, definition 1.1](#)). It was shown in [\[23\]](#) that the shape derivative of  $P(\Omega)$  as in [\(4.1.1\)](#) has the form

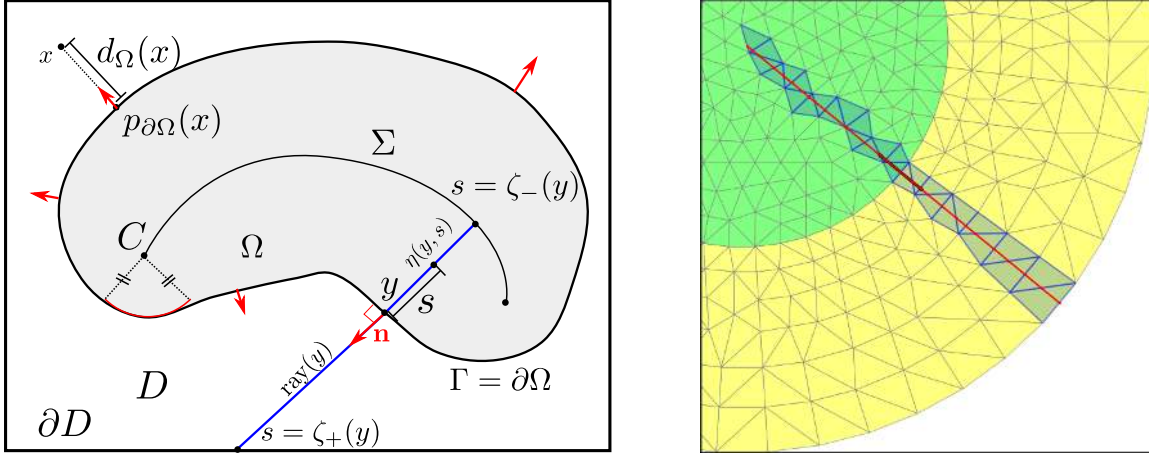
$$DP(\Omega)(\boldsymbol{\theta}) = \int_{\partial\Omega} u \boldsymbol{\theta} \cdot \mathbf{n} \, dy, \quad (4.1.2)$$

where  $dy$  refers to the surface measure on  $\partial\Omega$ , and the function  $u \in L^1(\partial\Omega)$  is given by the formula

$$u(y) = - \int_{x \in \text{ray}(y)} j'(d_\Omega(x)) \prod_{1 \leq i \leq n-1} (1 + \kappa_i(y) d_\Omega(x)) dx, \quad \forall y \in \partial\Omega. \quad (4.1.3)$$

In [\(4.1.2\)](#) and [\(4.1.3\)](#),  $\mathbf{n}$  and  $\kappa_i(y)$  stand for the unit normal vector to  $\partial\Omega$  pointing outward  $\Omega$  and the  $(n-1)$  principal curvatures of  $\partial\Omega$  respectively; the set  $\text{ray}(y)$  is the normal ray that originates from  $y \in \partial\Omega$  and that stops at either the boundary  $\partial D$  of the hold-all domain, or at the skeleton (or medial axis)  $\Sigma$  of  $\Omega$ ; see [Figure 4.2a](#) for an illustration. The normal velocity  $\boldsymbol{\theta} = -u\mathbf{n}$  is then exploited by optimization algorithms (such as our null space method of [chapter 3, algorithm 3.1](#)) as a constraint gradient to build a descent update for the shape  $\Omega$ . The numerical component of this program raises the need to calculate the quantity  $u : \partial\Omega \rightarrow \mathbb{R}$  defined in [\(4.1.3\)](#).

In principle, the formula [\(4.1.3\)](#) featuring integration along the normal rays to  $\partial\Omega$  can be implemented efficiently. In particular, its numerical evaluations at several points  $y \in \partial\Omega$  can be performed in parallel.



(a) The shape optimization setting: shape  $\Omega \subset D$ , skeleton  $\Sigma$  with normal rays  $(\text{ray}(y))_{y \in \partial\Omega}$ , local coordinate change  $\eta$ , outward normal vector  $\mathbf{n}$ , center of curvature  $C$ , orthogonal projection  $p_{\partial\Omega}(x)$  of a point  $x \in D$  onto  $\partial\Omega$ . (see [chapter 1, section 1.3](#) for the definitions). (b) Integration along rays requires the detection of the skeleton of the shape and the computation of the triangle angle path at the mesh level, that is the intersection points of each triangle with the ray.

Figure 4.2: Definition of geometric quantities and integration along normal rays.

However, the practical implementation when all the functions involved in the integrand of (4.1.3) are discretized on a computational mesh is not trivial: it requires the computation of (i) the ray trajectories on the discrete mesh (for instance by some variant of the method of characteristics), as exemplified on [Figure 4.2b](#), and (ii) the principal curvatures  $\kappa_i$  of the boundary  $\partial\Omega$  the numerical discretization of  $\Omega$ . As we shall illustrate in our numerical investigations, the accurate calculation of these quantities—even when  $\partial\Omega$  is endowed with a discretization as an explicit triangulated surface—is not trivial and is still an active research direction, already in the case of two space dimensions, and especially in the three-dimensional context [[233](#), [270](#), [143](#)].

The objective of this chapter is to introduce a robust and easy to implement method that allows to compute integral quantities such as (4.1.3) by solving a variational problem for advection-like operators which alleviates the need for resorting to direct integration along rays. This yields an efficient numerical method to calculate shape derivatives of geometric constraints, which is very simple to implement using a standard finite element software, and which does not raise any additional algorithmic difficulty in 3-d than in 2-d.

Our numerical method for the calculation of (4.1.3) is based on the following variational problem, set over a suitable Hilbert space  $V_\omega$ :

$$\text{Find } u \in V_\omega \text{ such that } \forall v \in V_\omega, \int_{\partial\Omega} uv ds + \int_{D \setminus \bar{\Sigma}} \omega(\nabla d_\Omega \cdot \nabla u)(\nabla d_\Omega \cdot \nabla v) dx = - \int_{D \setminus \bar{\Sigma}} j'(d_\Omega) v dx, \quad (4.1.4)$$

where  $\omega > 0$  is a rather arbitrary weight function (over which relevant assumptions shall be stated later on), and  $\bar{\Sigma}$  is the closure of the skeleton set  $\Sigma$  of  $\Omega$  (see [Figure 4.2a](#) and [section 4.3.1](#) below). Indeed, our key result is to show that (4.1.4) is well-posed, and that the trace of its unique solution  $u$  on  $\partial\Omega$  is exactly (4.1.3).

A formal insight that allows to see this is the following: assume that for any  $v_0 \in C^0(\partial\Omega)$ , the function  $v$  satisfying  $v = v_0$  on  $\partial\Omega$  and taking constant values along the rays normal to  $\partial\Omega$  (i.e.  $\nabla d_\Omega \cdot \nabla v = 0$  in  $D \setminus \bar{\Sigma}$ ) belongs to the space  $V_\omega$ . Using  $v$  as a test function in the variational formulation (4.1.4) then yields:

$$\begin{aligned} \int_{\partial\Omega} uv_0 ds &= \int_{\partial\Omega} uv ds + \int_{D \setminus \bar{\Sigma}} \omega(\nabla d_\Omega \cdot \nabla u)(\nabla d_\Omega \cdot \nabla v) dx \\ &= - \int_{D \setminus \bar{\Sigma}} j'(d_\Omega) v dx = - \int_{y \in \partial\Omega} \left( \int_{z \in \text{ray}(y)} j'(d_\Omega(z)) \prod_{1 \leq i \leq n-1} (1 + \kappa_i(y) d_\Omega(z)) ds \right) v_0(y) dy, \end{aligned} \quad (4.1.5)$$

where the last equality follows from a change of variables allowing to rewrite integration on  $D \setminus \bar{\Sigma}$  as a nested integration on points  $y \in \partial\Omega$  and on elements  $z$  in  $\text{ray}(y)$  (see (4.2.7) and (4.2.9) below). Since

$v_0 \in C^0(\partial\Omega)$  can be chosen arbitrarily in (4.1.5), the identity (4.1.3) follows. These considerations are made rigorous in section 4.2.5.

The variational formulation (4.1.4) makes it possible to compute integrals (4.1.3) along the normal rays to  $\partial\Omega$  without the need to calculate these rays or the curvatures  $\kappa_i$  explicitly on a discretization (i.e. a mesh) of the ambient space. This feature is especially convenient for shape optimization applications relying on the level set method, as described in section 4.4; there, the gradient of the signed distance function  $\nabla d_\Omega$  is easy to calculate on an unstructured mesh of the considered hold-all domain  $D$  from a  $\mathbb{P}_1$  approximation of  $d_\Omega$ . The variational formulation (4.1.4) can then be readily implemented in a finite element framework, even if the boundary  $\Gamma = \partial\Omega$  is not meshed explicitly. Our method requires only a rough estimate of the location of the skeleton  $\Sigma$ , upon a judicious choice of the weight  $\omega$  in (4.1.4); see the numerical examples in section 4.3. Let us emphasize that the variational approach (4.2.5) is equally simple to implement in any space dimension while the 3-d implementation of geometric integration along rays would require much more efforts than in 2-d.

Last, the previous arguments work identically when  $\nabla d_\Omega$  is replaced with an arbitrary  $C^1$  vector field  $\beta$ : we obtain that a variational formulation analogous to (4.1.4) (given in (4.2.5)) allows to compute integrals quantities along the characteristics curves of  $\beta$ , which is subject to offer wider applications than shape optimization.

With these perspectives in mind, the present chapter is organized as follows. Section 4.2 carefully discusses the mathematical setting that guarantees the existence and uniqueness of a solution to the variational problem (4.1.4) (in fact its generalization to arbitrary vector fields  $\beta$ ), and the justification of the key identity (4.1.3) satisfied by the trace of its solution. For this purpose, we provide a variational theory for the advection operator  $\beta \cdot \nabla$  associated to arbitrary  $C^1$  vector fields  $\beta$  on the weighted graph space  $V_\omega$ ; in particular, the existence of traces on  $\Gamma$  and an adapted Poincaré inequality for functions  $v \in V_\omega$  are obtained. Note that many related works are available about these matters (e.g. [131, 144]), however they usually rely on strong boundedness assumptions on the divergence of  $\beta$  (typically  $\operatorname{div}(\beta) \in L^\infty(D)$ ), which do not hold for our shape optimization applications (4.2.8) and (4.2.9). Indeed, in the latter situation where  $\beta = \nabla d_\Omega$ , the divergence  $\operatorname{div}(\nabla d_\Omega)$  typically blows up near  $\Sigma$ ; hence the need for our different approach. Our approach requires rather the existence of a flow  $\eta$  (see (4.2.1) below) associated with the vector field  $\beta$ .

In section 4.3, we investigate the numerical accuracy of our variational method for calculating integrals along characteristic curves in the shape optimization setting (4.2.8) where  $\beta = \nabla d_\Omega$ . After a short review of the properties of the signed distance function, we compare the direct numerical integration along rays with the use of our variational method on several numerical examples where the value of (4.1.3) is analytically known. We also consider “practical” cases where some of the regularity assumptions imposed by our framework are not fulfilled. Most importantly, we illustrate how the selection of an appropriate weight  $\omega$  in the variational formulation (4.1.4) allows to deal with the presence of cracks in the working domain, for instance the skeleton  $\Sigma$  in shape optimization when it is not explicitly meshed.

The last two sections are dedicated on numerical applications. Section 4.4 illustrates the simplicity and effectiveness of our method on practical shape optimization applications. Section 4.4 elaborates on the works [25, 30, 110] concerned with manufacturing constraints: we demonstrate that our variational method allows for a convenient and efficient implementation of maximum and minimum thickness constraints in structural design. Finally, section 5.1 applies the method in order to enforce a non-mixing constraint for liquid-liquid heat exchanger problems of such as the one depicted in the above Figure 4.1.

## 4.2 WEIGHTED GRAPH SPACE OF THE ADVECTION OPERATOR $\beta \cdot \nabla$ FOR VELOCITY FIELDS OF CLASS $C^1$

This section is concerned with the mathematical analysis of a slightly more general variational problem than (4.1.4):  $\nabla d_\Omega$  is replaced by a rather arbitrary vector field  $\beta$ . This setting and some technical assumptions are described in section 4.2.1. In order to obtain the well-posedness of the corresponding variational problem and the trace identity (4.1.3), suitable functional spaces  $V_\omega$  are introduced in section 4.2.2 in which the directional derivative  $\beta \cdot \nabla$  naturally makes sense.

Then, section 4.2.3 investigates the density of  $C^1$  functions in the graph space  $V_\omega$ . In section 4.2.4, we establish a trace theorem for functions in  $V_\omega$  and we provide a Poincaré-type inequality. These two ingredients allow us to prove in section 4.2.5 the well posedness of the variational problem (4.1.4) and of the identity in this more general context involving arbitrary fields  $\beta$ .



### 4.2.1 Preliminaries, notation and assumptions

#### Preliminaries and notation

Let  $U \subset \mathbb{R}^d$ ,  $\Gamma \subset \bar{U}$  and  $\beta : U \rightarrow \mathbb{R}^d$  be respectively a (possibly non smooth) bounded open set, a hypersurface and a vector field of class  $\mathcal{C}^1$ . Note we do not require  $\Gamma$  to be a compact manifold ( $\Gamma$  may differ from its closure  $\bar{\Gamma}$ ); in this case, we require  $\bar{\Gamma}$  to be a manifold with boundary, which in particular prevents  $\Gamma$  from showing spiralling patterns near its ends—an assumption which is needed for technical reasons (see e.g. the proof of [lemma 4.3](#) below). Two examples of admissible settings are represented on [Figure 4.3](#).

We assume that  $\Gamma$  is a Poincaré section or a *stream surface* for  $\beta$ , meaning that for any  $y \in \Gamma$ , there is a unique characteristic curve  $s \mapsto \eta(y, s)$  passing through  $y = \eta(y, 0)$  at time  $s = 0$ , and that lives in the domain  $U$  on some maximal interval  $s \in (\zeta_-(y), \zeta_+(y))$ . In other words, for any  $y \in \Gamma$ ,  $(\zeta_-(y), \zeta_+(y))$  is the maximum existence interval such that the solution  $s \mapsto \eta(y, s)$  of the ordinary differential equation

$$\begin{cases} \frac{d}{ds} \eta(y, s) = \beta(\eta(y, s)), \\ \eta(y, 0) = y, \end{cases} \quad (4.2.1)$$

remains in the domain  $U$  (note that by definition,  $\zeta_- \leq 0 \leq \zeta_+$ ). We assume that  $\zeta_-$  and  $\zeta_+$  are continuous, bounded functions on  $\Gamma$ , satisfying the following separation condition:

$$\exists \varepsilon > 0, \forall y \in \Gamma, \zeta_+(y) - \zeta_-(y) > \varepsilon. \quad (4.2.2)$$

The vector field  $\beta$  is required to be  $U$ -filling, in the sense that its related flow  $\eta$  realizes a  $\mathcal{C}^1$  diffeomorphism from the tensor product set

$$W = \{(y, s) \mid y \in \Gamma, s \in (\zeta_-(y), \zeta_+(y))\}, \quad (4.2.3)$$

onto  $U$  (see [\[212\]](#), Chap. IV), where  $W$  is a manifold obtained as an open subset of the tensor product of  $\Gamma$  with the real line  $\mathbb{R}$  (note that the “open” character of  $W$  comes from the continuity assumption on  $\zeta_-$  and  $\zeta_+$ ).

Finally, we denote by  $|D\eta|$  the Jacobian of the local coordinate change  $\eta$ :

$$\forall y \in \Gamma, \forall s \in (\zeta_-(y), \zeta_+(y)), |D\eta|(y, s) = |\det(\nabla\eta)|(y, s), \quad (4.2.4)$$

where the Jacobian matrix of  $\eta$  reads  $\nabla\eta = \begin{bmatrix} \partial_y \eta & \partial_s \eta \end{bmatrix}$  and  $\partial_y$  denotes the collection of derivatives with respect to the  $(n-1)$  tangential coordinates of  $\Gamma$ .

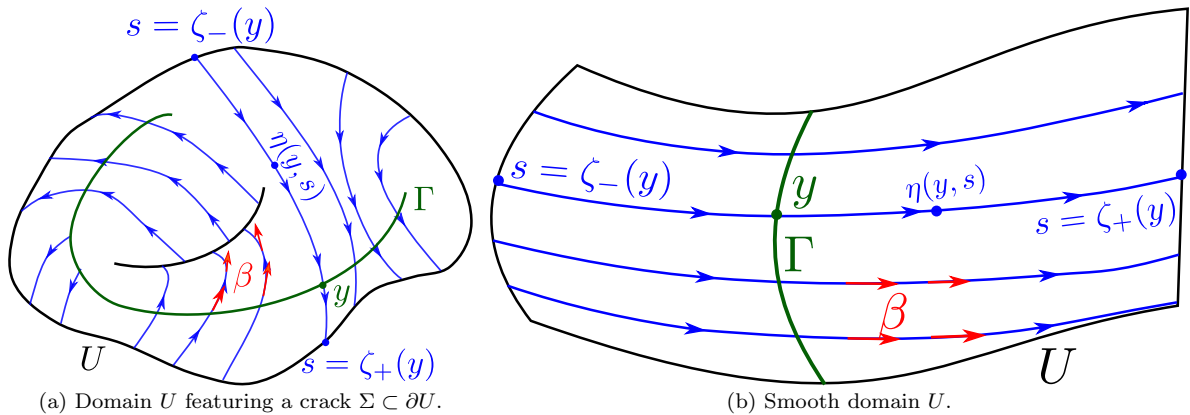


Figure 4.3: Two types of admissible domain  $U$ , with  $U$ -filling vector field  $\beta$  and stream surface  $\Gamma$ , in the framework of [section 4.2.1](#).

**Remark 4.1.** Let us consider a particular case where  $U$  is a smooth bounded domain in  $\mathbb{R}^d$ , and  $\Gamma \subset \partial U$  is defined as the inlet boundary of the flow field  $\beta$ , *i.e.*

$$\Gamma := \{x \in \partial U \mid \beta(x) \cdot \mathbf{n}(x) < 0\},$$

where  $\mathbf{n}$  is the unit normal vector to  $\partial U$ , pointing outward  $U$ , the separation condition (4.2.2) exactly requires that the inflow and outflow boundaries be separated by a positive minimum distance, which is a rather classical assumption in the study of the advection operator, see e.g. Section 2.1.3 in [125]. In our case, assumption (4.2.2) is required in the proof of proposition 4.2.

The main point of this section is to mathematically justify that given rather arbitrary function  $f : U \rightarrow \mathbb{R}$  and weight  $\omega \in \mathcal{C}^0(U, \mathbb{R}_+^*)$ , the trace of the solution  $u$  to the variational problem

$$\text{Find } u \in V_\omega \text{ such that } \forall v \in V_\omega, \int_\Gamma uv ds + \int_U \omega(\beta \cdot \nabla u)(\beta \cdot \nabla v) dx = \int_U f v dx \quad (4.2.5)$$

is given explicitly in terms of line integrals by the formula

$$\forall y \in \Gamma, \quad u(y) = \int_{\zeta_-(y)}^{\zeta_+(y)} (f \circ \eta |D\eta|)(y, s) ds, \quad (4.2.6)$$

which yields a numerical method for computing (4.2.6) by solving (4.2.5).

Throughout the chapter, the considered measure on  $W$  is that induced by the standard product measure of the surface measure  $dy$  on  $\Gamma$  and of the Lebesgue measure  $ds$  on  $\mathbb{R}$ . Thus, the space  $L^1(W)$  of integrable functions on  $W$  is defined as the space of measurable functions  $f : W \rightarrow \mathbb{R}$  such that  $\int_\Gamma \int_{\zeta_-(y)}^{\zeta_+(y)} |f(y, s)| ds dy < +\infty$ . Integration of  $f$  over  $W$  is then defined by:

$$\int_W f(y, s) ds dy := \int_\Gamma \int_{\zeta_-(y)}^{\zeta_+(y)} f(y, s) ds dy.$$

Under these notations, the classical change of variable formula between manifolds reads (see [212], Chap. XVI):

$$\forall f \in L^1(U), \quad \int_U f dx = \int_W f \circ \eta |D\eta| ds dy, \quad (4.2.7)$$

where the Jacobian  $|D\eta|$  is defined by (4.2.4).

### The shape optimization setting as a particular case

The shape optimization setting outlined in the introduction, which is exemplified on Figure 4.2a, reduces to the particular case

$$U = D \setminus \bar{\Sigma}, \quad \Gamma = \partial\Omega, \quad \beta = \nabla d_\Omega. \quad (4.2.8)$$

where  $\Omega \subset D$  are bounded Lipschitz domains of  $\mathbb{R}^d$  (note that in order for  $\beta = \nabla d_\Omega$  to be a  $\mathcal{C}^1$  vector field,  $\Omega$  is assumed to be in fact a  $\mathcal{C}^2$  domain in all what follows).

The bound functions  $\zeta_-(y)$  and  $\zeta_+(y)$  are the distances at which ray( $y$ ) hits either the skeleton  $\Sigma$  of  $\Omega$  or the boundary  $\partial D$  of the hold-all domain  $D$ .

In this situation, the local coordinate change  $\eta$  and its Jacobian  $|D\eta|$  are explicitly given by (see section 4.3.1)

$$\eta(y, s) = y + s \nabla d_\Omega(y), \quad |D\eta|(y, s) = \prod_{i=1}^{d-1} (1 + \kappa_i(y)s), \quad \forall y \in \partial\Omega, \forall s \in (\zeta_-(y), \zeta_+(y)), \quad (4.2.9)$$

so that the function  $u$  of (4.1.3) coincides with the expression (4.2.6) for  $f = -j'(d_\Omega)$ . Let us emphasize that, in this context, the open set  $U$  is not smooth because it features a “crack”, namely the skeleton  $\Sigma$  (we call it a “cracked domain” in Figure 4.2a); in particular,  $U$  is not located on one side of its boundary. This “lack of smoothness” of  $U$  prevents from using many convenient results from functional analysis [302], and thus raises the need for several technical ingredients in the sequel, which otherwise would have been fairly classical.

### Assumptions on weight functions $\omega$

Let us now consider a weight function  $\omega : U \rightarrow \mathbb{R}$ , which will be one of the key ingredients of our variational formulation (4.2.5), satisfying the following assumptions:

**(H1)**  $\omega \in \mathcal{C}^0(U, \mathbb{R}_+^*)$  is a positive, continuous function on  $U$ .

**(H2)** The transported weight  $\alpha := \omega \circ \eta |D\eta| \in \mathcal{C}^0(W, \mathbb{R}_+^*)$  over  $W$  is such that the function

$$\begin{aligned} g_\alpha &: \Gamma \longrightarrow \mathbb{R}_+ \cup \{+\infty\} \\ y &\longmapsto \int_{\zeta_-(y)}^{\zeta_+(y)} \alpha(y, s) ds \end{aligned} \quad (4.2.10)$$

is uniformly bounded, i.e.  $g_\alpha \in L^\infty(\Gamma)$ .

**(H3)** The transported weight  $\alpha = \omega \circ \eta |D\eta|$  over  $W$  is such that the function

$$\begin{aligned} h_\alpha &: \Gamma \longrightarrow \mathbb{R}_+ \cup \{+\infty\} \\ y &\longmapsto \int_{\zeta_-(y)}^{\zeta_+(y)} \left[ \alpha(y, s) \int_0^s \alpha^{-1}(y, t) dt \right] ds \end{aligned} \quad (4.2.11)$$

is uniformly bounded:  $h_\alpha \in L^\infty(\Gamma)$ .

Note that, above and in the sequel, the notations  $\alpha^{-1}(y, s) = 1/\alpha(y, s)$  and  $\omega^{-1}(y, s) = 1/\omega(y, s)$  are used for the inverses of the scalar weights  $\alpha$  and  $\omega$  respectively, whereas the notation  $\eta^{-1} : U \mapsto W$  shall stand for the reciprocal mapping of the diffeomorphism  $\eta$ .

Hypothesis **(H1)** is essentially a regularity assumption, positivity being no surprise for  $\omega$  to be an admissible weight in the variational formulation (4.2.5). Notice that no assumption is made about the behavior of  $\omega$  near the boundary of  $U$ ; in particular,  $\omega(x)$  may tend to 0 as  $x$  approaches  $\partial U$ . Hypothesis **(H2)** is an upper-boundedness assumption for  $\alpha$ . The weights that we are going to consider in our practical applications in section 4.3.3 will satisfy **(H1)** and **(H2)** almost automatically. Finally, **(H3)** is roughly a monotonicity constraint for the decay of  $\alpha$  towards 0 as  $s \rightarrow \zeta_\pm(y)$ . In practice, we will rely on the following lemma which provides a simple monotonicity condition under which the condition **(H3)** is fulfilled, indicating that the class of weights satisfying **(H3)** is large enough.

**Lemma 4.1.** *Let  $\alpha \in \mathcal{C}^0(W, \mathbb{R}_+^*)$  be a weight of the product form:*

$$\forall (y, s) \in W, \alpha(y, s) = f(y, s)g(y, s) \quad (4.2.12)$$

where  $f$  and  $g$  are two positive functions on  $W$  satisfying:

- (i) There exist two constants  $g_-, g_+ \in \mathbb{R}$  such that for all  $(y, s) \in W$ ,  $0 < g_- \leq g(y, s) \leq g_+$ .
- (ii) For any  $y \in \Gamma$ , the real function  $s \mapsto f(y, s)$  is non increasing on  $(0, \zeta_+(y))$  and non decreasing on  $(\zeta_-(y), 0)$ .

Then  $\alpha$  satisfies the condition **(H3)**.

*Proof.* Under the above conditions, it holds, for any  $y \in \Gamma$  and for  $s \in (0, \zeta_+(y))$ :

$$\left| \alpha(y, s) \int_0^s \alpha^{-1}(y, t) dt \right| \leq \frac{g_+}{g_-} f(y, s) \int_0^s f^{-1}(y, t) dt \leq \frac{g_+}{g_-} \zeta_+(y) \leq \frac{g_+}{g_-} \|\zeta_+ - \zeta_-\|_{L^\infty(\Gamma)}.$$

Arguing in the same fashion allows to prove a similar estimate when  $s \in (\zeta_-(y), 0)$ , which finally implies that  $\|h_\alpha\|_{L^\infty(\Gamma)} \leq \frac{g_+}{g_-} \|\zeta_+ - \zeta_-\|_{L^\infty(\Gamma)}$ ; this allows to conclude.  $\square$

**Remark 4.2.** The statement of lemma 4.1 does not require  $f$  or  $g$  to be continuous on  $W$ .

**Remark 4.3.** From the Liouville theorem for ordinary differential equations [95], it holds

$$\forall (y, s) \in W, |D\eta|(y, s) = \exp \left( \int_0^s \operatorname{div}(\beta) \circ \eta(y, t) dt \right). \quad (4.2.13)$$

Therefore, it is straightforward to verify that **(H1)** to **(H3)** are satisfied for the constant weight  $\omega = 1$  whenever  $\operatorname{div}(\beta) \in L^\infty(U)$ —a customary assumption in the study of advection operators; see e.g. [51, 195, 125, 257, 67].

Our setting, based on **(H1)** to **(H3)**, allows to handle more general velocity fields  $\beta$ , with unbounded divergence, which leaves room for the Jacobian  $|D\eta|$  to vanish on the boundary  $\partial U$ . This feature is crucial

to deal with the shape optimization setting (4.2.8); in there, the divergence  $\operatorname{div}(\beta)$  of  $\beta = \nabla d_\Omega$  blows up near *centers of curvatures*  $C \in \bar{\Sigma}$  (see Figure 4.2a and section 4.3.1 below) where one of the principal curvatures  $\kappa_i$  is such that  $-\kappa_i(p_{\partial\Omega}(x))d_\Omega(x) \rightarrow 1$  as  $x \rightarrow C$ . For instance, consider the following very simple situation where  $\Omega = D$  is the unit ball in two space dimensions. Its skeleton reduces to a single point, its center. Then,

$$\forall x \in D \setminus \{0\}, d_\Omega(x) = \|x\| - 1, \Delta d_\Omega(x) = \frac{1}{\|x\|},$$

where  $\|x\|$  is the euclidean norm in  $\mathbb{R}^2$ . Therefore  $\Delta d_\Omega(x)$  blows up at the center  $x = 0$  which implies  $\operatorname{div}(\nabla d_\Omega) \notin L^\infty(D \setminus \bar{\Sigma})$ . However, in this case, (H1) to (H3) hold for  $\omega = 1$ : from

$$\forall (y, s) \in \partial\Omega \times (-1, 0), \eta(y, s) = (1 + s)y, |\operatorname{D}\eta(y, s)| = 1 + s,$$

we obtain that  $\alpha(y, s) = 1 + s$  and  $\alpha(y, s) \int_0^s \alpha^{-1}(y, t) dt = (1 + s) \log(1 + s)$  are uniformly bounded on  $W = \partial\Omega \times (-1, 0)$ .

In the following, for a weight  $\omega$  satisfying (H1), we denote by  $L_\omega^2(U)$  and  $L_\alpha^2(W)$  the weighted spaces

$$L_\omega^2(U) = \left\{ v \text{ measurable} \mid \int_U \omega v^2 dx < +\infty \right\}, \quad L_\alpha^2(W) = \left\{ \tilde{v} \text{ measurable} \mid \int_W \alpha \tilde{v}^2 dx < +\infty \right\}, \quad (4.2.14)$$

with respective corresponding  $L^2$  norms  $\|v\|_{L_\omega^2(U)} := (\int_U \omega v^2 dx)^{1/2}$  and  $\|\tilde{v}\|_{L_\alpha^2(W)} = (\int_W \alpha \tilde{v}^2 ds dy)^{1/2}$ .

#### 4.2.2 Definition of the graph space $V_\omega$ of the advection operator $\beta \cdot \nabla$

We recall that the assumption  $\beta \in \mathcal{C}^1(U, \mathbb{R}^d)$  implies that  $\operatorname{div}(\beta)$  belongs to  $L_{loc}^\infty(U)$ .

**Definition 4.1** (Derivative along characteristic curves). Let  $v \in L_{loc}^1(U)$  be a locally integrable function on  $U$ . The directional derivative  $\beta \cdot \nabla v \in \mathcal{D}'(U)$  is the distribution on  $U$  defined by

$$\forall \phi \in \mathcal{C}_c^\infty(U), \int_U (\beta \cdot \nabla v) \phi dx = \int_U (-\beta \cdot \nabla \phi) v - \operatorname{div}(\beta) \phi v dx. \quad (4.2.15)$$

**Remark 4.4.** The definition (4.2.15) of  $\beta \cdot \nabla$  mimics the integration by part formula that holds for functions  $v$  of class  $\mathcal{C}^1$ . Considering test functions  $\phi \in \mathcal{C}_c^\infty(U)$  allows to avoid imposing classical regularity requirements on the open domain  $U$  such as that being locally located on one side of its boundary [302].

**Definition 4.2** (Graph space of the operator  $\beta \cdot \nabla$ ). Let  $\omega \in \mathcal{C}^0(U, \mathbb{R}_+^*)$  be a positive weight on  $U$ . The weighted graph space  $V_\omega$  of the advection operator  $\beta \cdot \nabla$  is defined by:

$$V_\omega = \{v \in L_\omega^2(U) \mid \beta \cdot \nabla v \in L_\omega^2(U)\}. \quad (4.2.16)$$

$V_\omega$  is a Hilbert space when it is endowed with the norm

$$\|v\|_{V_\omega} := \left( \int_U \omega v^2 dx + \int_U \omega |\beta \cdot \nabla v|^2 dx \right)^{1/2}.$$

**Remark 4.5.** The completeness of  $V_\omega$  follows from very standard arguments, see e.g. [306, 125]. It is also easy to see from the assumption  $\omega \in \mathcal{C}^0(U, \mathbb{R}_+^*)$  that the inclusion  $L_\omega^2(U) \subset L_{loc}^1(U)$  holds.

Let us then introduce corresponding notions of derivative with respect to the  $s$  variable and of graph space on the set  $W$ :

**Definition 4.3** (Graph space of the derivative  $\partial_s$ ). Let  $\alpha \in \mathcal{C}^0(W, \mathbb{R}_+^*)$  be a positive weight on  $W$ . The weak derivative  $\partial_s \tilde{v} \in \mathcal{D}'(W)$  of a function  $\tilde{v} \in L_{loc}^1(W)$  is the distribution defined by

$$\forall \phi \in \mathcal{C}_c^\infty(W), \int_W \partial_s \tilde{v} \phi dy ds = - \int_W \tilde{v} \partial_s \phi dy ds. \quad (4.2.17)$$

The weighted graph space of the operator  $\partial_s$  is defined by:

$$\tilde{V}_\alpha = \{\tilde{v} \in L_\alpha^2(W) \mid \partial_s \tilde{v} \in L_\alpha^2(W)\}. \quad (4.2.18)$$

This space is a Hilbert space when it is equipped with the norm

$$\|\tilde{v}\|_{\tilde{V}_\alpha} := \left( \int_W \alpha \tilde{v}^2 ds dy + \int_W \alpha |\partial_s \tilde{v}|^2 ds dy \right)^{1/2}.$$

**Remark 4.6.** The defining identity (4.2.15) of the weak derivative  $\beta \cdot \nabla v$  also holds for test functions  $\phi \in C_c^1(U)$ , as follows from a standard density argument. Likewise, (4.2.17) holds for test functions  $\phi$  which are only continuously differentiable with respect to the  $s$  variable.

The motivation behind the introduction of the spaces  $\tilde{V}_\alpha$  is that they allow to transfer the difficulty of studying the “curvilinear” directional derivative  $\beta \cdot \nabla$  over the domain  $U$ , to that of the more standard, “flat” derivative  $\partial_s$  over  $W$ . The price to pay is the need to account for the weight of the Jacobian  $|\mathbf{D}\eta|$  resulting from the change of variables. This point is made precise by the following technical lemma.

**Lemma 4.2.** *Let  $\omega \in C^0(U, \mathbb{R}_+^*)$  be a positive weight on  $U$ , and let  $\alpha = \omega \circ \eta |\mathbf{D}\eta|$ . The following equivalence holds true:*

$$v \in V_\omega \Leftrightarrow v \circ \eta \in \tilde{V}_\alpha.$$

Additionally, for any  $v \in V_\omega$ , the function  $\tilde{v} = v \circ \eta$  satisfies:

$$(i) \quad \partial_s \tilde{v} = (\beta \cdot \nabla v) \circ \eta,$$

$$(ii) \quad \|v\|_{V_\omega} = \|v \circ \eta\|_{\tilde{V}_\alpha}.$$

*Proof.* Let  $v \in V_\omega$ ; it follows from the change of variables (4.2.7) that  $\tilde{v} := v \circ \eta$  belongs to  $L_\alpha^2(W)$ . Moreover, we observe that  $\tilde{v}$  has a weak derivative  $\partial_s \tilde{v}$  given by:

$$\partial_s \tilde{v} = (\beta \cdot \nabla v) \circ \eta. \quad (4.2.19)$$

Indeed, using the definition (4.2.15) of the weak derivative along  $\beta$  with a test function  $\phi \in C_c^1(U)$  of the form  $\phi = \tilde{\phi} \circ \eta^{-1}$  for an arbitrary  $\tilde{\phi} \in C_c^1(W)$  (see remark 4.6), then using (4.2.7), we obtain:

$$\forall \tilde{\phi} \in C_c^1(W), \quad \int_W (\beta \cdot \nabla v) \circ \eta \tilde{\phi} |\mathbf{D}\eta| \, dsdy = \int_W (-\partial_s \tilde{\phi} \tilde{v} - \operatorname{div}(\beta) \circ \eta \tilde{\phi} \tilde{v}) |\mathbf{D}\eta| \, dsdy, \quad (4.2.20)$$

where we have used the equality  $\partial_s \tilde{\phi} = (\beta \cdot \nabla \phi) \circ \eta$ . Now using that  $|\mathbf{D}\eta|$  is continuously differentiable with respect to  $s$  with  $\partial_s |\mathbf{D}\eta| = \operatorname{div}(\beta) \circ \eta |\mathbf{D}\eta|$ , as a consequence of the Liouville formula (4.2.13), (4.2.20) rewrites:

$$\forall \tilde{\phi} \in C_c^1(W), \quad \int_W (\beta \cdot \nabla v) \circ \eta \tilde{\phi} |\mathbf{D}\eta| \, dsdy = \int_W (-\partial_s \tilde{\phi} \tilde{v} |\mathbf{D}\eta| - \partial_s |\mathbf{D}\eta| \tilde{\phi} \tilde{v}) \, dsdy = - \int_W \tilde{v} \partial_s (\tilde{\phi} |\mathbf{D}\eta|) \, dsdy.$$

By a standard density argument, the above equality holds more generally for functions  $\tilde{\phi}$  that are only continuously differentiable with respect to  $s$  on  $W$ . Therefore, taking  $\tilde{\phi} = \tilde{\psi} / |\mathbf{D}\eta|$  as a test function in the above equation for arbitrary  $\tilde{\psi} \in C_c^1(W)$  yields (4.2.19). The change of variables (4.2.7) now yields  $\|v\|_{V_\omega} = \|v \circ \eta\|_{\tilde{V}_\alpha}$ , and so  $v \circ \eta \in \tilde{V}_\alpha$ .

Conversely, one proves in a similar way that if  $\tilde{v}$  is in  $\tilde{V}_\alpha$ , the function  $v = \tilde{v} \circ \eta^{-1}$  belongs to  $V_\omega$ , which terminates the proof.  $\square$

**Remark 4.7.** In the following we prove a density result, a trace theorem and a Poincaré inequality in  $\tilde{V}_\alpha$  and we then obtain the direct counterparts of these results in the setting of the graph space  $V_\omega$  thanks to the above lemma 4.2. There exists actually a wide literature about weighted Sobolev spaces such as  $\tilde{V}_\alpha$  [148, 207, 209, 306], in which analogous results are proved for weights  $\alpha$  in the so-called Muckenhoupt class  $A_2$  (see e.g. [136] for a definition of the class  $A_p$ ). Our working assumptions however are of a different nature: for instance, the hypothesis (H3) essentially requires that the inverse weight  $\alpha^{-1}$  belong to the Muckenhoupt class  $A_1$ .

### 4.2.3 Density of functions of class $\mathcal{C}^1$ in the weighted space $V_\omega$

We now examine the density of  $C^\infty(W) \cap \tilde{V}_\alpha$  in  $\tilde{V}_\alpha$ , whence we shall infer the density of  $C^1(U) \cap V_\omega$  in  $V_\omega$ —the loss of regularity between both statements coming from the fact that the coordinate change  $\eta$  is only of class  $\mathcal{C}^1$  on  $W$ .

Our study classically involves mollifying functions; since the space  $\tilde{V}_\alpha$  of interest contains functions defined on the manifold  $W$ , a little treatment is in order. In particular, we shall need the so-called *Ahlfors regularity* of  $\Gamma$  (see [115]); this is the purpose of the next lemma, whose proof is outlined for the convenience of the reader. Here and throughout the chapter,  $B(y, h) \subset \mathbb{R}^d$  stands for the open ball with center  $y$  and radius  $h$ . The measure of a Lebesgue measurable set  $A \subset \mathbb{R}^d$  is denoted by  $|A|$ .

**Lemma 4.3** (Area covered by extrinsic balls on an embedded manifold). *Let  $\Gamma \subset \mathbb{R}^d$  be a  $\mathcal{C}^1$  hypersurface such that either  $\Gamma$  is compact ( $\Gamma = \bar{\Gamma}$ ) or  $\bar{\Gamma}$  is a compact manifold with boundary. Then there exists  $h_0 > 0$  and constants  $m > 0$  and  $M > 0$  depending only on  $\Gamma$  such that for any  $0 < h < h_0$ ,*

$$\forall y \in \Gamma, \quad mh^{d-1} \leq |\Gamma \cap B(y, h)| \leq Mh^{d-1}. \quad (4.2.21)$$

*Proof.* Since  $\bar{\Gamma}$  is compact, there exist finitely many open subset  $V_i \subset \mathbb{R}^d$ ,  $i = 1, \dots, N$ , such that  $\bar{\Gamma} \subset \bigcup_i V_i$ , and for each  $i = 1, \dots, N$ , there exists a local coordinate chart  $\phi_i : U_i \subset \mathbb{R}^{d-1} \rightarrow \phi_i(U_i) \subset \bar{\Gamma}$  such that  $\phi_i(U_i) = V_i \cap \bar{\Gamma}$ . The set  $U_i$  is a convex open subset of  $\mathbb{R}^{d-1}$  if  $V_i \cap \partial\Gamma = \emptyset$ ; when the latter intersection is not empty,  $U_i$  is of the form  $U_i = \tilde{U}_i \cap \mathbb{R}_+^{d-1}$ , where  $\tilde{U}_i \subset \mathbb{R}^{d-1}$  is a convex open subset, and  $\mathbb{R}_+^{d-1}$  is the upper half-space of  $\mathbb{R}^{d-1}$ . Let  $h_0 > 0$  be a Lebesgue's number associated with this cover, that is:

$$\forall y \in \Gamma, \exists i \in \{1, \dots, N\}, \quad B(y, h_0) \subset \subset V_i. \quad (4.2.22)$$

Now, given  $y \in V_i$ , one has for any  $0 < h < h_0$ :

$$|\Gamma \cap B(y, h)| = \int_{\Gamma} \mathbf{1}_{B(y, h)} dz = \int_{\phi_i(U_i)} \mathbf{1}_{B(y, h)} dz = \int_{U_i} \mathbf{1}_{\phi_i^{-1}(B(y, h))} |\mathbf{D}\phi_i| dx,$$

where  $i$  is the index supplied by (4.2.22) and  $\mathbf{1}_A$  denotes the characteristic function of a subset  $A \subset \mathbb{R}^d$ , and  $|\mathbf{D}\phi_i|$  is the Jacobian associated to the change of variables between manifolds induced by  $\phi_i$ .

Let  $\sigma_{d-1}(\nabla\phi_i(x)) \leq \sigma_1(\nabla\phi_i(x))$  be respectively the smallest and the largest singular values of the  $n \times (n-1)$  Jacobian matrix  $\nabla\phi_i(x)$ , and let  $0 < \sigma_- \leq \sigma_+$  be defined by:

$$\sigma_- = \min_{i=1, \dots, N} \inf_{x \in U_i} \sigma_{d-1}(\nabla\phi_i(x)), \quad \sigma_+ = \max_{i=1, \dots, N} \sup_{x \in U_i} \sigma_1(\nabla\phi_i(x)).$$

Applying the Taylor formula to  $\phi_i$  yields:

$$\forall i = 1, \dots, N, \quad \forall (x_0, x_1) \in U_i, \quad \sigma_- \|x_1 - x_0\| \leq \|\phi_i(x_1) - \phi_i(x_0)\| \leq \sigma_+ \|x_1 - x_0\|.$$

Therefore, possibly shrinking the value of the constant  $h_0$  supplied by (4.2.22) and taking  $x_0 = \phi_i^{-1}(y)$ , we obtain:

$$\forall 0 < h < h_0, \quad B\left(x_0, \frac{h}{\sigma_+}\right) \subset \phi_i^{-1}(B(y, h)) \subset B\left(x_0, \frac{h}{\sigma_-}\right) \subset U_i.$$

Finally, denoting by  $|B_{d-1}|$  the volume of the unit ball in  $\mathbb{R}^{d-1}$ , we obtain:

$$\begin{aligned} \frac{1}{2} |B_{d-1}| \left(\frac{\sigma_-}{\sigma_+}\right)^{d-1} h^{d-1} &\leq \int_{U_i} \mathbf{1}_{B(x_0, h/\sigma_+)} |\mathbf{D}\phi_i| dx \\ &\leq \int_{U_i} \mathbf{1}_{\phi_i^{-1}(B(y, h))} |\mathbf{D}\phi_i| dx = |\Gamma \cap B(y, h)| \\ &\leq \int_{U_i} \mathbf{1}_{B(x_0, h/\sigma_-)} |\mathbf{D}\phi_i| dx \leq |B_{d-1}| \left(\frac{\sigma_+}{\sigma_-}\right)^{d-1} h^{d-1}, \end{aligned}$$

which completes the proof.  $\square$

In what follows, the hypersurface  $\Gamma$  is the one introduced by section 4.2.1. In the next lemma, we construct the kernels  $\rho_h$  and  $\xi_h$  which shall be used for mollification purposes in  $W$ .

**Lemma 4.4.** *For any  $h > 0$ , there exist two positive, smooth functions  $\rho_h \in \mathcal{C}_c^\infty(\mathbb{R}^d \times \mathbb{R}^d, \mathbb{R})$  and  $\xi_h \in \mathcal{C}_c^\infty(\mathbb{R}, \mathbb{R})$  satisfying the conditions:*

- (i) For all points  $y \in \mathbb{R}^d$ ,  $\text{supp}(\rho_h(y, \cdot)) \subset B(y, h)$ .
- (ii) For any  $y \in \mathbb{R}^d$ ,  $\int_{\Gamma} \rho_h(y, z) dz = 1$ .
- (iii) There exist constants  $C > 0$  and  $h_0 > 0$  depending only on  $\Gamma$  such that

$$\forall 0 < h < h_0, \quad \forall z \in \Gamma, \quad \int_{\Gamma} \rho_h(y, z) dy \leq C. \quad (4.2.23)$$

- (iv)  $\text{supp}(\xi_h) \subset [-h, h]$  and  $\int_{\mathbb{R}} \xi_h(s) ds = 1$ .

*Proof.* Let  $\rho \in C_c^\infty(\mathbb{R}^d, \mathbb{R})$  be a smooth, positive function such that  $\text{supp}(\rho) \subset B(0, 1)$ ,  $\rho \leq 1$  on  $\mathbb{R}^d$  and  $\rho|_{B(0, 1/2)} = 1$ ; we define  $\rho_h$  by:

$$\forall h > 0, \forall (y, z) \in \mathbb{R}^d \times \mathbb{R}^d, \rho_h(y, z) := \frac{\rho\left(\frac{y-z}{h}\right)}{\int_\Gamma \rho\left(\frac{y-z}{h}\right) dz}; \quad (4.2.24)$$

note that  $\rho_h(y, z)$  is not a function of  $(y - z)$ . The conditions (i) and (ii) of the statement are obviously satisfied by (4.2.24). The condition (iii) is a consequence of lemma 4.3, which implies:

$$\forall 0 < h < h_0, \forall y \in \Gamma, \int_\Gamma \rho\left(\frac{y-z}{h}\right) dz \geq \left| \Gamma \cap B\left(y, \frac{h}{2}\right) \right| \geq m \left(\frac{h}{2}\right)^{d-1},$$

$$\forall 0 < h < h_0, \forall z \in \Gamma, \int_\Gamma \rho\left(\frac{y-z}{h}\right) dy \leq |\Gamma \cap B(z, h)| \leq Mh^{d-1},$$

so that (4.2.23) holds with  $C = 2^{d-1}M/m$ .

Eventually, a function  $\xi_h$  satisfying (iv) is constructed from any positive function  $\xi \in C_c^\infty(\mathbb{R}, \mathbb{R})$  with compact support inside  $[-1, 1]$  and unit integral over  $\mathbb{R}$ , by setting  $\xi_h = h^{-1}\xi(\cdot/h)$ .  $\square$

**Definition 4.4** (Mollification on the tensor product manifold  $W$ ). For  $h > 0$ , let  $\rho_h$  and  $\xi_h$  be two functions as in the statement of lemma 4.4. For any  $u \in L_{loc}^1(W)$ ,  $(y, s) \in W$  and  $h > 0$  sufficiently small (depending on  $(y, s)$ ), the mollification of  $u$  is the function  $u_h = \rho_h \xi_h * u$  defined by

$$(\rho_h \xi_h * u)(y, s) = \int_\Gamma \int_{\mathbb{R}} \rho_h(y, z) \xi_h(s - t) u(z, t) dt dz. \quad (4.2.25)$$

Note that for a given open subdomain  $W_1 \subset\subset W$  and because  $\bar{\Gamma}$  is compact, there exists  $h_0 > 0$ , depending on  $W_1$ , sufficiently small so that (4.2.25) makes sense for any  $(y, s) \in W_1$  and  $0 < h < h_0$ .

**Lemma 4.5.** *The following properties hold true:*

- (i) *If  $u \in L_{loc}^1(W)$ , for any subdomain  $W_1 \subset\subset W$  and for  $h > 0$  sufficiently small, the convolution  $\rho_h \xi_h * u$  is of class  $C^\infty$  on  $W_1$ .*
- (ii) *If  $u \in \tilde{V}_\alpha$ , then for any subdomain  $W_1 \subset\subset W$  and for  $h$  sufficiently small,  $\partial_s(\rho_h \xi_h * u) = \rho_h \xi_h * \partial_s u$ .*
- (iii) *If  $\phi \in C^0(W, \mathbb{R})$  then  $\rho_h \xi_h * \phi \rightarrow \phi$  in  $L_{loc}^\infty(W)$  as  $h \rightarrow 0$ .*

*Proof.* (i) This results from the Lebesgue dominated convergence theorem and the smoothness of  $\rho_h$  and  $\xi_h$ .

(ii) This is again a consequence of the Lebesgue dominated theorem and of an integration by parts:

$$\begin{aligned} \forall (y, s) \in W_1, (\rho_h \xi_h * \partial_s u)(y, s) &= \int_\Gamma \int_{\mathbb{R}} \rho_h(y, z) \xi_h(s - t) \partial_s u(z, t) dt dz \\ &= \int_\Gamma \int_{\mathbb{R}} \rho_h(y, z) \partial_s \xi_h(s - t) u(z, t) dt dz = \partial_s(\rho_h \xi_h * u). \end{aligned}$$

(iii) For a given subset  $W_1 \subset\subset W$ , let  $\varepsilon > 0$  and  $h_1 > 0$  be uniform continuity constants for  $\phi$  on  $W_1$ , i.e.

$$\forall 0 < h < h_1, (y, s) \in W_1, (z, t) \in W_1, (||y - z|| < h \text{ and } |s - t| < h) \Rightarrow |\phi(z, t) - \phi(y, s)| < \varepsilon.$$

Then for any  $0 < h < h_1$ :

$$\begin{aligned} \forall (y, s) \in W_1, |\rho_h \xi_h * \phi(y, s) - \phi(y, s)| &= \left| \int_\Gamma \int_{\mathbb{R}} \rho_h(y, z) \xi_h(s - t) (\phi(z, t) - \phi(y, s)) dt dz \right| \\ &\leq \varepsilon \int_\Gamma \int_{\mathbb{R}} \rho_h(y, z) \xi_h(s - t) dt dz = \varepsilon. \end{aligned}$$

$\square$

**Proposition 4.1.** *Let  $\alpha \in \mathcal{C}^0(W, \mathbb{R}_+^*)$ . Then for any function  $u \in L_{\alpha, \text{loc}}^2(W)$ ,  $u_h := \rho_h \boldsymbol{\xi}_h * u$  belongs to  $L_{\alpha, \text{loc}}^2(W) \cap \mathcal{C}^\infty(W)$  and  $u_h \rightarrow u$  in  $L_{\alpha, \text{loc}}^2(W)$ .*

*Proof.* We adapt the proof of Prop. 1.18, p.30 in [268]. Let  $W_1 \subset\subset W$  be an open subset of  $W$  and consider another open subset  $W_2$  such that  $W_1 \subset\subset W_2 \subset\subset W$ . For  $\varepsilon > 0$ , let  $\phi \in \mathcal{C}^0(W)$  be such that (see [269]):

$$\|u - \phi\|_{L_\alpha^2(W_2)} < \varepsilon. \quad (4.2.26)$$

From lemma 4.5,  $\rho_h \boldsymbol{\xi}_h * \phi \rightarrow \phi$  in  $L^\infty(W_1)$  as  $h \rightarrow 0$  and so, for  $h$  small enough,  $\|\rho_h \boldsymbol{\xi}_h * \phi - \phi\|_{L_\alpha^2(W_1)} \leq \varepsilon$ . Then,

$$\|u - u_h\|_{L_\alpha^2(W_1)} \leq \|u - \phi\|_{L_\alpha^2(W_1)} + \|\phi - \rho_h \boldsymbol{\xi}_h * \phi\|_{L_\alpha^2(W_1)} + \|\rho_h \boldsymbol{\xi}_h * \phi - \rho_h \boldsymbol{\xi}_h * u\|_{L_\alpha^2(W_1)}. \quad (4.2.27)$$

The first two terms in the right-hand side of (4.2.27) are controlled by  $\varepsilon$  owing to the previous discussion; as for the last term, we get from the Cauchy-Schwarz inequality:

$$\begin{aligned} \forall (y, s) \in W_1, |\rho_h \boldsymbol{\xi}_h * (\phi - u)(y, s)|^2 &= \left| \int_\Gamma \int_{\mathbb{R}} \rho_h(y, z) \boldsymbol{\xi}_h(s-t) (\phi(z, t) - u(z, t)) dt dz \right|^2 \\ &\leq \left( \int_\Gamma \int_{\mathbb{R}} \rho_h(y, z) \boldsymbol{\xi}_h(s-t) \alpha^{-1}(z, t) dt dz \right) \left( \int_\Gamma \int_{\mathbb{R}} \rho_h(y, z) \boldsymbol{\xi}_h(s-t) \alpha(z, t) |\phi(z, t) - u(z, t)|^2 dt dz \right). \end{aligned}$$

Multiplying both sides by  $\alpha$  and integrating over  $W_1$  yields:

$$\begin{aligned} \|\rho_h \boldsymbol{\xi}_h * (\phi - u)\|_{L_\alpha^2(W_1)}^2 &\leq \|\alpha(\rho_h \boldsymbol{\xi}_h * \alpha^{-1})\|_{L^\infty(W_1)} \int_{(y,s) \in W_1} \left( \int_\Gamma \int_{\mathbb{R}} \rho_h(y, z) \boldsymbol{\xi}_h(s-t) \alpha(z, t) |\phi(z, t) - u(z, t)|^2 dt dz \right) ds dy \\ &\leq C \|\alpha(\rho_h \boldsymbol{\xi}_h * \alpha^{-1})\|_{L^\infty(W_1)} \|\phi - u\|_{L_\alpha^2(W_2)}^2, \end{aligned}$$

where  $C$  is the constant supplied by the condition (iii) in the statement of lemma 4.4. By assumption,  $\alpha^{-1}$  is a continuous function on  $W_2$ , which implies by lemma 4.5 that  $\alpha(\rho_h \boldsymbol{\xi}_h * \alpha^{-1}) \rightarrow 1$  in  $L^\infty(W_1)$ . In particular,  $\|\alpha(\rho_h \boldsymbol{\xi}_h * \alpha^{-1})\|_{L^\infty(W_1)}$  is bounded by some constant  $C'$ . Finally, using (4.2.26), we obtain from (4.2.27) that for  $h > 0$  small enough:

$$\|u - u_h\|_{L_\alpha^2(W_1)} \leq (CC' + 2)\varepsilon,$$

which is the desired result.  $\square$

We conclude this subsection with the desired density result of  $\mathcal{C}^1$  functions in  $V_\omega$ :

**Corollary 4.1.** (i) *Let  $\alpha \in \mathcal{C}^0(W, \mathbb{R}_+^*)$  be a positive weight on  $W$ ; the space  $\mathcal{C}^\infty(W) \cap \tilde{V}_\alpha$  is dense in  $\tilde{V}_\alpha$ .*

(ii) *Let  $\omega \in \mathcal{C}^0(U, \mathbb{R}_+^*)$  be a positive weight on  $U$ ; the space  $\mathcal{C}^1(U) \cap V_\omega$  is dense in  $V_\omega$ .*

*Proof.* The proof of the density (i) of  $\mathcal{C}^\infty(W) \cap \tilde{V}_\alpha$  in  $\tilde{V}_\alpha$  relies on a partition of unity argument and on the properties of lemma 4.5 and proposition 4.1, exactly along the lines of the proof of Theorem 5.15 in [268], to which the reader is referred for details.

The density (ii) of  $\mathcal{C}^1(U)$  in  $V_\omega$  follows then from the density of  $\mathcal{C}^\infty(W) \cap \tilde{V}_\alpha$  in  $\tilde{V}_\alpha$  with  $\alpha = \omega \circ \eta|_{D\eta}$  and by composition with the  $\mathcal{C}^1$  diffeomorphism  $\eta$ .  $\square$

**Remark 4.8.** Corollary 4.1 does not imply the density of  $\mathcal{C}^1(\bar{U})$  in  $V_\omega$ . A result of this kind would require careful regularity assumptions on  $\partial U$  and on the behavior of  $\omega$  near  $\partial U$ .

#### 4.2.4 Trace theorem and Poincaré inequality in $V_\omega$

In this section, the trace operator on  $\Gamma$  is defined and studied for functions in the weighted space  $V_\omega$ , or equivalently for functions in  $\tilde{V}_\alpha$  on  $\Gamma \times \{0\} = \{(y, 0) \mid y \in \Gamma\}$ . In the sequel, with a little abuse of notations, the latter set  $\Gamma \times \{0\}$  is identified with  $\Gamma$ .



**Proposition 4.2** (Trace theorem). *Let  $\omega \in \mathcal{C}^0(U, \mathbb{R}_+^*)$  be a positive weight on  $U$ . The trace operator*

$$\begin{aligned} \gamma &: \mathcal{C}^1(U) &\rightarrow & L^2(\Gamma) \\ v &\mapsto \gamma(v) &= & v|_\Gamma \end{aligned} \quad (4.2.28)$$

*induces a bounded operator  $V_\omega \rightarrow L^2(\Gamma)$ ; there exists a constant  $C > 0$  (possibly depending on the weight  $\omega$ ) such that*

$$\forall v \in V_\omega, \|\gamma(v)\|_{L^2(\Gamma)} \leq C\|v\|_{V_\omega}. \quad (4.2.29)$$

*Proof.* Introducing  $\alpha = \omega \circ \eta|D\eta|$ , using the change of variables (4.2.7) and the density result of [corollary 4.1](#), it is enough to prove that there exists a constant  $C > 0$  such that:

$$\forall \tilde{v} \in \mathcal{C}^\infty(W) \cap \tilde{V}_\alpha, \|\gamma(\tilde{v})\|_{L^2(\Gamma)} \leq C\|\tilde{v}\|_{\tilde{V}_\alpha}. \quad (4.2.30)$$

Let us consider the following partition of  $\Gamma$ :  $\Gamma = \Gamma_+ \cup \Gamma_-$ , where  $\Gamma_+ = \{y \in \Gamma \mid \zeta_-(y) \geq -\varepsilon/2\}$ ,  $\Gamma_- = \Gamma \setminus \Gamma_+$ , and  $\varepsilon$  is the parameter featured in the separation condition (4.2.2). Then

$$\forall y \in \Gamma_+, 0 < \varepsilon/2 < \zeta_+(y),$$

$$\forall y \in \Gamma_-, \zeta_-(y) < -\varepsilon/2 < 0.$$

Now, let  $a = \varepsilon/2$  and  $\xi \in \mathcal{C}_c^\infty(\mathbb{R})$  be such that  $\xi(0) = 1$  and  $\xi(-a) = \xi(a) = 0$ . Let  $K$  be the bounded measurable set defined by  $K := (\Gamma_- \times [-a, 0]) \cup (\Gamma_+ \times [0, a]) \subset W$ . Then for any function  $\tilde{v} \in \mathcal{C}^\infty(W) \cap \tilde{V}_\alpha$ , it holds:

$$\begin{aligned} \int_\Gamma |\tilde{v}(y, 0)|^2 dy &= \int_\Gamma |\tilde{v}(y, 0)|^2 \xi(0) dy = \int_{\Gamma_-} \int_{-a}^0 \partial_s(\tilde{v}^2 \xi) ds dy + \int_{\Gamma_+} \int_a^0 \partial_s(\tilde{v}^2 \xi) ds dy \\ &\leq \int_K |2\tilde{v} \partial_s \tilde{v} \xi + \tilde{v}^2 \partial_s \xi| ds dy \\ &\leq 2\|\alpha^{-1} \xi\|_{L^\infty(K)} \|\tilde{v}\|_{L_\alpha^2(W)} \|\partial_s \tilde{v}\|_{L_\alpha^2(W)} + \|\alpha^{-1} \partial_s \xi\|_{L^\infty(K)} \|\tilde{v}\|_{L_\alpha^2(W)}^2 \\ &\leq (2\|\alpha^{-1} \xi\|_{L^\infty(K)} + \|\alpha^{-1} \partial_s \xi\|_{L^\infty(K)}) \|\tilde{v}\|_{\tilde{V}_\alpha}^2, \end{aligned}$$

which implies (4.2.30) and therefore terminates the proof of [proposition 4.2](#).  $\square$

**Remark 4.9.** The proof of [proposition 4.2](#) supplies the existence of the trace on  $\Gamma$  of an arbitrary function  $\tilde{v} \in \tilde{V}_\alpha$ , which we shall also denote by  $\tilde{v}|_\Gamma$ .

For later purposes (see [section 4.2.5](#)), we shall need the surjectivity of the above trace operator; this is the purpose of the next proposition.

**Proposition 4.3** (Surjectivity of traces). *Let  $\omega \in \mathcal{C}^0(U, \mathbb{R}_+^*)$  be a positive continuous weight.*

(i) *The trace operator defined by (4.2.28) is surjective from  $V_\omega$  onto  $L^2(\Gamma)$ :*

$$L^2(\Gamma) = \{v|_\Gamma \mid v \in V_\omega\}.$$

(ii) *If  $\omega$  additionally satisfies (H2), then any function  $v_0 \in L^2(\Gamma)$  can be extended constantly along the characteristics of  $\beta$ : there exists  $v \in V_\omega$  such that  $v|_\Gamma = v_0$  and  $\beta \cdot \nabla v = 0$ .*

*Proof.* (i) We rather prove that  $L^2(\Gamma) = \{\tilde{v}|_\Gamma \mid \tilde{v} \in \tilde{V}_\alpha\}$ , where  $\alpha = \omega \circ \eta|D\eta|$ ; see [remark 4.9](#). To this end, consider  $\xi$  and  $K$  be as in the proof of [proposition 4.2](#). For an arbitrary function  $v_0 \in L^2(\Gamma)$ , we define  $\tilde{v}$  by the formula

$$\tilde{v}(y, s) := v_0(y) \xi(s) \text{ a.e. in } W.$$

Obviously,  $\tilde{v}(y, 0) = v_0(y)$  and  $\partial_s \tilde{v}(y, s) = v_0(y) \partial_s \xi(s)$ , whence the Cauchy-Schwarz inequality yields:

$$\int_W \alpha \tilde{v}^2 ds dy = \int_\Gamma \int_{\zeta_-(y)}^{\zeta_+(y)} \alpha(y, s) \xi(s)^2 v_0(y)^2 ds dy < \|\alpha\|_{L^\infty(K)} \|\xi\|_{L^2(\mathbb{R})} \|v_0\|_{L^2(\Gamma)}^2 < +\infty,$$

$$\int_W \alpha |\partial_s \tilde{v}|^2 ds dy = \int_\Gamma \int_{\zeta_-(y)}^{\zeta_+(y)} \alpha(y, s) v_0(y)^2 |\partial_s \xi(s)|^2 ds dy \leq \|\alpha\|_{L^\infty(K)} \|\partial_s \xi\|_{L^2(\mathbb{R})} \|v_0\|_{L^2(\Gamma)}^2 < +\infty.$$

Hence  $\tilde{v}$  is a function in  $\tilde{V}_\alpha$  such that  $\tilde{v}|_\Gamma = v_0$ , which is the desired result.

(ii) If (H2) is satisfied, then for an arbitrary function  $v_0 \in L^2(\Gamma)$ , we simply define  $\tilde{v}$  by the formula:

$$\tilde{v}(y, s) := v_0(y) \text{ a.e. in } W.$$

Then, clearly  $\partial_s \tilde{v} = 0$ ; what's more, one has  $\tilde{v} \in L^2_\alpha(W)$  as follows from the following estimate:

$$\int_W \alpha \tilde{v}^2 ds dy = \int_\Gamma \int_{\zeta_-(y)}^{\zeta_+(y)} \alpha(y, s) v_0(y)^2 ds dy \leq \|g_\alpha\|_{L^\infty(\Gamma)} \|v_0\|_{L^2(\Gamma)}^2, \quad (4.2.31)$$

where  $g_\alpha$  is as in the statement of (H2). Hence the function  $\tilde{v}$  belongs to  $\tilde{V}_\alpha$ , and so  $v := \tilde{v} \circ \eta^{-1}$  is an element of  $V_\omega$  which has the desired properties owing to lemma 4.2.  $\square$

We now prove a Poincaré-type inequality in the spaces  $V_\omega$  under the additional assumptions (H2) and (H3) about the weight  $\omega$ .

**Proposition 4.4** (Poincaré type inequality on  $V_\omega$ ). *Let  $\omega \in C^0(U, \mathbb{R}_+^*)$  be a weight satisfying the assumptions (H1) to (H3). Then there exists a constant  $C > 0$  (depending on  $\omega$ ) such that:*

$$\forall v \in V_\omega, \int_U \omega v^2 dx \leq C \left[ \int_\Gamma v^2 dy + \int_U \omega |\beta \cdot \nabla v|^2 dx \right]. \quad (4.2.32)$$

*Proof.* Introducing again  $\alpha = \omega \circ \eta |D\eta| \in C^0(W, \mathbb{R}_+^*)$  and using the change of variables (4.2.7), we rather prove the analogous Poincaré inequality in  $W$ , that is:

$$\forall \tilde{v} \in \tilde{V}_\alpha, \int_W \alpha \tilde{v}^2 ds dy \leq C \left[ \int_\Gamma \tilde{v}^2 ds + \int_W \alpha |\partial_s \tilde{v}|^2 ds dy \right]. \quad (4.2.33)$$

Furthermore, since  $C^\infty(W) \cap \tilde{V}_\alpha$  is dense in  $\tilde{V}_\alpha$ , it is enough to prove that (4.2.33) holds for  $\tilde{v} \in C^\infty(W) \cap \tilde{V}_\alpha$ , which we now do. To this end, for arbitrary  $\tilde{v} \in C^\infty(W) \cap \tilde{V}_\alpha$ , a use of Taylor's formula yields:

$$\forall (y, s) \in W, \tilde{v}(y, s) = \tilde{v}(y, 0) + \int_0^s \partial_s \tilde{v}(y, t) dt. \quad (4.2.34)$$

In (4.2.34), the Cauchy-Schwarz inequality implies that, for  $(y, s) \in W$ :

$$\int_0^s |\partial_s \tilde{v}(y, t)| dt \leq \left( \int_0^s \alpha^{-1}(y, t) dt \right)^{1/2} \left( \int_{\zeta_-(y)}^{\zeta_+(y)} \alpha |\partial_s \tilde{v}|^2(y, t) dt \right)^{1/2}. \quad (4.2.35)$$

Now squaring (4.2.34), using the Young's inequality ( $\forall a, b \in \mathbb{R}, (a + b)^2 \leq 2a^2 + 2b^2$ ) together with (4.2.35), then multiplying by  $\alpha(y, s)$ , we obtain:

$$\forall (y, s) \in W, \alpha(y, s) |\tilde{v}(y, s)|^2 \leq 2\alpha(y, s) |\tilde{v}(y, 0)|^2 + 2\alpha(y, s) \int_0^s \alpha^{-1}(y, t) dt \int_{\zeta_-(y)}^{\zeta_+(y)} \alpha(y, t) |\partial_s \tilde{v}|^2(y, t) dt. \quad (4.2.36)$$

Integrating (4.2.36) over  $W$  now results in:

$$\begin{aligned} \int_W \alpha \tilde{v}^2 ds dy &\leq 2 \int_\Gamma \tilde{v}^2(y, 0) \left( \int_{\zeta_-(y)}^{\zeta_+(y)} \alpha(y, s) ds \right) dy + 2 \int_\Gamma h_\alpha(y) \left( \int_{\zeta_-(y)}^{\zeta_+(y)} \alpha |\partial_s \tilde{v}|^2 ds \right) dy \\ &\leq 2 \|g_\alpha\|_{L^\infty(\Gamma)} \int_\Gamma \tilde{v}^2 dy + 2 \|h_\alpha\|_{L^\infty(\Gamma)} \int_W \alpha |\partial_s \tilde{v}|^2 ds dy, \end{aligned}$$

where  $g_\alpha$  and  $h_\alpha$  are the functions featured in (H2) and (H3). This completes the proof of (4.2.33), and so that of proposition 4.4.  $\square$

**Remark 4.10.** This Poincaré type inequality is close in spirit to the ‘‘curvilinear’’ Poincaré inequality of Azerad [50, 51], who considered the weight  $\omega = 1$  and vector fields  $\beta$  satisfying  $\operatorname{div} \beta = 0$ .

### 4.2.5 Well-posedness of the variational problem (4.2.5)

We are now in a position to state and prove the main result of this section.

**Proposition 4.5.** *Let  $\omega \in \mathcal{C}^0(U, \mathbb{R}_+^*)$  be a positive weight on  $U$ , satisfying the assumptions (H1) to (H3). Then:*

(i) *For any function  $f \in L_{\omega^{-1}}^2(U)$ , there exists a unique solution  $u \in V_\omega$  to the variational problem*

$$\text{Find } u \in V_\omega \text{ such that } \forall v \in V_\omega, \int_{\Gamma} uv ds + \int_U \omega(\beta \cdot \nabla u)(\beta \cdot \nabla v) dx = \int_U f v dx. \quad (4.2.37)$$

(ii) *The trace  $u|_{\Gamma} \in L^2(\Gamma)$  of the solution  $u$  to (4.2.37), is independent of the weight  $\omega$ ; it is given by:*

$$u(y) = \int_{\zeta_-(y)}^{\zeta_+(y)} f \circ \eta |D\eta| ds \quad \text{a.e. on } \Gamma. \quad (4.2.38)$$

*Proof.* (i) The assumption  $f \in L_{\omega^{-1}}^2(U)$  ensures that  $v \mapsto \int_U f v dx$  is a continuous linear form on  $V_\omega$  owing to the Cauchy-Schwartz inequality:

$$\forall v \in V_\omega, \left| \int_U f v dx \right| \leq \|f\|_{L_{\omega^{-1}}^2(U)} \|v\|_{L_\omega^2(U)}.$$

Moreover, the Poincaré inequality of proposition 4.4 ensures the coercivity on  $V_\omega$  of the bilinear form

$$(u, v) \mapsto \int_{\Gamma} uv dy + \int_U \omega(\beta \cdot \nabla u)(\beta \cdot \nabla v) dx.$$

Hence, the classical Lax-Milgram theorem (see e.g. [144]) yields the existence and uniqueness of a solution  $u \in V_\omega$  to (4.2.37).

(ii) Since  $\omega$  satisfies (H2), it follows from proposition 4.3 (ii) that for any  $v_0 \in \mathcal{C}^0(\Gamma)$ , there exists a function  $v \in V_\omega$  such that  $v|_{\Gamma} = v_0$  and  $\beta \cdot \nabla v = 0$ . Taking  $v$  as an admissible test function in (4.2.37) and using once again the change of variables (4.2.7) yields:

$$\forall v_0 \in \mathcal{C}^0(\Gamma), \int_{\Gamma} uv_0 dy = \int_U f v dx = \int_{\Gamma} v_0(y) \left( \int_{\zeta_-(y)}^{\zeta_+(y)} f \circ \eta |D\eta| ds \right) dy, \quad (4.2.39)$$

whence (4.2.38). □

**Remark 4.11.** Problem (4.2.5) or (4.2.37) can be solved, for an arbitrary choice of weight  $\omega$  satisfying (H1) to (H3), if the right hand side satisfies  $f \in L_{\omega^{-1}}^2(U)$ . This holds true, for example when  $f$  belongs to  $L^\infty(U)$ , as soon as the weight  $\omega$  satisfies:  $\omega^{-1} \in L^1(U)$ . The latter property is not a consequence of (H1) to (H3) as explained in Remark 4.13.

**Remark 4.12.** If (H2) is not satisfied, which is the case if for example  $\omega$  blows up “too fast” near some part of the boundary of  $U$ , functions  $v \in V_\omega$  are expected to vanish near this part of  $\partial U$  and then  $V_\omega$  may not contain all functions which are constant along the characteristic curves of  $\beta$ .

## 4.3 NUMERICAL METHODS FOR INTEGRATION ALONG NORMAL RAYS

This section provides recipes for the implementation of the variational formulation (4.2.5) with the finite element method, and numerical comparisons with the line integral formula (4.2.6). For simplicity and because we are motivated by our shape optimization applications, this comparison is considered in the context where these lines are normal rays to the shape.

The general setting of section 4.2 to the shape optimization context outlined in the introduction; in section 4.3.1. We then clarify, for comparison purposes, some of the algorithmic stages required by the direct integration of (4.2.6) along rays in section 4.3.2—such as the numerical computation of the principal curvatures  $\kappa_i$  of the considered shapes, and the delicate detection of their skeleton on unstructured meshes—. We discuss next in section 4.3.3 the construction of suitable weights  $\omega$  that allow to solve accurately the variational problem (4.2.5) by means of  $\mathbb{P}_1$  conforming finite elements. Finally, in section 4.3.4, we compare on several 2-d analytical examples the numerical calculations of (4.2.6) by means of our variational formulation (4.2.5) with those produced by direct integration of this formula along rays.

### 4.3.1 Shape optimization context: normal rays and flow map of the signed distance function gradient

This section explicits the shape optimization context in the framework of the previous [section 4.2](#) and in the perspective of applying the variational problem [\(4.2.5\)](#) for enforcing distance constraints. A fixed, bounded and Lipschitz ‘hold-all’ open domain  $D \subset \mathbb{R}^d$  is considered, as well as a bounded  $\mathcal{C}^2$  open subdomain  $\Omega \subset D$ .

The vector field  $\beta$  is given by the gradient of the signed distance function:

$$U = D \setminus \bar{\Sigma}, \quad \Gamma = \partial\Omega, \quad \beta = \nabla d_\Omega. \quad (4.3.1)$$

The bound functions  $\zeta_-(y)$  and  $\zeta_+(y)$  correspond to the extremities of the normal rays as follows:

**Definition 4.5** (Normal rays). For  $y \in \partial\Omega$ , the ray emerging from  $y$  is the one-dimensional segment

$$\text{ray}(y) := \{x \in D \setminus \bar{\Sigma}, p_{\partial\Omega}(x) = y\} = \{y + s\nabla d_\Omega(y) \mid s \in (\zeta_-(y), \zeta_+(y))\}. \quad (4.3.2)$$

where  $\zeta_-(y)$  and  $\zeta_+(y)$  are the maximum distances at which  $\text{ray}(y)$  hits either the skeleton  $\Sigma$  or the boundary  $\partial D$  of the hold-all domain:

$$\forall y \in \partial\Omega, \quad \zeta_+(y) = \sup\{s \geq 0 \mid \{y + t\nabla d_\Omega(y) \mid t \in [0, s]\} \cap (\bar{\Sigma} \cup \partial D) = \emptyset\}, \quad (4.3.3)$$

$$\forall y \in \partial\Omega, \quad \zeta_-(y) = \inf\{s \leq 0 \mid \{y + t\nabla d_\Omega(y) \mid t \in (s, 0]\} \cap (\bar{\Sigma} \cup \partial D) = \emptyset\}. \quad (4.3.4)$$

The functions  $\zeta_-$  and  $\zeta_+$  are continuous on  $\partial\Omega$  (see [\[82\]](#) or [\[215\]](#)).

Finally, the next proposition explicits the local coordinate change  $\eta$  and its Jacobian  $|\text{D}\eta|$  featured in the identity [\(4.2.6\)](#) (see [\[64\]](#) and [chapter 1, section 1.3](#)). Because this will be useful when investigating the formulation of relevant weights  $\omega$  in [lemma 4.6](#), we also recall the value of  $\Delta d_\Omega$ .

**Proposition 4.6** (Shape optimization setting). *Let  $\Omega \subset D$  be a domain of class  $\mathcal{C}^2$ ; then the signed distance function  $d_\Omega$  is of class  $\mathcal{C}^2$  on the open set  $U := D \setminus \bar{\Sigma}$ . Hence  $\beta := \nabla d_\Omega$  is a vector field of class  $\mathcal{C}^1$  on  $U$ ; the associated flow map  $\eta : W \rightarrow D \setminus \bar{\Sigma}$  is a diffeomorphism of class  $\mathcal{C}^1$ , whose expression reads:*

$$\forall (y, s) \in W, \quad \eta(y, s) = y + s\nabla d_\Omega(y), \quad (4.3.5)$$

where  $W$  is the set defined by [\(4.2.3\)](#). The inverse flow mapping  $\eta^{-1} : U \rightarrow W$  is given by

$$\forall x \in U, \quad \eta^{-1}(x) = (p_{\partial\Omega}(x), d_\Omega(x)). \quad (4.3.6)$$

The divergence of the vector field  $\beta = \nabla d_\Omega$  and the Jacobian  $|\text{D}\eta|$  of the flow map  $\eta$  are respectively given by

$$\forall x \in D \setminus \bar{\Sigma}, \quad \text{div}(\nabla d_\Omega)(x) = \Delta d_\Omega(x) = \sum_{i=1}^{d-1} \frac{\kappa_i(p_{\partial\Omega}(x))}{1 + d_\Omega(x)\kappa_i(p_{\partial\Omega}(x))}, \quad (4.3.7)$$

$$\forall (y, s) \in W, \quad |\text{D}\eta|(y, s) = |\det(\nabla\eta)|(y, s) = \prod_{i=1}^{d-1} (1 + s\kappa_i(y)). \quad (4.3.8)$$

### 4.3.2 Computing curvatures and detecting the skeleton for direct integration along the rays

Before going to the numerical aspects of our variational method, we clarify important practical details which are required in the implementation of the direct integration along characteristics involved in the calculation of [\(4.2.6\)](#).

We first discuss the delicate issue of detecting mesh triangles crossing the skeleton  $\Sigma$  of  $\Omega$  when traveling along the rays (note that this step is not required by our variational method). Then, we detail the method we used to compute the curvature  $\kappa$  (there is only one curvature  $\kappa := \kappa_1$  in 2-d) required in the line integral formula [\(4.1.3\)](#). Note that these steps serve only for comparison purposes with our variational method. These are fairly classical numerical issues which could otherwise be addressed with more sophisticated techniques, see e.g. [\[270, 124\]](#).

For our present numerical applications, the hold-all domain  $D$  is equipped with a simplicial mesh  $\mathcal{T}$  featuring a discretization of the domain  $\Omega$  as a submesh. The only information we use about  $\Omega$  is an accurate approximation  $d_h$  of the signed distance function  $d_\Omega$  as an element of the space  $V_h$  of Lagrange  $\mathbb{P}_1$  finite element, where  $h$  is the maximum mesh element size. In our context, this approximation  $d_h$  is obtained by the algorithm of [\[111\]](#) implemented in the software program `mshdist`.

### Detection of the skeleton $\Sigma$ of $\Omega$ and identification of normal rays

The numerical detection of  $\Sigma$  in the course of the identification of the set  $\text{ray}(y)$  for some given point  $y \in \partial\Omega$  is achieved by assessing the following criterion (independently of the dimension), holding at the continuous level (see [122]):

$$\forall y \in \partial\Omega, \forall z_0, z_1 \in \text{ray}(y), \frac{d_\Omega(z_1) - d_\Omega(z_0)}{(z_1 - z_0) \cdot \nabla d_\Omega(y)} = 1.$$

In our implementation, when computing the ray emerging at some point  $y \in \partial\Omega$  (which is detected by the fact that  $d_h(y) = 0$ ), we travel the triangles in the mesh  $\mathcal{T}$  in the normal direction  $\mathbf{n} = \nabla d_h(y)$ , and we stop the calculation of the ray in the triangle  $T \in \mathcal{T}$  where the entering and exiting points  $z_0$  and  $z_1$  satisfy:

$$\left| \frac{d_h(z_1) - d_h(z_0)}{\|z_1 - z_0\|} - \text{sign}(d_h(z_0)) \right| \geq \text{tolRay}, \quad (4.3.9)$$

where  $\text{tolRay}$  a small tolerance (set to 0.3 in our implementation). This provides meanwhile an approximate location of the skeleton  $\Sigma$ , up to a tolerance of the order of the mesh size.

Our criterion (4.3.9) differs from that used in the related works [30, 234]. In there, the authors detect  $\Sigma$  by looking at changes in the monotonicity of the signed distance function  $d_h$  along the ray, i.e. they rely on the following property of the (continuous) signed distance function  $d_\Omega$ :

$$\forall y \in \partial\Omega, \forall z_0, z_1 \in \text{ray}(y), (d_\Omega(z_1) - d_\Omega(z_0))((z_1 - z_0) \cdot \nabla d_\Omega(y)) \geq 0. \quad (4.3.10)$$

Our personal experiment with the above criterion suggests that it may sometimes fail to detect the skeleton  $\Sigma$  accurately, because such a change in monotonicity may simply not occur when the ray is supposed to cross  $\Sigma$  in the neighborhood of center of curvatures (see remark 1.10). Our criterion (4.3.9) may also fail depending on the chosen tolerance parameter, but it offered visible improvements (Figure 4.4) in our academic test-cases. Note that when integrating along the ray, the last triangle, where the skeleton is hit, is included in the integration.

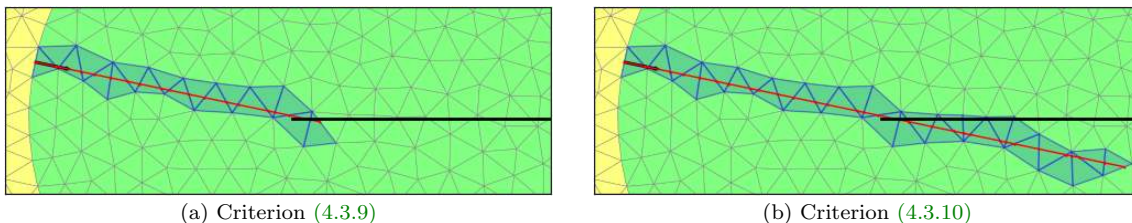


Figure 4.4: Comparison between the two criteria of section 4.3.2 for the detection of  $\Sigma$  when travelling along rays in an unstructured mesh (Skeleton displayed in the black line).

### Estimating the curvature $\kappa$ of a 2-d subdomain based on its signed distance function

In this part, we detail our method for the numerical approximation, in 2-d, of the principal curvature  $\kappa$  (there is only one in 2-d) from the knowledge of a  $\mathbb{P}^1$  discretization  $d_h$  of the signed distance function  $d_\Omega$  at the nodes of the mesh  $\mathcal{T}$ . We essentially rely on the fact that in 2-d,  $\kappa$  is given by the trace of  $\Delta d_\Omega$  on the boundary  $\partial\Omega$  (in view of (4.3.7)). In 3-d, the estimation of  $\Delta d_\Omega$  would not be sufficient to evaluate the values of both principal curvatures  $\kappa_1$  and  $\kappa_2$ : these could e.g. be computed from the eigenvalues of the Hessian matrix  $\nabla^2 d_\Omega$ .

Our first step towards estimating  $\Delta d_\Omega$  consists in calculating a  $\mathbb{P}_1$  interpolation  $g_h \in V_h \times V_h$  of the piecewise constant gradient  $\nabla d_h$  by solving the following variational problem:

$$\forall \psi_h \in V_h \times V_h, \int_D g_h \cdot \psi_h dx = \int_D \nabla d_h \cdot \psi_h dx. \quad (4.3.11)$$

The approximation of the divergence  $\text{div}(\nabla d_\Omega)$  is then calculated as the (piecewise constant) divergence of the reconstructed field  $g_h$ . Unfortunately, this procedure generally produces a very noisy approximation characterized by a lot of spurious oscillations when the mesh resolution increases (see Figure 4.5). In

order to overcome this difficulty, we calculate a regularization  $\kappa_h$  of this noisy estimation with a Laplace kernel, namely we solve:

$$\text{Find } \kappa_h \in V_h \text{ such that } \forall \psi_h \in V_h, \int_D (\gamma_h^2 \nabla \kappa_h \cdot \nabla \psi_h + \kappa_h \psi_h) dx = \int_D \text{div}(g_h) \psi_h dx, \quad (4.3.12)$$

where  $\gamma_h > 0$  is a regularization length scale (equal to  $3h_{\max}$  where  $h_{\max}$  is the maximum edge length in the mesh). This procedure yields satisfying results in practice (even with shapes  $\Omega$  characterized by discontinuous curvatures, up to some over smoothing near the discontinuities), although we do not have a proof of convergence of the approximation  $\kappa_h$  towards the exact function  $\Delta d_\Omega$ .

Let us illustrate our method by considering the example of an ellipse  $\Omega$  inside a square-shaped hold-all domain  $D$ : let  $D$  and  $\Omega$  be defined by

$$D = \{(x_1, x_2) \in \mathbb{R}^2 \mid |x_1| < 2 \text{ and } |x_2| < 2\}, \quad \Omega = \left\{ (x_1, x_2) \in D \mid \frac{x_1^2}{a^2} + \frac{x_2^2}{b^2} < 1 \right\}, \quad (4.3.13)$$

where  $a = 1.5$  and  $b = 1$ . The skeleton of  $\Omega$  is explicitly known in this case:  $\Sigma = \{(x_1, 0) \mid |x_1| < a - b^2/a\}$  and the curvature  $\kappa$  of  $\partial\Omega$  at a point  $y = (y_1, y_2)$  is given by  $\kappa(y) = ab/\gamma^3$  with  $\gamma = \sqrt{\frac{b^2}{a^2}y_1^2 + \frac{a^2}{b^2}y_2^2}$ . The difference between the exact curvature  $\kappa(y)$  of  $\partial\Omega$  and its reconstruction  $\kappa_h$  (at the boundary nodes discretizing  $\partial\Omega$ ) using both procedures (4.3.11) and (4.3.12) is represented in Figure 4.5.

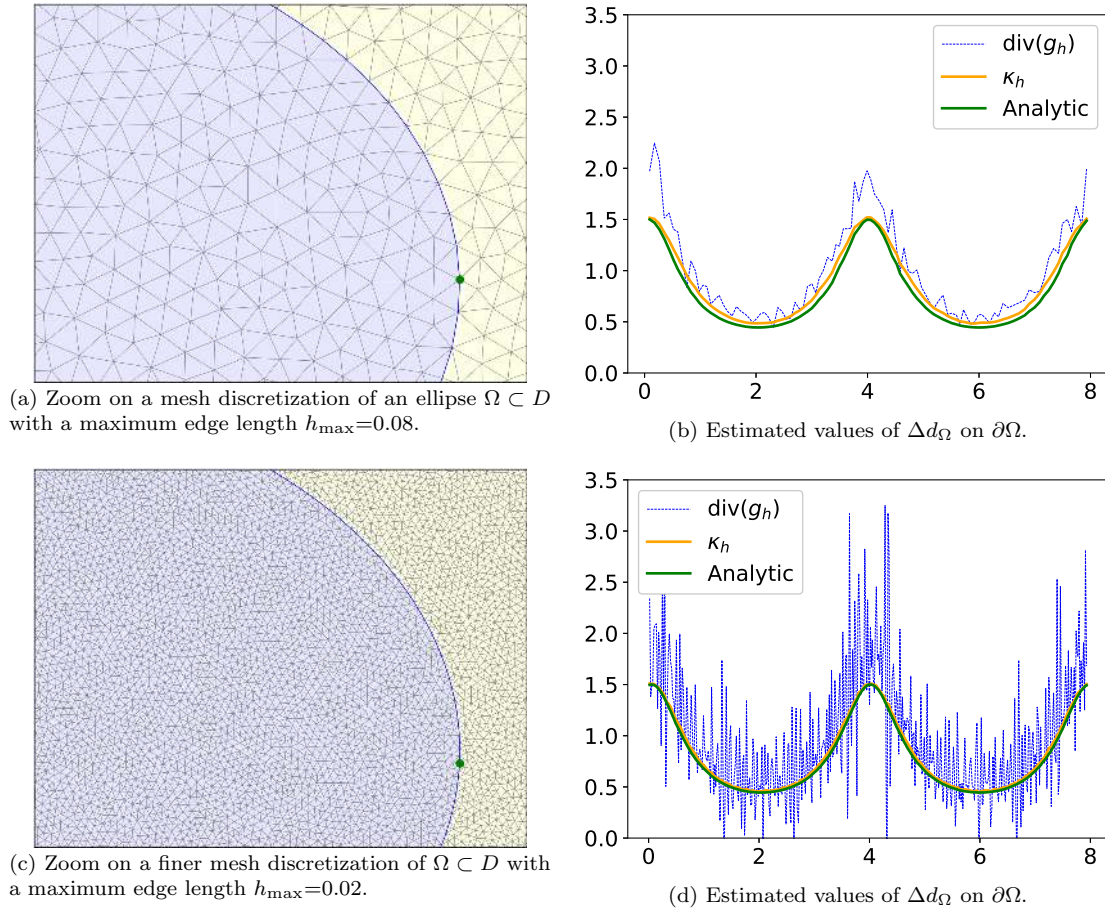


Figure 4.5: Estimation of  $\Delta d_\Omega$  on the mesh  $\mathcal{T}$  for the shape  $\Omega$  in (4.3.13) and for two different mesh resolutions. The  $x$  coordinate on the right-hand graphs represents the arc length coordinate on  $\partial\Omega$  when the starting point is the green reference point. Estimates of the mean curvature of  $\partial\Omega$ ,  $\text{div}(g_h)$  and  $\kappa_h$  (see (4.3.11) and (4.3.12) for the definitions) are compared to the analytical value  $\kappa(y)$ .

### 4.3.3 Admissible numerical weights built upon the signed distance function

In this section, we discuss the numerical resolution of the variational problem (4.2.5) in the shape optimization context (4.2.8) and (4.2.9) (see also proposition 4.6). The numerical setting is the same

as that of the previous [section 4.3.2](#): the hold-all domain  $D$  is equipped with a simplicial mesh  $\mathcal{T}$  (i.e. composed of triangles in 2-d, or tetrahedra in 3-d, although our method would work with other kinds of meshes) featuring a discretization of the domain  $\Omega$  as a submesh, and we rely on the Lagrange  $\mathbb{P}_1$  finite element method for the discretization of [\(4.2.5\)](#) on account of its robustness and ease of implementation. In other terms, the approximation  $u_h$  to the solution  $u \in V_\omega$  of [\(4.2.5\)](#) is sought in the space  $V_h$  of continuous piecewise linear functions on each simplex of  $\mathcal{T}$ .

We start by introducing some motivations for the use in [\(4.2.5\)](#) of weights  $\omega$  vanishing on the skeleton, and we provide (in [lemma 4.6](#) below) a formula for analytical and admissible weights satisfying approximately this property. We then discuss the issue of computing numerically these weights. Finally, we perform a few numerical experiments where we show that weights vanishing on the skeleton make  $\mathbb{P}_1$  finite elements able to capture discontinuous test functions across the skeleton, which allows to confirm numerically the latter motivations.

### Motivations for weights vanishing on the skeleton

As already mentioned, our final target is to calculate the trace [\(4.2.6\)](#) on  $\Gamma = \partial\Omega$  of  $u$ . In the continuous setting of [section 4.2](#), this trace does not depend on the choice of the weight  $\omega$  as long as it fulfills [\(H1\)](#) to [\(H3\)](#), so that in principle, any such weight could be used. However, when  $U$  is a cracked domain (typically, the skeleton  $\Sigma$  is a crack in the working domain, i.e.,  $U = D \setminus \bar{\Sigma}$ ) and the crack is not explicitly discretized in the mesh  $\mathcal{T}$ , then the most simple choice  $\omega = 1$  (which is an admissible weight on account of [lemma 4.6](#) below) might not work well in practice. Indeed and as we shall illustrate below in this section, test functions of [\(4.2.5\)](#) which belong to  $V_\omega$  with  $\omega = 1$  are in general discontinuous across  $\Sigma$  and not well captured by  $\mathbb{P}_1$  finite elements. In many shape optimization applications, discretizing the skeleton  $\Sigma$  of the current domain  $\Omega$  at every iteration of the optimization process or resorting to discontinuous finite elements (in order to be able to set  $\omega = 1$ ) is very inconvenient, for instance if working with fixed structured meshes as it is performed in many applications built on the level set method [\[26, 311\]](#).

Therefore, we shall be interested in determining weights  $\omega$  adapted to our commitment to use the space  $V_h$  of Lagrange  $\mathbb{P}_1$  finite elements for the resolution of [\(4.2.5\)](#) without the need for an accurate discretization of the skeleton  $\Sigma$  (alternative approaches could be to use discontinuous finite elements close to the skeleton, or to have a zero weight on the degrees of freedom corresponding to modes close to the skeleton and to remove the null space in the corresponding linear system, but they seemed to be more complicated to implement, at least to us). The weight  $\omega$  should be chosen in such a way that arbitrary functions  $v \in V_\omega$  are well approximated (in the norm  $\|\cdot\|_{V_\omega}$ ) by functions  $v_h \in V_h$ , as is reflected by the classical Céa's lemma (see e.g. [\[144\]](#)):

$$\|u - u_h\|_{V_\omega} \leq C \inf_{v_h \in V_h} \|u - v_h\|_{V_\omega}, \quad (4.3.14)$$

for a constant  $C > 0$  (which possibly depends on  $\omega$ ). The space  $V_h$  of Lagrange  $\mathbb{P}_1$  elements is a conformal finite element space in the sense that the inclusion  $V_h \subset V_\omega$  always holds (because functions of  $V_h$  are smoother than those of  $V_\omega$ ), however  $V_h$  may be “too small” to guarantee a correct approximation of discontinuous solutions  $u \in V_\omega$  in the sense of [\(4.3.14\)](#). Heuristically, and without looking for a very precise statement, these considerations call for the choice of a weight  $\omega$  almost vanishing on  $\bar{\Sigma}$ , so that the approximation error  $\|u - v_h\|_{V_\omega}$  in [\(4.3.14\)](#) attributes a lesser weight to a neighborhood of  $\bar{\Sigma}$  where  $u$  is expected to be discontinuous while the functions  $v_h \in V_h$  are continuous.

We now provide explicit candidates for weights  $\omega$  which fulfill the conditions [\(H1\)](#) to [\(H3\)](#) while taking small values near the skeleton, that we are going to use in our practical implementations.

**Lemma 4.6.** *For any real numbers  $q \geq 0, r \geq 0$ , the weight*

$$\omega = \frac{1}{1 + |d_\Omega|^q |\Delta d_\Omega|^r} \quad (4.3.15)$$

*satisfies the conditions [\(H1\)](#) to [\(H3\)](#) (this includes in particular the constant weight  $\omega = 1/2$  for  $q = r = 0$ ).*

*Proof.* At first, it follows readily from the definition that  $\omega$  belongs to  $\mathcal{C}^0(U, \mathbb{R}_+^*)$ , and is uniformly bounded on  $U$ , so that [\(H1\)](#) and [\(H2\)](#) are trivially satisfied. We then define  $\kappa_-(y)$  and  $\kappa_+(y)$  by:

$$\kappa_-(y) = \min(0, \min_i \kappa_i(y)), \quad \kappa_+(y) = \max(0, \max_i \kappa_i(y)), \quad (4.3.16)$$

as well as the corresponding multiplicities  $m_-(y)$  and  $m_+(y)$ :

$$m_{\pm}(y) = \text{Card}(\{i \in \{1, \dots, n-1\} \mid \kappa_i(y) = \kappa_{\pm}(y)\}).$$

Using formula (4.3.7) for  $\Delta d_{\Omega}$ , the weight  $\alpha = \omega \circ \eta |D\eta|$  is decomposed as

$$\begin{aligned} \alpha(y, s) &= \frac{\prod_{i=1}^{d-1} (1 + \kappa_i(y)s)}{1 + |s|^q \left| \sum_i \frac{\kappa_i(y)}{1 + \kappa_i(y)s} \right|^r} = \frac{\prod_{i=1}^{d-1} (1 + \kappa_i(y)s)^{r+1}}{\prod_{i=1}^{d-1} (1 + \kappa_i(y)s)^r + |s|^q \left| \sum_i \kappa_i(y) \prod_{j \neq i} (1 + \kappa_j(y)s) \right|^r} \\ &= f(y, s)g(y, s) \end{aligned} \quad (4.3.17)$$

where  $f$  and  $g$  are the following functions:

$$f(y, s) = (1 + s\kappa_-(y))^{m_-(y)(r+1)} \mathbf{1}_{s \geq 0}(s) + (1 + s\kappa_+(y))^{m_+(y)(r+1)} \mathbf{1}_{s < 0}(s), \quad (4.3.18)$$

$$g(y, s) = \frac{\prod_{i, \kappa_i \neq \kappa_-} (1 + \kappa_i(y)s)^{r+1} \mathbf{1}_{s \geq 0}(s) + \prod_{i, \kappa_i \neq \kappa_+} (1 + \kappa_i(y)s)^{r+1} \mathbf{1}_{s < 0}(s)}{\prod_{i=1}^{d-1} (1 + \kappa_i(y)s)^r + |s|^q \left| \sum_i \kappa_i(y) \prod_{j \neq i} (1 + \kappa_j(y)s) \right|^r}. \quad (4.3.19)$$

Then,  $f$  satisfies clearly the monotonicity condition (ii) in the statement of lemma 4.1 and  $g$  is a continuous function on each of the domains  $W_- = \{(y, s) \in W \mid s \leq 0\}$  and  $W_+ = \{(y, s) \in W \mid s \geq 0\}$ , that does not vanish on the compact sets  $\overline{W_-}$  and  $\overline{W_+}$ . The assumptions of lemma 4.1 are therefore fulfilled, so that  $\omega$  satisfies (H3). This terminates the proof.  $\square$

**Remark 4.13.** • When it comes to solving (4.2.5) by relying on proposition 4.5 with some data  $f \in L^\infty(U)$ , it is useful to observe that  $f$  belongs to  $L^2_{\omega^{-1}}(U)$  as soon as the weight  $\omega$  satisfies  $\omega^{-1} \in L^1(U)$  (see Remark 4.11), which is the case if it is of the form (4.3.15) with  $r < 2$ . In the following numerical experiments, we shall see however that using values for  $r$  which are larger than 2 still provides good results in practice: in general, taking higher values of  $q$  and  $r$  yields a faster decay of  $\omega$  near the skeleton.

- Taking  $r > 0$  ensures that the weight (4.3.15) will vanish at points  $x \in \overline{\Sigma}$  that are centers of curvatures (for which  $\Delta d_{\Omega}$  blows up). Taking  $q > 0$  allows to make sure that  $\omega = 1/2$  on  $\partial\Omega$  whatever the value of  $r$ , and to accentuate the decay of  $\omega$  near the skeleton. Note that in general,  $\omega$  will not vanish on points  $x \in \overline{\Sigma}$  that are not centers of curvatures. However, it still takes very small values on  $\overline{\Sigma}$  and it is convenient to use in the implementation. There is no unicity of weights appropriate for the numerical computation and variants can easily be imagined.
- For the most general setting where  $U$  is an arbitrary open set, the methodology of this section extends naturally by considering weights  $\omega$  which vanish on the cracked parts of  $\partial U$  that are not explicitly meshed.

### Numerical computations of the weight based on the Laplacian of the signed distance function

For our numerical applications below, we shall use the weights of lemma 4.6 in the definition of our variational formulation (4.2.5). This requires the computation of the Laplacian  $\Delta d_{\Omega}$  on the triangulated mesh  $\mathcal{T}$  based on the  $\mathbb{P}^1$  estimation of the signed distance function  $d_{\Omega}$ . For this purpose, we use the same regularization method outlined in section 4.3.2 for the computation of the numerical curvature.

Importantly, from proposition 4.5, the variational formulation (4.2.5) is rather insensitive to the choice of the weight  $\omega$ , and as a result the estimation of  $\Delta d_{\Omega}$  does not need to be very accurate as long as it takes large values near the skeleton (as we shall illustrate below in section 4.3.3). In contrast, the estimation of the principal curvatures  $\kappa_i$  for the direct method would *need* to be accurate. From a numerical standpoint, the formula (4.3.15) featuring  $\Delta d_{\Omega}$  at the denominator is convenient to obtain numerically vanishing weights near the skeleton (even if (4.3.15) truly vanish for  $r > 0$  at centers of curvatures). Indeed,  $\nabla d_{\Omega}$  is discontinuous across the skeleton, which should reflect in high numerical values of  $\Delta d_{\Omega}$  when computing numerically the divergence  $\text{div}(\nabla d_{\Omega})$  with the method of section 4.3.2.



### Assessing the choice of the weight $\omega$ when using a $\mathbb{P}_1$ discretization: generating numerical test functions constant along rays

In this section, we perform several numerical experiments about the influence of the choice of the weight function  $\omega$  in the resolution of the variational problem (4.2.5) using the Lagrange  $\mathbb{P}_1$  finite element method. With this perspective in mind, and in the 2-d numerical setting described in sections 4.3.2 and 4.3.3, we consider the issue of generating numerical functions  $v \in V_\omega$  which are constant along the normal rays to the shape  $\Omega$ . Namely, we solve the boundary-value problem

$$\begin{cases} \nabla d_\Omega \cdot \nabla v = 0 & \text{in } D \setminus \bar{\Sigma}, \\ v = v_0 & \text{on } \partial\Omega, \end{cases} \quad (4.3.20)$$

for given data  $v_0 \in L^2(\partial\Omega)$ .

Using the variational setting of section 4.2, we show that it is possible to obtain the solution  $v$  to (4.3.20) by solving a variational problem of the same nature of (4.2.5).

**Proposition 4.7.** *Let  $\omega \in C^0(D \setminus \bar{\Sigma}, \mathbb{R}_+^*)$  be a weight satisfying (H1) to (H3). There exists a unique solution  $v \in V_\omega$  to the following variational problem:*

$$\text{Find } v \in V_\omega \text{ such that } \forall w \in V_\omega, \int_{\partial\Omega} v w ds + \int_{D \setminus \bar{\Sigma}} \omega (\nabla d_\Omega \cdot \nabla v) (\nabla d_\Omega \cdot \nabla w) dx = \int_{\partial\Omega} v_0 w ds. \quad (4.3.21)$$

The solution  $v$  is independent of the choice of  $\omega$  as long as (H1) to (H3) are satisfied, and it is given by the formula:

$$v(x) = v_0(p_{\partial\Omega}(x)), \text{ a.e. } x \in D \setminus \bar{\Sigma}. \quad (4.3.22)$$

*Proof.* To see that (4.3.20) and (4.3.21) are equivalent, it is sufficient to take  $v \in V_\omega$  constant along rays in (4.3.21) as in the proof of proposition 4.5, which yields  $v = v_0$  on  $\partial\Omega$ , and then  $\nabla d_\Omega \cdot \nabla v = 0$ .  $\square$

The formulation (4.3.21) is to be compared to the so-called Galerkin Least Square formulation and SUPG methods for advection-reaction problems [144], with the difference that usual assumptions of uniformly bounded divergence (which do not hold in our applications) are replaced with the regularity assumptions of section 4.2.1.

**Remark 4.14.** The problem of building constant functions along normal rays of the form (4.3.20) may be solved on unstructured meshes owing to variants of the fast marching algorithm; see e.g. [97].

We now verify that a good approximation of the solution  $v$  to (4.3.20) (or more precisely (4.3.21)), which is in particular discontinuous across the skeleton  $\Sigma$  of  $\Omega$ , can be obtained either by truncating manually the skeleton from the computational mesh  $\mathcal{T}$  and taking the weight  $\omega = 1$ , or by selecting a weight  $\omega$  taking “small” values near  $\Sigma$ . For the latter experiment, we use the weight  $\omega = 1/(1 + |d_\Omega|^{3.5} |\Delta d_\Omega|^2)$  (see also remark 4.13 about the “small” values of  $\omega$  near  $\Sigma$ ). Of course, removing the skeleton from the mesh is not a straightforward task in full generality but it is performed here for the sake of comparison.

Let us consider again the ellipse example of (4.3.13). We consider two different computational meshes  $\mathcal{T}$  and  $\mathcal{T}'$ . The former is a triangular mesh of  $D$ , and the latter  $\mathcal{T}'$  is a triangular mesh of  $D \setminus \bar{\Sigma}$  (i.e. the skeleton  $\bar{\Sigma}$  has been manually removed). In both  $\mathcal{T}$  and  $\mathcal{T}'$ , the considered shape  $\Omega$  is explicitly discretized as a submesh; see Figure 4.6 for an illustration. The variational problem (4.3.21) is numerically solved for a boundary datum  $v_0 \in L^2(\partial\Omega)$  given by (see Figure 4.7d):

$$\forall (y_1, y_2) \in \partial\Omega, v_0(y_1, y_2) = \cos(3y_1)^2 + 20y_2, \quad (4.3.23)$$

and the computed finite element solution is plotted on Figure 4.7 in the following three situations.

- The mesh  $\mathcal{T}'$  is used, in which  $D$  and  $\Omega \subset D$  are meshed explicitly, and where a thin layer around  $\bar{\Sigma}$  has been manually removed (see Figure 4.6c). The solution  $v$  to (4.3.21) is computed using the constant weight  $\omega = 1$  and the result is displayed on Figure 4.7a. As expected, the fact that  $\bar{\Sigma}$  is absent from  $\mathcal{T}'$  allows the numerical solution  $v$  to have very different values on either sides of  $\bar{\Sigma}$ .
- The mesh  $\mathcal{T}$  of  $D$  (where  $\bar{\Sigma}$  has not been removed) is used (see Figure 4.6b), and  $v$  is computed with the constant weight  $\omega = 1$ ; the result is represented on Figure 4.7b. The formulation (4.3.21) proves numerically stable with the choice  $\omega = 1$ , but it tends to over smoothen the sharp discontinuities of  $v$  near the skeleton  $\Sigma$ , which results in a loss of accuracy for the extension problem (4.3.20).

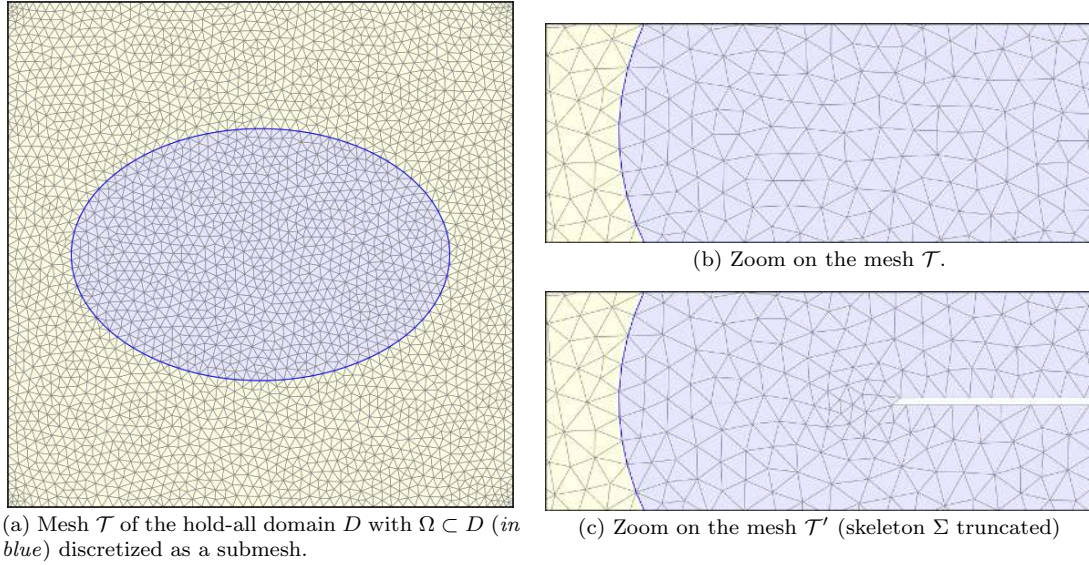


Figure 4.6: The two meshes  $\mathcal{T}$  and  $\mathcal{T}'$  used for the numerical example of section 4.3.3.

- The mesh  $\mathcal{T}$  of  $D$  is used again, but the solution  $v$  to (4.3.21) is now computed by using the weight  $\omega = 2/(1 + |d_\Omega|^{3.5} |\Delta d_\Omega|^2)$ ; the result is represented on Figure 4.7c. The obtained numerical solution is much closer to the expected result (4.3.22): the values of the solution function  $v$  look constant along the normal rays up to a small neighborhood of the skeleton (of the size of the mesh element size), where sharp variations are observed, as expected. The numerical procedure in this case seems therefore to achieve the same order of accuracy than in the experiment using the truncated mesh  $\mathcal{T}'$ .

#### 4.3.4 Numerical comparisons between the variational method and direct integration along rays

We now investigate the numerical evaluation of the function  $u \in L^2(\partial\Omega)$  in (4.2.6) for several 2-d academic configurations of domains  $D$ ,  $\Omega$  and functions  $f$ . In particular, we compare the evaluation of  $u$  obtained by direct integration along rays (*i.e.* implementing directly the formula (4.2.6)) to that obtained by solving the variational formulation (4.2.5) on meshes  $\mathcal{T}$  (resp.  $\mathcal{T}'$ ) of  $D$  in which  $\Omega$  is explicitly discretized and the skeleton  $\Sigma$  of  $\Omega$  is not removed (resp. is removed).

##### A domain with trivial skeleton: the case of a circle

We first consider the case where  $\Omega$  is a disk enclosed in a larger disk  $D$ :

$$D = \{(x_1, x_2) \in \mathbb{R}^2 \mid x_1^2 + x_2^2 < 4\}, \text{ and } \Omega = \{(x_1, x_2) \in \mathbb{R}^2 \mid x_1^2 + x_2^2 < 1\}; \quad (4.3.24)$$

see Figure 4.8. In this case, the skeleton  $\Sigma$  is reduced to the point 0. The considered function  $f$  is:

$$\forall x = (x_1, x_2) \in D, \quad f(x_1, x_2) = x_2,$$

which belongs to the finite element space  $V_h$ , and is therefore amenable to an exact integration when the travel procedure along rays of section 4.3.2 is used.

In this situation, the sought function  $u$ , given by (4.2.6), is known analytically; a calculation in polar coordinates indeed yields:

$$\forall (y_1, y_2) \in \partial\Omega, \quad u(y_1, y_2) = \frac{8}{3} y_2. \quad (4.3.25)$$

Comparisons are displayed on Figure 4.9 between the version of  $u$  obtained after direct integration along rays, and the numerical solutions of (4.2.5) obtained for various choices of weight functions  $\omega$  on the mesh  $\mathcal{T}$  of  $D$  (where  $\Sigma$  has not been removed). We observe a good match between the variational and the direct method. As expected, the solutions computed thanks to our variational method are less accurate when the constant weight  $\omega = 1$  is chosen. A significant increase in accuracy is achieved by

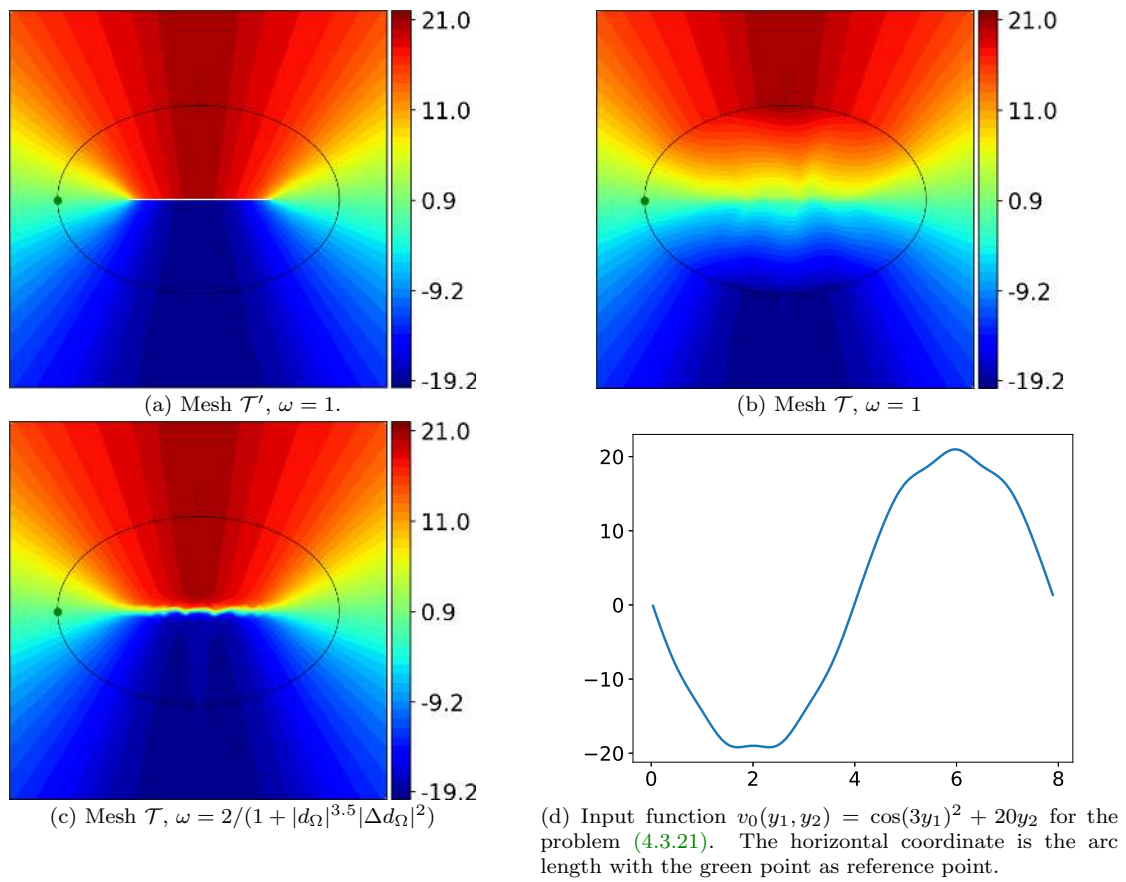


Figure 4.7: Numerical resolution of the problem (4.3.20) using the variational problem (4.3.21) for various weights, with or without removing the skeleton  $\Sigma$  from the computational mesh.

selecting a weight  $\omega$  vanishing at the center of the circle, i.e.  $\omega = 1/(1 + |d_\Omega|^{3.5}|\Delta d_\Omega|^{3.5})$ . Note that this weight does not fulfill the condition  $\omega \in L^1_{\omega^{-1}}(U)$  (see [remark 4.12](#)), but works well in numerical practice.

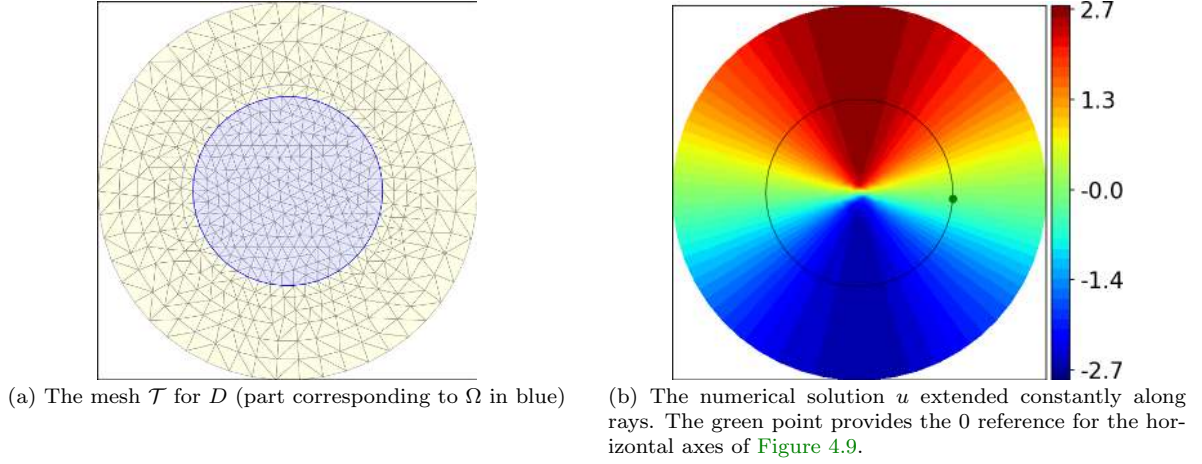


Figure 4.8: Setting of the disk test case of (4.3.24) in section 4.3.4.

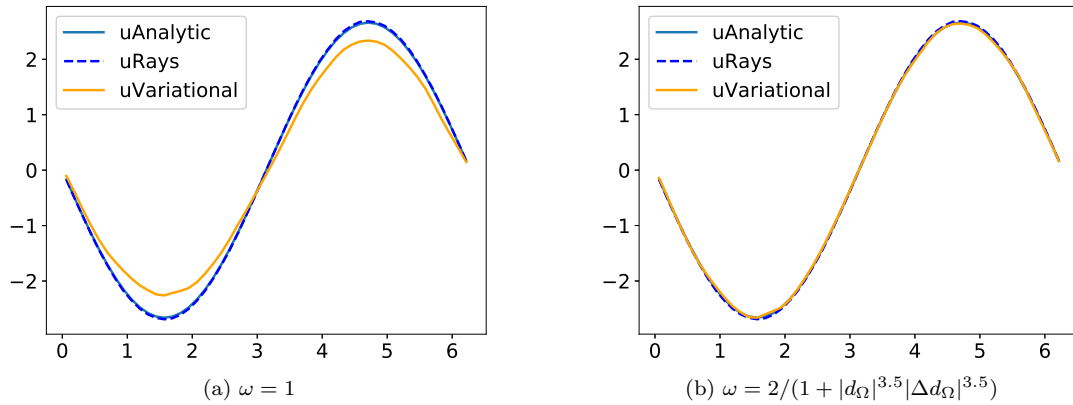


Figure 4.9: Comparison between direct integration along rays, analytical formula (4.1.3), and finite element solution of (4.2.5) for two choices of weights  $\omega$  for the example of the test case of (4.3.24) in section 4.3.4. In the legends of this figure and all those to follow, `uAnalytic`, `uRays` and `uVariational` refer respectively to the analytical value of  $u$ , its numerical estimation using (4.2.38), and the trace of the variational solution of (4.2.5).

### A $\mathcal{C}^2$ domain with non trivial skeleton

We take on the example (4.3.13) where  $\Omega$  is an ellipse inside a square-shaped hold-all domain  $D$ , reusing both meshes  $\mathcal{T}$  and  $\mathcal{T}'$  depicted on [Figure 4.6](#). To evaluate the influence of the skeleton on our numerical method, we compute the quantity (4.2.5) for a function  $f$  which is supported in  $\Omega$  only (see [Figure 4.10](#)):

$$f(x_1, x_2) = (1 + x_2)\mathbf{1}_\Omega(x_1, x_2). \quad (4.3.26)$$

The presence of the constant 1 in (4.3.26) avoids the “simplification” of having  $f$  vanishing on  $\Sigma$ . After an explicit calculation based on (4.2.6), the exact solution is given by:

$$\forall y = (y_1, y_2) \in \partial\Omega, \quad u(y) = (1 + y_2)\zeta_- - \frac{\zeta_-^2}{2} \left( \frac{y_2}{\zeta_-} + \kappa(1 + y_2) \right) + \frac{y_2}{3}\kappa\zeta_-^2, \quad (4.3.27)$$

with  $\kappa = \frac{ab}{\gamma^3}$ ,  $\zeta_- = \frac{\gamma b}{a}$ ,  $\gamma = \sqrt{\frac{b^2}{a^2}y_1^2 + \frac{a^2}{b^2}y_2^2}$ .

The numerical trace  $u$  (which is extended along rays using (4.3.21) to ease the visualization) obtained by solving the variational formulation (4.2.5) on  $\mathcal{T}$  or  $\mathcal{T}'$ , and for various choices of weights, is represented

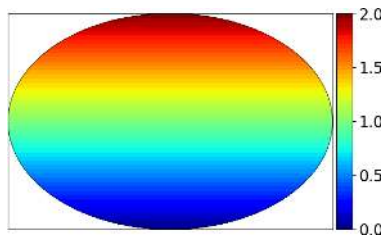


Figure 4.10: The function  $f$  for the identification problem (4.2.6) for the ellipse test case (4.3.13) of section 4.3.4.

on Figure 4.11. On Figure 4.12, we compare for each of these strategies the boundary values of the numerical trace  $u$  to that computed with (4.1.3) by direct integration along rays, and to the analytical expression (4.3.27).

We note that the solutions obtained by integrating along rays are generally characterized by small spurious oscillations. We explain this error by the fact that our criterion (4.3.9) detects the skeleton up to an error proportional to the mesh size (the actual location of the skeleton is not estimated within the last travelled triangle; see Figure 4.13). Such a procedure could be improved by using more accurate skeleton detection methods; see for instance the works [37, 45] in the field of computer graphics. We verify that the amplitude of these oscillations decreases with the size of the mesh, as shown on Figure 4.12d.

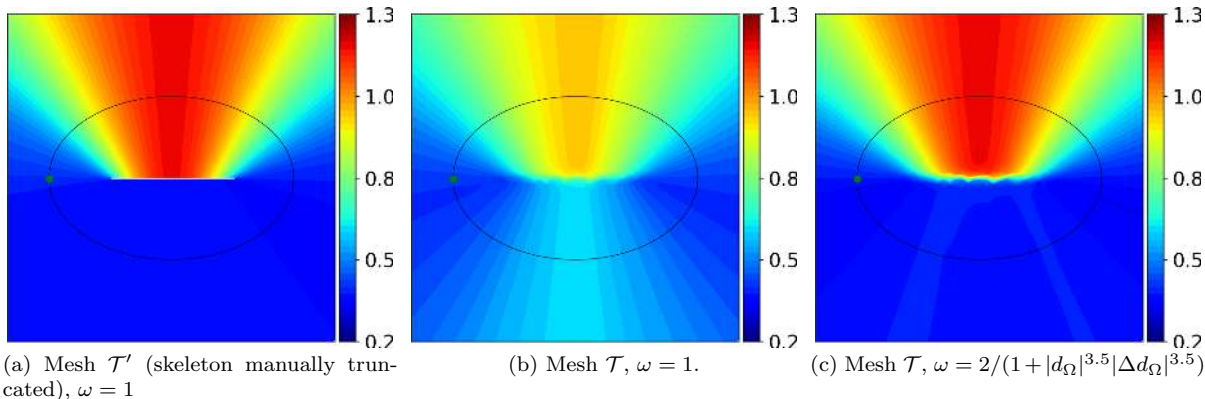


Figure 4.11: The numerical solution  $u$  of the variational problem (4.2.5) extended constantly along rays (for visualization purpose). Inaccuracies occur if the skeleton  $\Sigma$  is not removed from the mesh or if the weight  $\omega$  does not vanish in its vicinity.

### A $C^1$ domain, with discontinuous curvature

We now consider the case where  $\Omega$  is a stadium, i.e. the reunion of a rectangle and two half-disks. Define (see Figure 4.14)

$$D = \{(x_1, x_2) \in \mathbb{R}^2 \mid |x_1| \leq 2 \text{ and } |x_2| < 2\} \tag{4.3.28}$$

and

$$\Omega = \{(x_1, x_2) \in D \mid (|x_1| < 0.5 \text{ and } |x_2| < 0.5) \text{ or } (x_1 - 0.5)^2 + x_2^2 < 0.25 \text{ or } (x_1 + 0.5)^2 + x_2^2 < 0.25\}. \tag{4.3.29}$$

The domain  $\Omega$  is not of class  $C^2$ , and the curvature  $\kappa$  of the boundary  $\partial\Omega$  is discontinuous at the points where  $x_1 = \pm 0.5$ . Hence, this example does not fall into the admissible setting of proposition 4.6, and there is, in principle, no guarantee that our variational method based on (4.2.5) should still work.

In this example, the skeleton  $\Sigma$  of  $\Omega$  is explicitly given by  $\Sigma = \{(x_1, 0) \mid |x_1| \leq 0.5\}$ . We calculate the function  $u$  in (4.2.6) associated to the datum function  $f$  defined by (4.3.26); the analytical solution  $u$  for this problem is given by

$$u(y_1, y_2) = \begin{cases} \frac{1}{2} + \frac{1}{8} \text{sign}(y_1) & \text{if } |y_1| < 0.5, \\ \frac{1}{4} + \frac{y_2}{6} & \text{if } |y_1| > 0.5. \end{cases}$$

In particular,  $u$  is ill-defined at the points  $x \in \partial\Omega$  where the curvature is discontinuous.

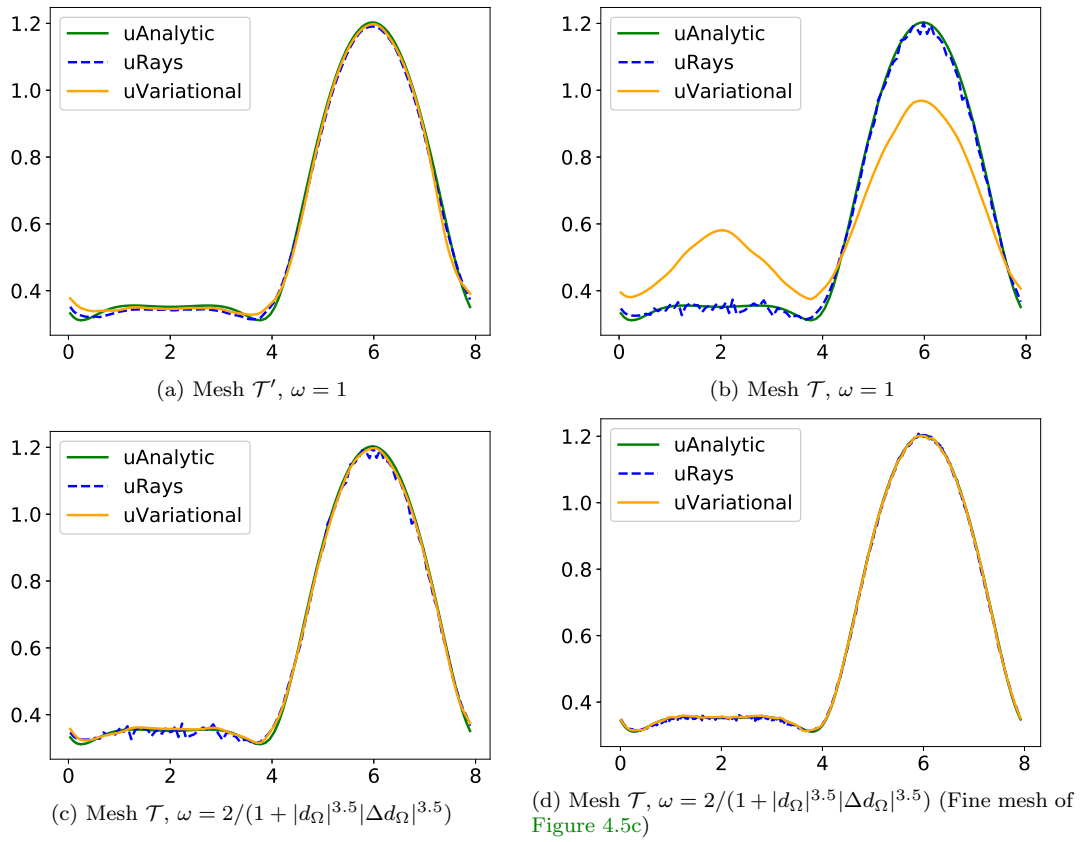


Figure 4.12: Comparison between the results of a direct integration along rays and our variational method for the ellipse example of section 4.3.4.

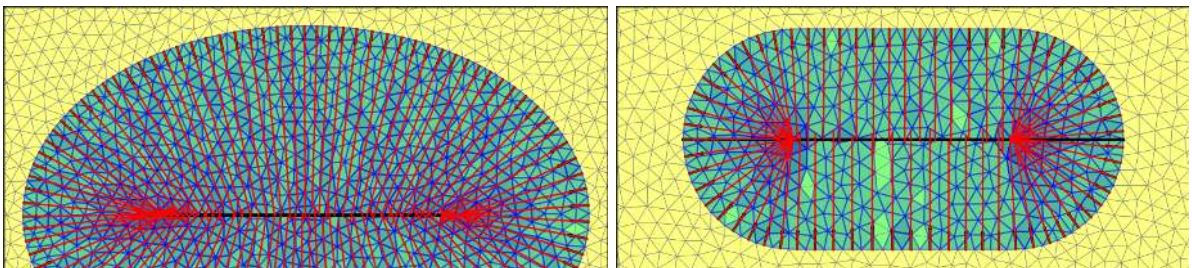


Figure 4.13: Numerical rays for the ellipse and stadium examples of section 4.3.4. The inaccuracy of the direct integration along rays is related to that of the skeleton detection within the last triangle. The triangle paths are colored in transparent blue (the darker, the more often visited); depicting that some mesh triangles might not be crossed by all the rays or some triangles may be visited more often than others.

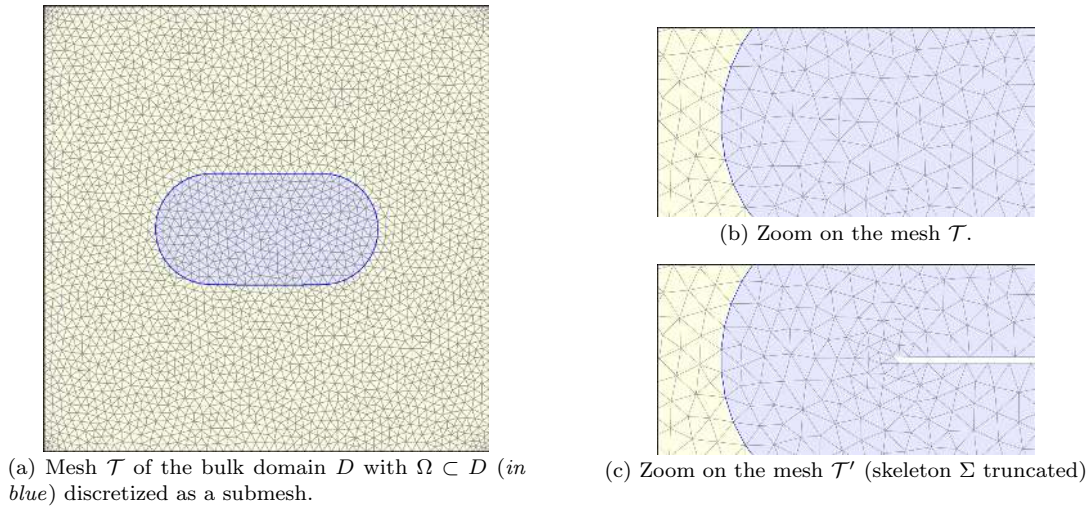


Figure 4.14: Meshes used for the stadium example of section 4.3.4.

Both methods for the numerical evaluation of  $u$  (i.e. the direct integration along rays and our variational method) are applied, and the results are represented on Figure 4.15. Even though this example does not fit into the admissible setting of section 4.2, the results indicate that our numerical method still works, up to over-smoothing inaccuracies of the computed solution  $u$  near the discontinuities of the curvature  $\kappa$ .

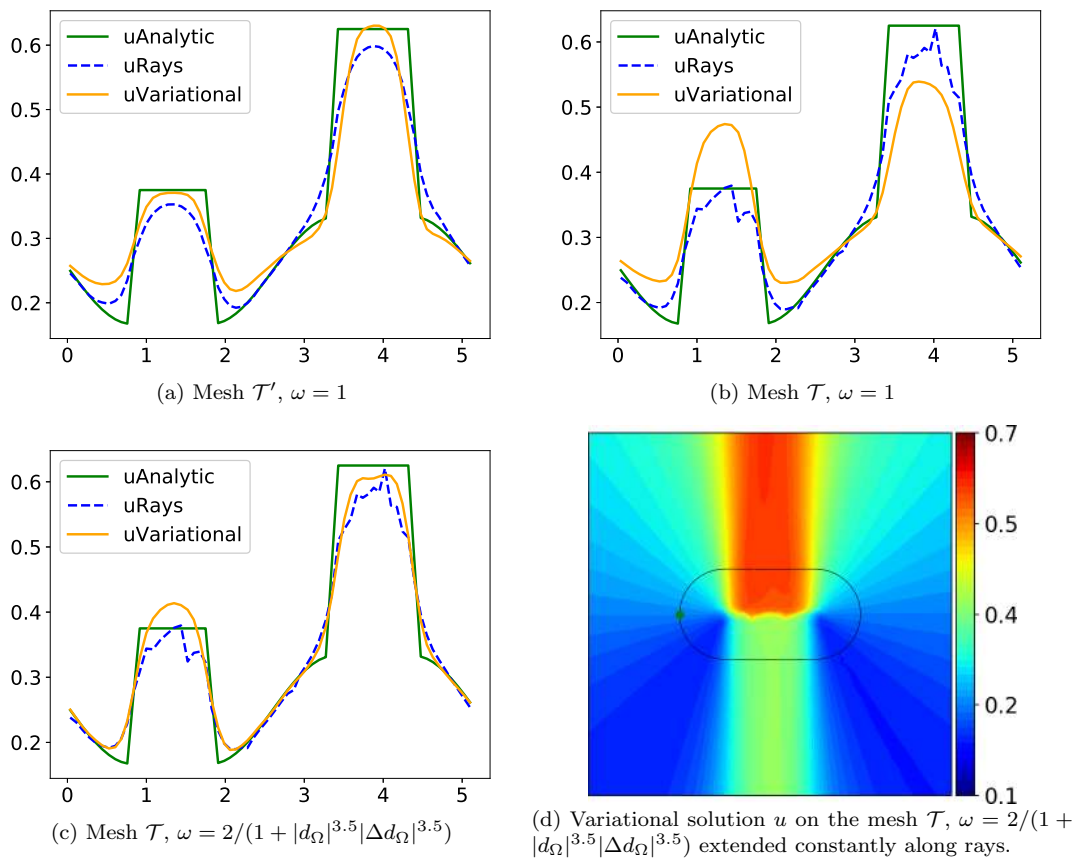


Figure 4.15: Comparison between the direct integration along rays and our variational method for the stadium example of section 4.3.4.

### A test case involving a Lipschitz domain $\Omega$ with angular corners

In this last example, we consider the following situation:

$$D = \{(x_1, x_2) \in \mathbb{R}^2 \mid |x_1| \leq 2 \text{ and } |y| \leq 2\}, \Omega = \{(x_1, x_2) \in D \mid \max(x_1^2/2 - x_2 - 1, x_2) \leq 0\}, \quad (4.3.30)$$

together with the datum function  $f$ :

$$f(x_1, x_2) = \cos(6x_1)^2,$$

as illustrated on [Figure 4.16](#). Again, the theoretical framework of [section 4.2](#) does not apply to the present situation because:

- The vector field  $\nabla d_\Omega$  is not of class  $\mathcal{C}^1$  on  $U$  and is not  $U$ -filling: the reunion of all the rays emerging from points  $y \in \partial\Omega$  does not cover the whole set  $U = D \setminus \bar{\Sigma}$ ;
- The flow  $\eta$  of this vector field is not of class  $\mathcal{C}^1$ , because the curvature  $\kappa$  of the boundary (and even the normal vector) is discontinuous at the corners. Numerically,  $|\Delta d_\Omega|$  blows up in the whole area filled by the skeleton and the set of points that are not covered by normal rays.

Still, the quantity  $u$  given by [\(4.2.6\)](#) can be calculated analytically wherever it makes sense. A few elementary, albeit technical computations yield:

$$u(y_1, y_2) = \begin{cases} \cos(6y_1)^2(2 - \lambda_-(s^{-1}(y_1))) & \text{if } y_2 \geq 0, \\ \phi(\lambda_+(y_1)) - \phi(\lambda_-(y_1)) & \text{if } y_2 \leq 0 \end{cases} \quad (4.3.31)$$

where

$$\lambda_-(t) = \frac{t^2/2 - 1}{1 + 1/\sqrt{1+t^2}}, \quad \lambda_+(t) = \min((2/|t| - 1), t^2/2)\sqrt{1+t^2}, \quad (4.3.32)$$

$s^{-1}$  is the reciprocal function of  $t \mapsto s(t) = \left(1 + \frac{\lambda_-(t)}{\sqrt{1+t^2}}\right)t$  and  $\phi$  is a primitive function of  $\lambda \mapsto \cos(6t(1 + \lambda/\sqrt{1+t^2}))^2(1 + \lambda/(1+t^2)^{3/2})$ .

A numerical approximation to this function  $u$  is computed by using either a direct integration along rays or our variational method based on [\(4.2.5\)](#), on a single mesh  $\mathcal{T}$  where the skeleton  $\Sigma$  is not removed, and using three possible choices of the weight  $\omega$ . The results are represented on [Figure 4.17](#). In all cases, the observed numerical inaccuracies, characterized by very high values, are concentrated near the angular corners of  $\Omega$ , while the method yields satisfying accuracy on the remaining smooth parts of  $\partial\Omega$ . Again, we observe the poor accuracy associated to the choice of a constant weight  $\omega = 1$ . For this example, a uniformly small weight  $\omega = 1e^{-10}$  seems to yield a better accuracy than the previous choice  $\omega = 2/(1 + |d_\Omega|^{3.5}|\Delta d_\Omega|^{3.5})$ .

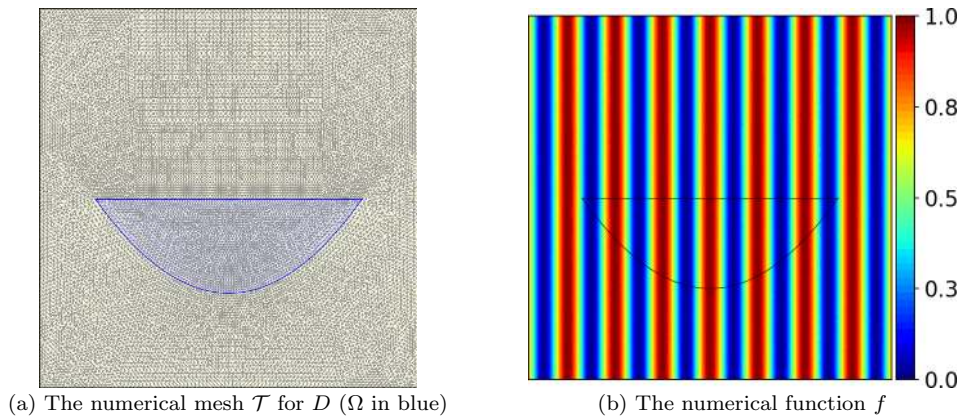


Figure 4.16: Setting for the non smooth example of [section 4.3.4](#).



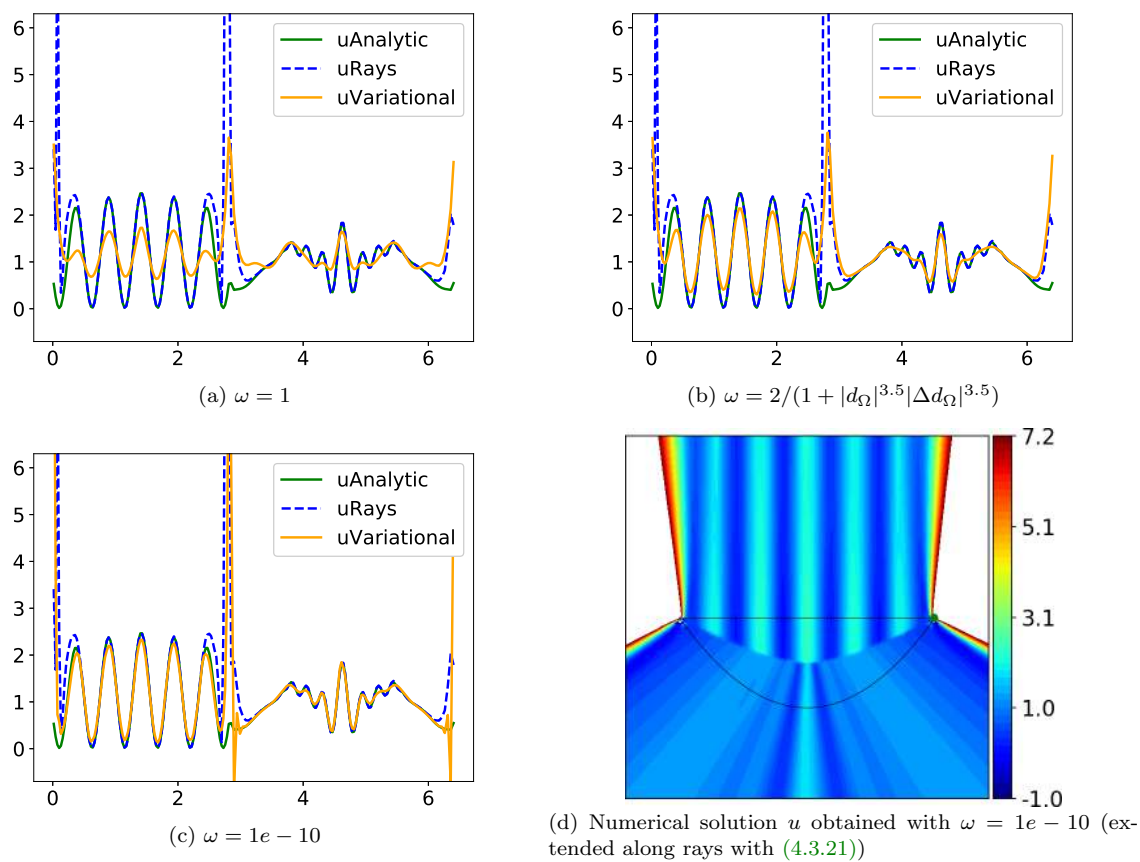


Figure 4.17: Comparison between direct integration along normal rays and our variational method for the non smooth example of section 4.3.4.

#### 4.4 APPLICATIONS TO MAXIMUM AND MINIMUM THICKNESS CONSTRAINTS IN SHAPE OPTIMIZATION

We now show how the general method proposed in [section 4.2](#) and the inferred numerical methods in [section 4.3](#) allow to efficiently implement geometric constraints, namely maximum and minimum thickness constraints in shape optimization, as the original motivation for our work. Here no comparison is made between our new variational approach and the previous method (using direct integration along rays) to evaluate the shape derivatives of the related shape functionals. The optimized shapes and topologies resulting from the variational method are very similar to those obtained in the previous works [[30](#), [234](#)]. There is no clear gain in computational time but there is a very substantial simplification of the implementation (which would be tremendous in 3-d).

##### 4.4.1 Shape optimization setting for linearly elastic structures

In this section, we consider situations where only the mechanical displacement  $\mathbf{u}$  of a structure comes into play. For compatibility with the notation of the present chapter, the domain to be optimized associated to the solid structure is denoted  $\Omega \subset D$  (and not  $\Omega_s$  as in [chapter 2](#)). Any such shape  $\Omega$  is clamped on a part  $\Gamma_D$  of its boundary, and traction loads  $\mathbf{g} \in L^2(\Gamma_N)$  are applied on a disjoint region of  $\partial\Omega$ ; the complement  $\Gamma := \partial\Omega \setminus (\overline{\Gamma_D} \cup \overline{\Gamma_N})$  is traction-free and body forces are omitted for simplicity;  $\Gamma$  is the only region of the boundary  $\partial\Omega$  which is subject to optimization.

In this situation, the displacement  $\mathbf{u}_\Omega$  of the shape is characterized as the unique solution in  $H^1(\Omega, \mathbb{R}^d)$  to the linearized elasticity system:

$$\begin{cases} -\operatorname{div}(Ae(\mathbf{u}_\Omega)) = 0 & \text{in } \Omega, \\ Ae(\mathbf{u}_\Omega) \cdot \mathbf{n} = \mathbf{g} & \text{on } \Gamma_N, \\ Ae(\mathbf{u}_\Omega) \cdot \mathbf{n} = 0 & \text{on } \Gamma, \\ \mathbf{u}_\Omega = 0 & \text{on } \Gamma_D, \end{cases} \quad (4.4.1)$$

where  $e(\mathbf{u}) := \frac{1}{2}(\nabla\mathbf{u} + \nabla\mathbf{u}^T)$  is the strain tensor associated to the displacement  $\mathbf{u}$  and  $A$  is the Hooke's law, defined for any symmetric  $n \times n$  matrix by  $Ae(\mathbf{u}) = 2\mu e(\mathbf{u}) + \lambda \operatorname{Tr}(e(\mathbf{u}))I$  involving the Lamé coefficients  $\lambda, \mu$  which characterize the physical properties of the constituent material.

In this context, we consider structural optimization problems of the form

$$\min_{\Omega \subset D} J(\Omega), \quad \text{s.t. } P(\Omega) \leq 0, \quad (4.4.2)$$

where  $J(\Omega)$  is a performance criterion, which will typically be the volume  $\operatorname{Vol}(\Omega)$  or the compliance  $C(\Omega)$  of shapes:

$$\operatorname{Vol}(\Omega) = \int_{\Omega} dx, \quad C(\Omega) = \int_{\Omega} Ae(\mathbf{u}_\Omega) : e(\mathbf{u}_\Omega) dx, \quad (4.4.3)$$

and  $P(\Omega)$  is a geometric constraint

$$P(\Omega) = \int_D j(d_\Omega) dx, \quad (4.4.4)$$

involving the signed distance function  $d_\Omega$  to  $\Omega$  and a given smooth function  $j : \mathbb{R} \rightarrow \mathbb{R}$ .

The shapes derivatives of the volume and compliance functions are classically given by (see [chapter 2](#)):

$$\operatorname{DVol}(\Omega)(\boldsymbol{\theta}) = \int_{\Gamma} \boldsymbol{\theta} \cdot \mathbf{n} dy, \quad \text{and } \operatorname{DC}(\Omega)(\boldsymbol{\theta}) = - \int_{\Gamma} Ae(\mathbf{u}_\Omega) : e(\mathbf{u}_\Omega)(\boldsymbol{\theta} \cdot \mathbf{n}) dy. \quad (4.4.5)$$

where we recall that  $dy$  stands for the surface measure on  $\Gamma$ . The shape derivative of the geometric constraint  $P(\Omega)$  in [\(4.4.4\)](#) is given by (see the reminders of [chapter 1, section 1.3.2](#))

$$\operatorname{DP}(\Omega)(\boldsymbol{\theta}) = \int_{D \setminus \overline{\Sigma}} j'(d_\Omega(x)) d'_\Omega(\boldsymbol{\theta})(x) dx = \int_{\Gamma} u \boldsymbol{\theta} \cdot \mathbf{n} dy. \quad (4.4.6)$$

The ‘‘Eulerian’’ derivative  $d'_\Omega(\boldsymbol{\theta})$  of the signed distance function  $d_\Omega$  is defined on  $U = D \setminus \overline{\Sigma}$  by:

$$\forall x \in D \setminus \overline{\Sigma}, \quad d'_\Omega(\boldsymbol{\theta})(x) = -\boldsymbol{\theta}(p_{\partial\Omega}(x)) \cdot \mathbf{n}(p_{\partial\Omega}(x)), \quad (4.4.7)$$

and the scalar function  $u : \Gamma \rightarrow \mathbb{R}$  reads:

$$u(y) = - \int_{s \in \text{ray}(y)} j'(d_\Omega(s)) \prod_{1 \leq i \leq n-1} (1 + \kappa_i(y)d_\Omega(s)) ds,$$

as follows from an application of the coarea formula with the help of the material recalled in [chapter 1, section 1.3.2](#).

Using the conclusions of [section 4.2](#), the function  $u$  in (4.4.6) may be conveniently evaluated by solving the variational problem:

$$\text{Find } u \in V_\omega \text{ such that } \forall v \in V_\omega, \int_{\partial\Omega} uvdy + \int_{D \setminus \bar{\Sigma}} \omega(\nabla d_\Omega \cdot \nabla u)(\nabla d_\Omega \cdot \nabla v)dx = - \int_{D \setminus \bar{\Sigma}} j'(d_\Omega(x))vdx, \quad (4.4.8)$$

for a suitable weight  $\omega$  satisfying (H1) to (H3), as discussed in [section 4.3](#). Note how easily (4.4.6) is retrieved by taking the test function  $v = -d'_\Omega(\boldsymbol{\theta})$  in (4.4.8). Since  $\boldsymbol{\theta} \cdot \mathbf{n} \in L^\infty(\partial\Omega)$ , the derivative given by (4.4.7) is indeed an admissible test function  $d'_\Omega(\boldsymbol{\theta}) \in V_\omega$ , for it belongs to  $L^\infty(D \setminus \bar{\Sigma})$  and is constant along normal rays.

In our numerical implementation, we set  $\omega = 1/(1 + 100|d_\Omega \Delta d_\Omega|^{3.5})$  in order to solve (4.4.8) with  $\mathbb{P}_1$  finite elements, the constant 100 being selected to increase the slope of the weight near the skeleton.

#### 4.4.2 Shape optimization under a maximum thickness constraint

We first consider structural optimization problems featuring a maximum thickness constraint. Following the work in [30, 234], a shape  $\Omega \subset D$  is said to have maximum thickness lower than  $d_{\max} > 0$  provided:

$$\forall x \in \Omega, d_\Omega(x) \geq -d_{\max}/2. \quad (4.4.9)$$

Loosely speaking, this amounts to saying that the skeleton  $\Sigma$  of  $\Omega$  lies at a distance of at most  $d_{\max}$  from the boundary  $\partial\Omega$ . Following closely [30, 234], the pointwise constraint (4.4.9) is relaxed into a single integral constraint formulated in terms of the following penalty functional  $P_{\text{MaxT}}(\Omega)$ :

$$P_{\text{MaxT}}(\Omega) \leq \frac{d_{\max}}{2}, \text{ where } P_{\text{MaxT}}(\Omega) := \left( \frac{\int_{\Omega} h(d_\Omega) d_\Omega^2 dx}{\int_{\Omega} h(d_\Omega) dx} \right)^{1/2}, \quad (4.4.10)$$

and  $h$  is a regularized Heaviside function centered at  $d_{\max}/2$ :

$$\forall x \in \mathbb{R}, h(x) = \frac{1}{2} \left( 1 + \tanh \left( \frac{x - d_{\max}/2}{\alpha_f d_{\max}/2} \right) \right). \quad (4.4.11)$$

The parameter  $\alpha_f$  in (4.4.11) tunes the level of regularization; in our context, it is set to  $\alpha_f = \frac{4h_{\max}}{5d_{\max}}$ , where  $h_{\max}$  is the maximum edge length of the computational mesh of  $D$ .

A simple calculation yields the shape derivative of  $P_{\text{MaxT}}(\Omega)$ :

$$\begin{aligned} DP_{\text{MaxT}}(\Omega)(\boldsymbol{\theta}) = & - \frac{1}{2P_{\text{MaxT}}(\Omega)} \frac{\int_{\Omega} h(d_\Omega) d_\Omega^2 dx}{\left( \int_{\Omega} h(d_\Omega) dx \right)^2} h(0) \int_{\Gamma} \boldsymbol{\theta} \cdot \mathbf{n} dy \\ & + \frac{1}{2P_{\text{MaxT}}(\Omega)} \int_{\Omega} \left[ \frac{h'(d_\Omega) d_\Omega^2 + 2h(d_\Omega) d_\Omega}{\int_{\Omega} h(d_\Omega) dz} - \frac{\int_{\Omega} h(d_\Omega) d_\Omega^2 dz}{\left( \int_{\Omega} h(d_\Omega) dz \right)^2} h'(d_\Omega) \right] d'_\Omega(\boldsymbol{\theta}) dx, \end{aligned} \quad (4.4.12)$$

where the term involving  $d'_\Omega(\boldsymbol{\theta})$  is then computed using the variational formulation (4.4.8).

In the forthcoming examples of [sections 4.4.2 and 4.4.2](#), we solve the optimization problem of minimizing the volume  $\text{Vol}(\Omega)$  of the structure while imposing that the compliance  $C(\Omega)$  do not exceed a given threshold  $g_{\max}$ , as well as the maximum thickness constraint (4.4.10), namely:

$$\begin{aligned} \min \quad & \text{Vol}(\Omega) \\ \text{s.t.} \quad & \begin{cases} C(\Omega) \leq g_{\max} \\ P_{\text{MaxT}}(\Omega) \leq \frac{d_{\max}}{2}. \end{cases} \end{aligned} \quad (4.4.13)$$

### Optimization of the shape of a two-dimensional arch

Our first test-case reproduces that considered in §8.1.1 of [30], whose setting is displayed on Figure 4.18a: in a square-shaped domain  $D$ , the considered shapes are clamped on both their bottom left and bottom right corners, and a vertical load  $\mathbf{g} = (0, -20)$  is applied at the middle of their bottom side.

Starting from an initial shape arbitrarily perforated with several holes, the optimization problem (4.4.13) is solved with and without including the maximum thickness constraint, using the numerical values  $g_{\max} = 7.00$  and  $d_{\max} = 0.12$  when they are relevant.

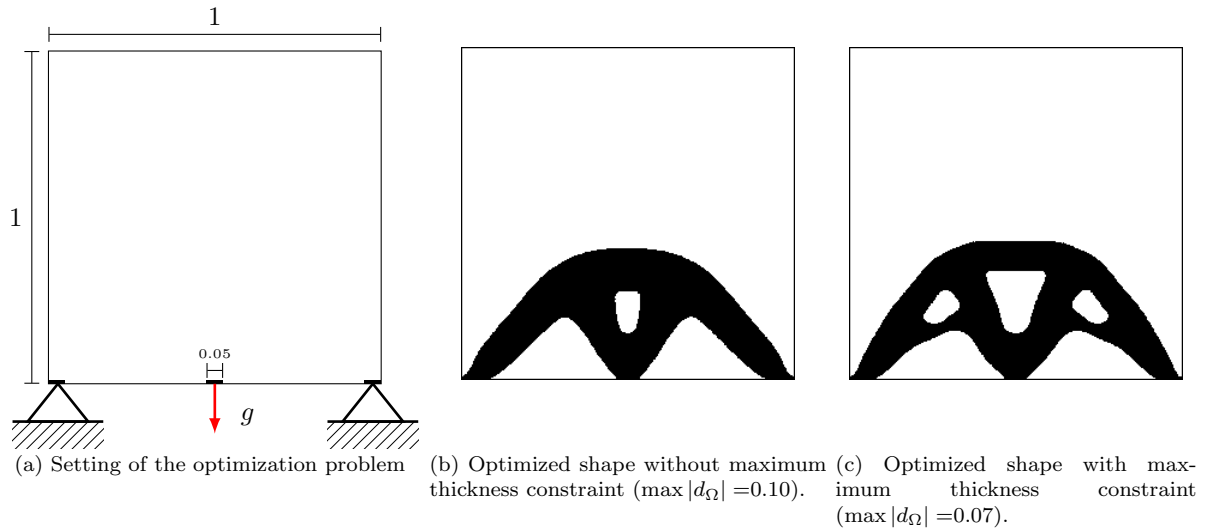


Figure 4.18: Physical setting and optimized shapes obtained in the two-dimensional arch optimization test-case of section 4.4.2.

The resulting optimized shapes in both cases are displayed on Figs. 4.18b and 4.18c. Several intermediate shapes as well as the convergence histories of the computation are displayed on Figs. 4.19 to 4.21. The obtained shapes are quite analogous to those obtained by [30, 234] where the calculation of the shape derivative of  $P_{\text{MaxT}}(\Omega)$  relied on a direct numerical integration along rays.

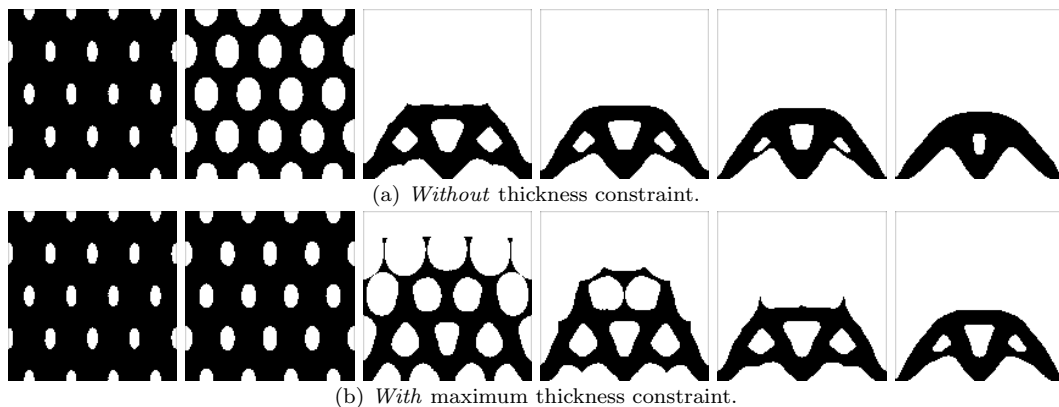


Figure 4.19: Iterations 0, 5, 40, 65, 100 and 200 in the shape optimization example of a 2-d arch of section 4.4.2.

### Optimization of a two-dimensional MBB-Beam

We consider now the classical MBB-Beam test-case depicted on Figure 4.22a: in a box  $D$  with dimensions  $3 \times 1$ , a material shape  $\Omega$  is constrained to no horizontal motion on the left boundary, and to no vertical motion on the bottom right corner. A vertical load  $\mathbf{g} = (0, -10)$  is applied on the top left corner, and the optimization problem (4.4.13) is considered again, with the numerical values for the thresholds  $g_{\max} = 30.00$  and  $d_{\max} = 0.16$ .

The resulting optimized shapes with and without including the maximum thickness constraint in (4.4.13) are represented on Figs. 4.22b and 4.22c and the convergence histories of the computation are

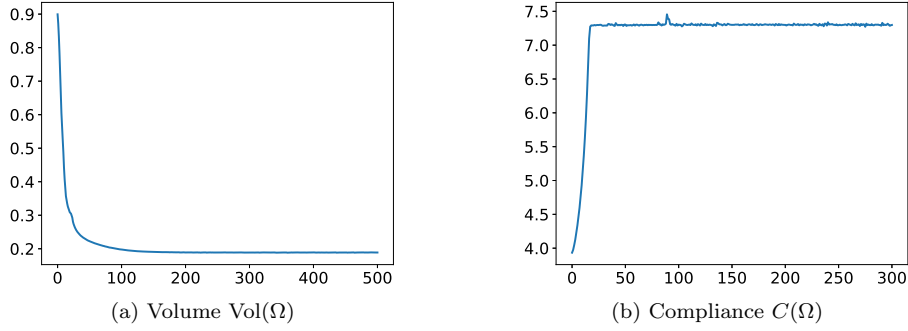


Figure 4.20: Optimization histories for the optimization problem (4.4.13) *without* maximum thickness constraint. Final values: volume  $\text{Vol}(\Omega) = 0.19$ , and compliance  $C(\Omega) = 7.00$ .

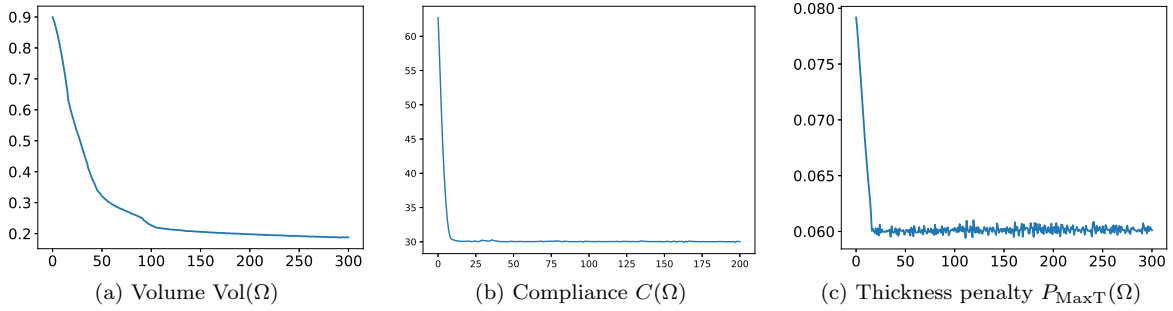


Figure 4.21: Optimization histories for the optimization problem (4.4.13) *with* maximum thickness constraint. Final values: volume  $\text{Vol}(\Omega) = 0.19$ , compliance  $C(\Omega) = 7.30$ , and penalty  $P_{\text{MaxT}}(\Omega) = 0.06$ .

shown on Figs. 4.23 to 4.24.

### 4.4.3 Shape optimization examples under a minimum thickness constraint

We now turn to the implementation of a minimum thickness constraint thanks to our variational method. Following [30, 234], we say that a shape  $\Omega$  has minimum thickness greater than  $d_{\min}$  if

$$\forall y \in \partial\Omega, \zeta_-(y) < -d_{\min}/2. \quad (4.4.14)$$

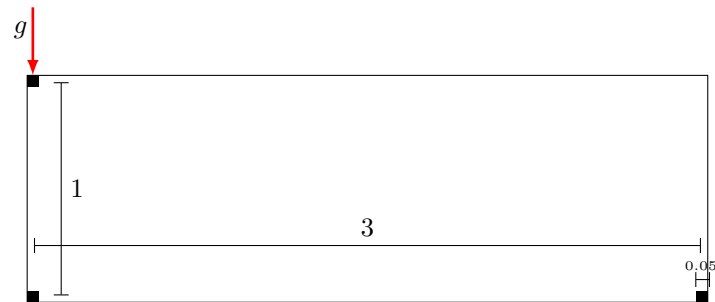
In other words, the boundary  $\partial\Omega$  is at a minimum distance  $d_{\min}/2$  of the part  $\Sigma \cap \Omega$  of the skeleton located inside the shape.

Enforcing a minimum thickness as a hard constraint (*i.e.* rather than a penalty term in the objective function) is by no means a straightforward task, because:

1. Our definition (4.4.14) of minimum thickness involve the distance to the skeleton  $\zeta_-$ , which is not differentiable with respect to the shape,
2. It is not clear how to formulate (4.4.14) by mean of a penalty functional such as (4.4.10) to penalize localizations on the shape that do not meet the thickness requirement,
3. Even if we were able to enforce the constraint at each iteration, such would prevent topology changes to occur naturally, which would be an issue in numerical practice.

We propose in the following a more flexible setting to enforce such a requirement in a structural optimization problem. Elaborating on ideas proposed in [234, 88, 92, 222], the minimum thickness requirement is implemented in the objective function rather than in the constraints: we minimize a penalty functional  $P_{\text{MinT}}(\Omega)$  for the minimum thickness under constraints on the volume and compliance of shapes, *i.e.* we solve:

$$\begin{aligned} \min \quad & P_{\text{MinT}}(\Omega) \\ \text{s.t.} \quad & \begin{cases} C(\Omega) \leq g_{\max} \\ \text{Vol}(\Omega) \leq V_{\max}. \end{cases} \end{aligned} \quad (4.4.15)$$



(a) Setting of the optimization problem. Small black areas correspond to non optimizable parts of the design domain.

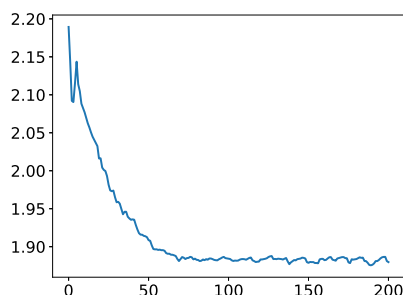


(b) Optimized shape without maximum thickness constraint ( $\max d_{\Omega} = 0.36$ ).

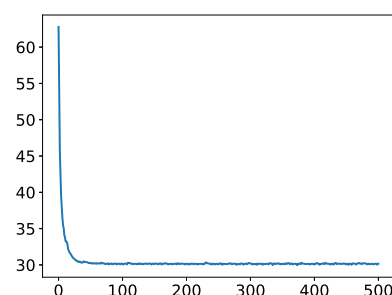


(c) Optimized shape with maximum thickness constraint.

Figure 4.22: Physical setting and obtained optimized shapes in the 2-d MBB beam optimization test-case of section 4.4.2 ( $\max |d_{\Omega}| = 0.26$ ).



(a) Volume  $\text{Vol}(\Omega)$



(b) Compliance  $C(\Omega)$

Figure 4.23: Optimization histories for the shape optimization problem (4.4.13) of the 2-d MBB-Beam of section 4.4.2 *without* including a maximum thickness constraint. Final values: volume  $\text{Vol}(\Omega) = 1.88$ , and compliance  $C(\Omega) = 29.96$ .

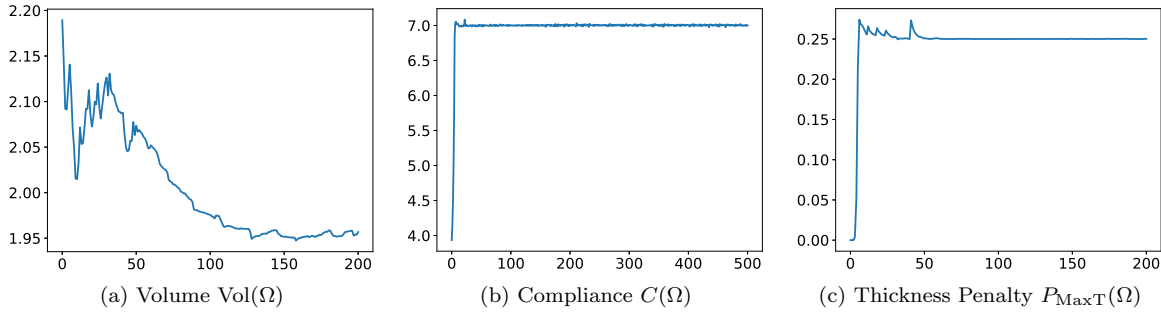


Figure 4.24: Optimization histories for the shape optimization problem (4.4.13) of the 2-d MBB-Beam of section 4.4.2 with maximum thickness constraint. Final values:  $\text{Vol}(\Omega) = 1.96$ , compliance:  $C(\Omega) = 30.03$ , and penalty  $P_{\text{MaxT}}(\Omega) = 0.25$ .

This strategy is expected to work because our optimization algorithm is designed to satisfy violated constraints *first*, before then attempting to reduce the objective function while maintaining the constraints respected.

The penalty functional  $P_{\text{MinT}}(\Omega)$  we considered in our case when solving (4.4.15) is taken from [234] and is specially designed to have a zero derivative when the constraint (4.4.14) is satisfied:

$$P_{\text{MinT}}(\Omega) = - \int_{\Omega} d_{\Omega}^2 \max(d_{\Omega} + d_{\text{min}}/2, 0)^2 dx. \quad (4.4.16)$$

The shape derivative of  $P_{\text{MinT}}(\Omega)$  is given by

$$DP_{\text{MinT}}(\Omega)(\boldsymbol{\theta}) = - \int_{\Omega} 2(d_{\Omega} \max(d_{\Omega} + d_{\text{min}}/2, 0)^2 + d_{\Omega}^2 \max(d_{\Omega} + d_{\text{min}}/2, 0)) d'_{\Omega}(\boldsymbol{\theta}) dx. \quad (4.4.17)$$

Note that an increase in perimeter entails a decrease in the value of  $P_{\text{MinT}}(\Omega)$ , but this behavior is tempered in our case by the volume constraint. Variants can be considered to address such an issue, and we refer to [234] for the details.

### Optimization of the shape of a 2-d cantilever beam with minimum thickness constraint

We first consider the classical two-dimensional cantilever benchmark example, as depicted on Figure 4.25: in a box  $D$  with size  $2 \times 1$ , shapes  $\Omega$  are clamped on their left-hand side, and a vertical load  $\mathbf{g} = (0, -10)$  is applied on the middle of their right-hand side.

The optimization problem (4.4.15) is solved with the parameter values  $g_{\text{max}} = 70.00$  and  $V_{\text{max}} = 0.80$ . The resulting optimized shapes are displayed on Figure 4.26 without including minimum thickness constraint (we minimize the volume  $\text{Vol}(\Omega)$  subject to the compliance constraint as in (4.4.13)), and with the minimum thickness for two different values of  $d_{\text{min}}$ . The corresponding convergence histories are shown on Figs. 4.27 to 4.29.

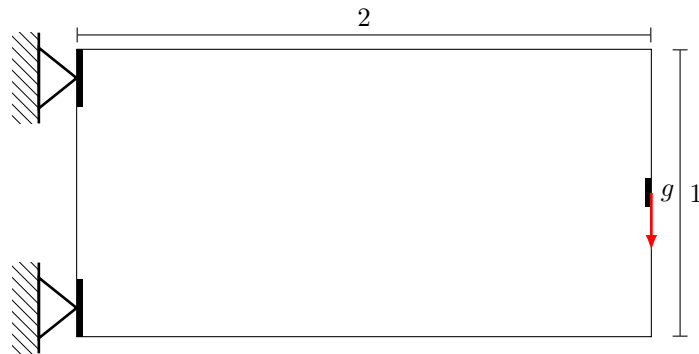


Figure 4.25: Setting of the cantilever test case of section 4.4.3.

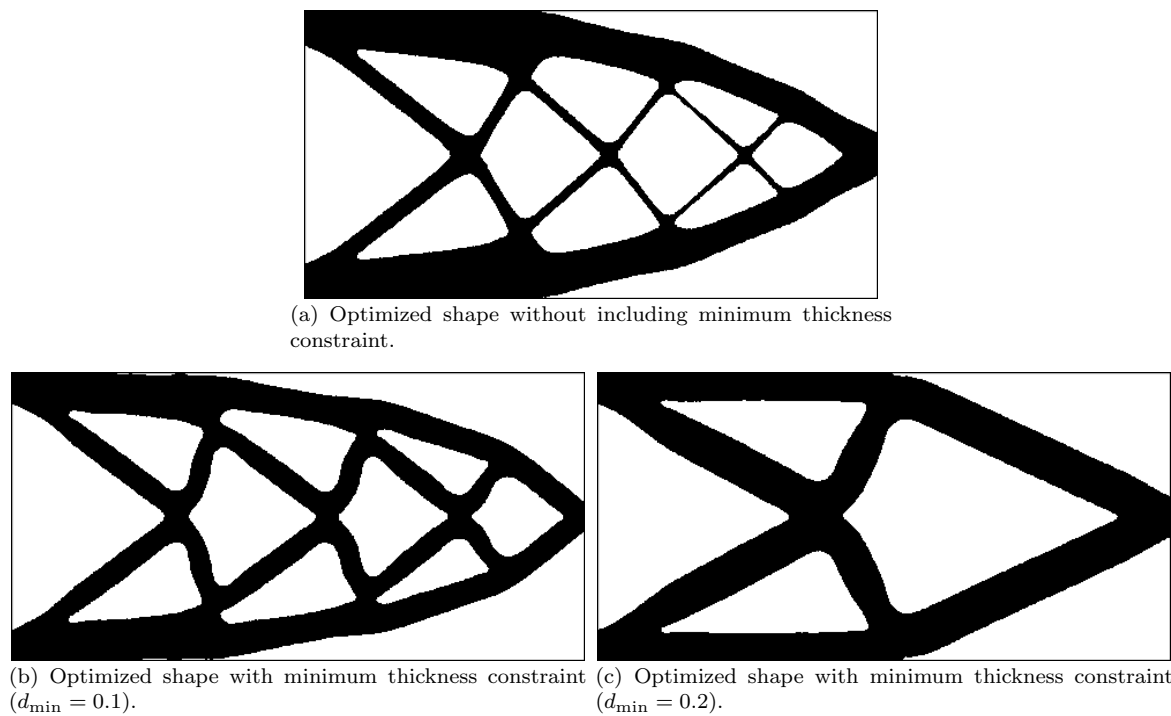


Figure 4.26: Optimization of the shape of the 2-d cantilever of [section 4.4.3](#) under minimum thickness constraint.

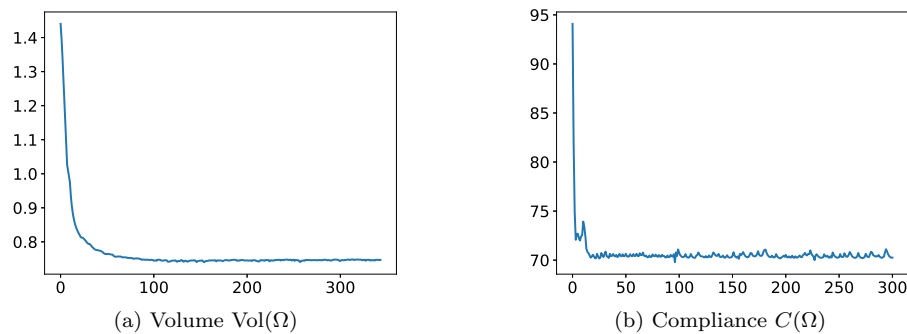


Figure 4.27: Optimization history for the 2-d cantilever optimization problem (4.4.13) of [section 4.4.3](#) without minimum thickness constraint. Final values: volume  $\text{Vol}(\Omega) = 0.75$ , and compliance  $C(\Omega) = 70.29$ .

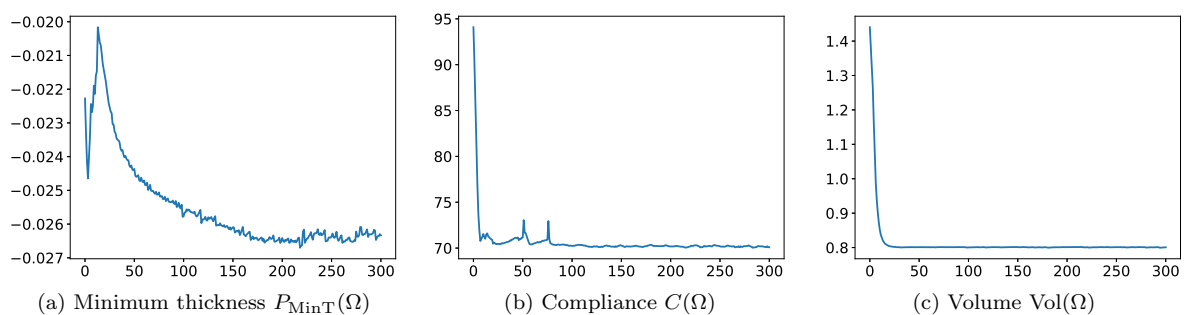


Figure 4.28: Optimization history for the 2-d cantilever optimization problem (4.4.13) of [section 4.4.3](#) with minimum thickness constraint. ( $d_{\min} = 0.1$ ) Final values : volume  $\text{Vol}(\Omega) = 0.80$ , compliance:  $C(\Omega) = 70.26$ .



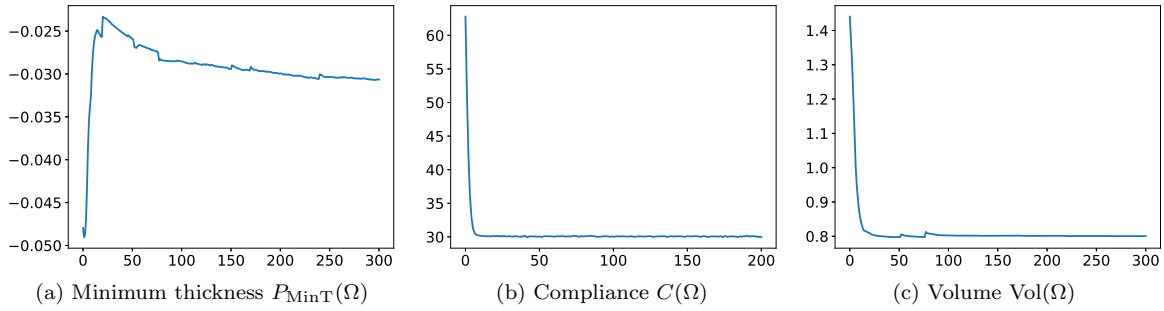


Figure 4.29: Optimization history for the 2-d cantilever optimization problem (4.4.13) of section 4.4.3 with minimum thickness constraint. ( $d_{\min} = 0.2$ ). Final values: volume  $\text{Vol}(\Omega) = 0.80$ , and compliance:  $C(\Omega) = 70.09$ .

### Shape optimization of a 2-d MBB Beam under a minimum thickness constraint

We now apply the same methodology on the MBB beam test-case of section 4.4.2. Optimized shapes are compared on Figure 4.30 without minimum thickness constraint (the result being that of Figure 4.22b), and for two values of  $d_{\min}$ . The corresponding convergence histories are shown on Figs. 4.23, 4.31 and 4.32. Finally, some intermediate shapes of the optimization process are reprinted on Figure 4.33.

One observes that for the last case of Figure 4.30c with  $d_{\min} = 0.2$ , the minimum thickness constraint is not satisfied everywhere but a substantial improvement is visible over the first design of Figure 4.30a. Notably, this approach is sufficiently flexible to guide the optimization path towards shapes with very different topologies.



(a) Optimized shape without minimum thickness constraint.



(b) Optimized shape with minimum thickness constraint ( $d_{\min} = 0.1$ ).



(c) Optimized shape with minimum thickness constraint ( $d_{\min} = 0.2$ ).

Figure 4.30: Minimum thickness optimization for a cantilever test case of section 4.4.3.

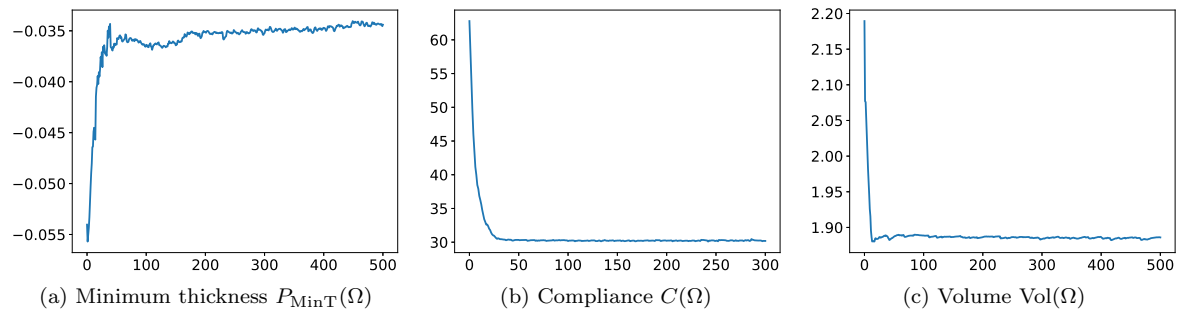


Figure 4.31: Convergence histories for the MBB-Beam test-case of section 4.4.3 with minimum thickness constraint. ( $d_{\min} = 0.1$ ). Final values: volume  $\text{Vol}(\Omega) = 1.89$ , and compliance  $C(\Omega) = 30.17$ .

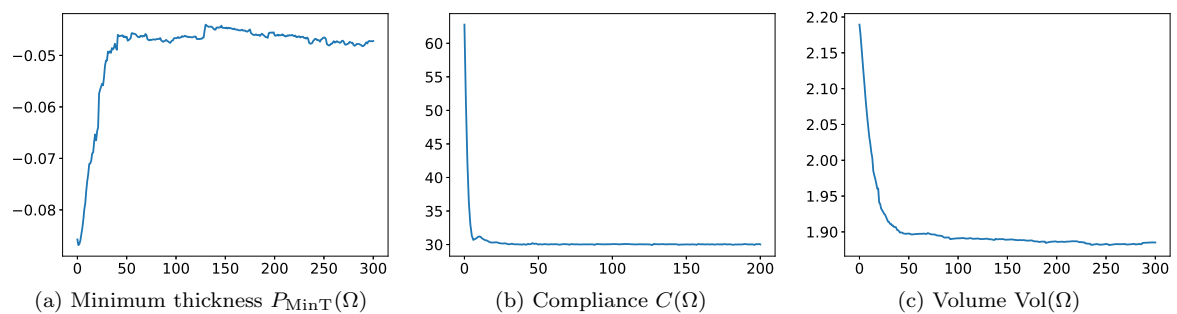


Figure 4.32: Convergence curves for the MBB-Beam test-case of section 4.4.3 with minimum thickness constraint ( $d_{\min} = 0.2$ ). Final values: volume  $\text{Vol}(\Omega) = 1.89$ , and compliance  $C(\Omega) = 30.18$ .

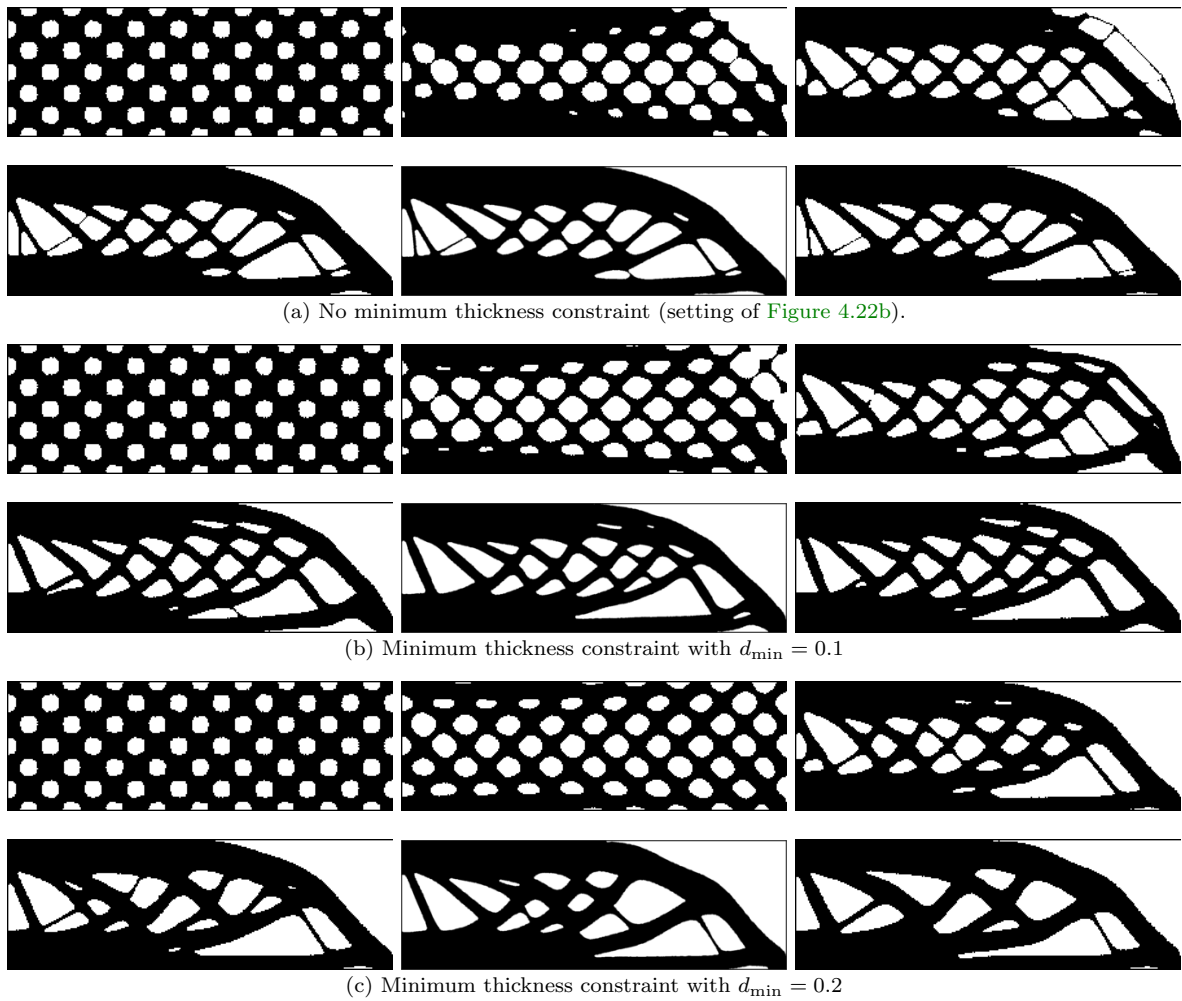


Figure 4.33: Optimization histories for the MBB beam test-case with minimum thickness constraints of section 4.4.3: Iterations 0, 10, 40, 100, 150 and 200.



## CHAPTER 5

# TOPOLOGY OPTIMIZATION OF 2-D HEAT EXCHANGERS

### Contents

---

<b>5.1 Design optimization of 2-d liquid-liquid heat exchangers with a non-mixing constraint</b>	<b>189</b>
5.1.1 Physical setting . . . . .	189
5.1.2 Formulation of the shape optimization problem . . . . .	191
5.1.3 Shape derivative of the non-mixing constraint $Q_{hot\leftrightarrow cold}(\Omega_f)$ . . . . .	193
5.1.4 Numerical results for two different configurations . . . . .	194
<b>5.2 Topology optimization of a 2-d air-oil heat exchanger</b>	<b>194</b>
5.2.1 Setting of the case study . . . . .	194
5.2.2 Physical modeling and formulation of optimization problems . . . . .	197
5.2.3 Numerical results . . . . .	199
5.2.4 An alternative model featuring a stagnation pressure boundary condition . . . . .	205
5.2.5 Conclusions . . . . .	208

---

Heat exchangers are a very topical issue in the topology optimization community with the publication of an increasing number of dedicated works [255, 258, 189, 271, 275, 205, 10, 104, 116, 225, 266, 303] This chapter reports on two independent case studies considering the shape optimization of two models of 2-d heat exchangers.

Section 5.1 considers a liquid-liquid heat exchanger featuring two fluid phases that must not interpenetrate. The cornerstone of the optimization problem is the treatment of the non-mixing condition: in our approach, it is imposed by using a minimum distance constraint between the two phases which lends itself to the variational method of chapter 4. This test case is inspired from a similar case study investigated by Papazoglou [255] with a very different (density based) method.

Section 5.2 then addresses the topology optimization of a different, air-oil heat exchanger. The aim is to determine the shape of the cross sections of oil pipes cooled down by an air flow. We rely on a different physical model, featuring a thermostatic boundary condition for the temperature on the optimized interface, which is very convenient to make the problem two dimensional. Since the optimization problem favors the apparition of very thin and elongated structures, a minimum thickness constraint is implemented to improve numerical convergence towards more manufacturable designs.

This part is an application of the material developed in the previous chapters: the formulas of chapter 2 for the shape derivatives of coupled thermal fluid problems, the null space algorithm of chapter 3 for optimization drives the resolution of the optimization problem, and we rely on the variational method of chapter 4 in order to enforce non penetration or minimum thickness constraints.

### 5.1 DESIGN OPTIMIZATION OF 2-D LIQUID-LIQUID HEAT EXCHANGERS WITH A NON-MIXING CONSTRAINT

This section investigates the shape optimization of a 2-d heat exchanger featuring two liquid phases that must not interpenetrate. The physical setting considered is described in section 5.1.1. The mathematical formulation of the optimization problem, including the modeling of the non-penetration constraint is outlined in section 5.1.2. Details about the computation of associated shape derivatives are provided in section 5.1.3. Finally, numerical results are presented in section 5.1.4.

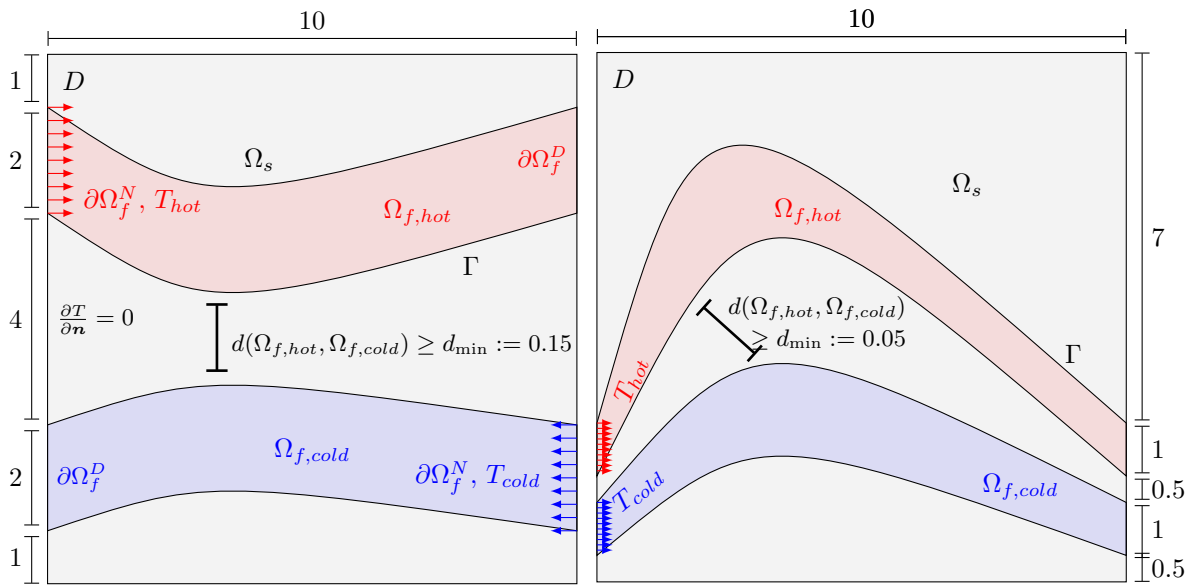
#### 5.1.1 Physical setting

We consider a ‘hold-all’ domain  $D = [0, 10] \times [0, 10] = \Omega_s \cup \Omega_f$  which is the disjoint union of a solid phase  $\Omega_s$  and a fluid phase  $\Omega_f$ . The fluid phase  $\Omega_f = \Omega_{f,hot} \cup \Omega_{f,cold}$  is itself constituted of two distinct channels  $\Omega_{f,hot}$  and  $\Omega_{f,cold}$  whose shapes are to be optimized. The phases  $\Omega_{f,hot}$  and  $\Omega_{f,cold}$  enter the domain  $D$  with respective inlet temperatures  $T_{hot} = 100$  and  $T_{cold} = 0$ . All other boundaries are adiabatic, corresponding to a zero Neumann condition  $\partial T / \partial \mathbf{n} = 0$  for the temperature field  $T$ . Two

possible configurations are considered about the location of the inlets and described in the schematic of Figure 5.1 below:

- Test case 1 (Figure 5.1a): the two liquid phases enter  $D$  in opposite directions. The inlet and outlet cross sections share a common size  $a = 2$ .
- Test case 2 (Figure 5.1b): the two liquid phases enter from the same side of  $D$  with inlet and outlet cross sections having a smaller common size  $a = 1$ .

Following the notation convention of chapter 2, the reunion of the two inlets, of the two outlets, and of total fluid-solid interface to be optimized are denoted respectively by  $\partial\Omega_f^D$ ,  $\partial\Omega_f^N$  and  $\Gamma$ . These test cases are very similar to those considered by Papazoglou [255].



(a) Test case 1: A hot fluid phase  $\Omega_{f,hot} \subset D$  is entering from the upper left side of  $D$  with a temperature  $T_{hot}$ , and a cold fluid phase  $\Omega_{f,cold}$  is entering in the reverse direction from the lower right inlet.  
(b) Test case 2: A hot fluid phase  $\Omega_{f,hot} \subset D$  is entering from the upper left side of  $D$  with a temperature  $T_{hot}$ , and a cold fluid phase  $\Omega_{f,cold}$  is entering in the same direction at the lower left inlet (boundary conditions not represented).

Figure 5.1: Settings of the two test cases considered in the optimal design of the heat exchangers of section 5.1 featuring the non-mixing condition  $d(\Omega_{f,hot}, \Omega_{f,cold}) \geq d_{min}$ .

The physics involved in this problem are those described in chapter 2, section 3.6.3: the velocity of the fluid velocity and pressure  $(\mathbf{v}, p)$  are characterized by the Navier-Stokes equations in the total fluid domain

$$\Omega_f := \Omega_{f,cold} \cup \Omega_{f,hot}$$

while the temperature field  $T$  in the whole domain  $D$  is determined by the equations of conduction-convection in both solid and liquid phases  $\Omega_s$  and  $\Omega_f$ :

$$\left\{ \begin{array}{ll} -\operatorname{div}(\sigma_f(\mathbf{v}, p)) + \rho \nabla \mathbf{v} \cdot \mathbf{v} = 0 & \text{in } \Omega_f \\ \operatorname{div}(\mathbf{v}) = 0 & \text{in } \Omega_f \\ \mathbf{v} = \mathbf{v}_0 & \text{on } \partial\Omega_f^D \\ \sigma_f(\mathbf{v}, p) \mathbf{n} = 0 & \text{on } \partial\Omega_f^N \\ \mathbf{v} = 0 & \text{on } \Gamma, \end{array} \right. \quad (5.1.1)$$

$$\left\{ \begin{array}{ll} -\operatorname{div}(k_f \nabla T_f) + \rho c_p \mathbf{v} \cdot \nabla T_f = 0 & \text{in } \Omega_f \\ -\operatorname{div}(k_s \nabla T_s) = 0 & \text{in } \Omega_s \\ T = 100 & \text{on } \partial\Omega_f^D \cap \partial\Omega_{f,hot} \\ T = 0 & \text{on } \partial\Omega_f^D \cap \partial\Omega_{f,cold} \\ -k_f \frac{\partial T_f}{\partial \mathbf{n}} = 0 & \text{on } \partial\Omega_f^N \cap \partial\Omega_f \\ -k_s \frac{\partial T_s}{\partial \mathbf{n}} = 0 & \text{on } \partial\Omega_f^N \cap \partial\Omega_s \\ T_f = T_s & \text{on } \Gamma \\ -k_f \frac{\partial T_f}{\partial \mathbf{n}} = -k_s \frac{\partial T_s}{\partial \mathbf{n}} & \text{on } \Gamma. \end{array} \right. \quad (5.1.2)$$

For each of the considered test cases, both fluid phases  $\Omega_{f,cold}$  and  $\Omega_{f,hot}$  are assumed to have the same density  $\rho = 1$  and thermal conductivity  $k_f = 10$ . They enter the domain  $D$  from the inlet boundaries  $\partial\Omega_f^D \subset \partial D$  with a parabolic velocity profile with maximum norm  $\|\mathbf{v}_0\|_\infty = 1$ , and they exit at outlet boundaries  $\partial\Omega_f^N \subset \partial D$  with the vanishing normal stress condition  $\sigma_f(\mathbf{v}, p) \cdot \mathbf{n} = 0$ , where we recall that the fluid stress tensor is defined by the Newton law

$$\sigma_f(\mathbf{v}, p) := 2\nu e(\mathbf{v}) - pI.$$

The viscosity  $\nu$  is the same for both fluids; it is computed by the formula  $\nu := \rho a \|\mathbf{v}_0\|_\infty / \operatorname{Re}$  where  $\operatorname{Re} = 60$  is the Reynolds number and  $a$  is the size of the inlet ( $a = 2$  or  $a = 1$  for the respective test cases 1 and 2). The capacity coefficient of the fluids is calculated by  $c_p := k_f \operatorname{Pe} / (\nu \operatorname{Re})$  with the Péclet number given by  $\operatorname{Pe} := 500$ . The solid phase  $\Omega_s$  is assumed to have a larger thermal conductivity coefficient  $k_s = 110$  than the fluid.

### 5.1.2 Formulation of the shape optimization problem

The aim of the optimal design problem is to find the shape of both fluid phases  $\Omega_f := \Omega_{f,cold} \cup \Omega_{f,hot}$  which maximizes the heat exchanged between the “hot” and “cold” phases under a maximal pressure drop constraint and a non-mixing condition. This exchanged heat is mathematically appraised by the opposite of an objective functional  $J(\Omega_f, \mathbf{v}(\Omega_f), T(\Omega_f))$  which is to be minimized:

$$J(\Omega_f, \mathbf{v}(\Omega_f), T(\Omega_f)) := - \left( \int_{\Omega_{f,cold}} \rho c_p \mathbf{v} \cdot \nabla T dx - \int_{\Omega_{f,hot}} \rho c_p \mathbf{v} \cdot \nabla T dx \right).$$

This quantity can indeed be interpreted as the heat transferred because an integration by part (see also [chapter 2, section 2.5.7](#)) implies that

$$-J(\Omega_f, \mathbf{v}(\Omega_f), T(\Omega_f)) = \int_{\partial\Omega_{f,cold}} \rho c_p T \mathbf{v} \cdot \mathbf{n} dy - \int_{\partial\Omega_{f,hot}} \rho c_p T \mathbf{v} \cdot \mathbf{n} dy.$$

The above expression turns out to be exactly the heat exiting the cold phase minus the one leaving the hot phase (to be maximized), up to additional constant terms depending on inlet boundary values.

Following [chapter 2, section 2.5.7](#), the pressure drop constraint reads

$$\operatorname{DP}(\Omega_f) := \int_{\partial\Omega_f^D} p dy - \int_{\partial\Omega_f^N} p dy \leq \operatorname{DP}_0$$

where  $\operatorname{DP}_0$  is a given threshold value. In our implementation, this value is set according to the initial domain  $\Omega_f^0$ :

$$\operatorname{DP}_0 := \begin{cases} 2\operatorname{DP}(\Omega_f^0) & \text{in test case 1,} \\ 5\operatorname{DP}(\Omega_f^0) & \text{in test case 2.} \end{cases}$$

The initial domain  $\Omega_f^0$  features two straight pipes which are displayed for each case 1 and 2 on [Figure 5.2](#) below.

Let us now discuss the modeling of the non-mixing condition between both fluid phases. In his thesis, Papazoglou [255] proposed a multi-material model suited to the use of density based topology

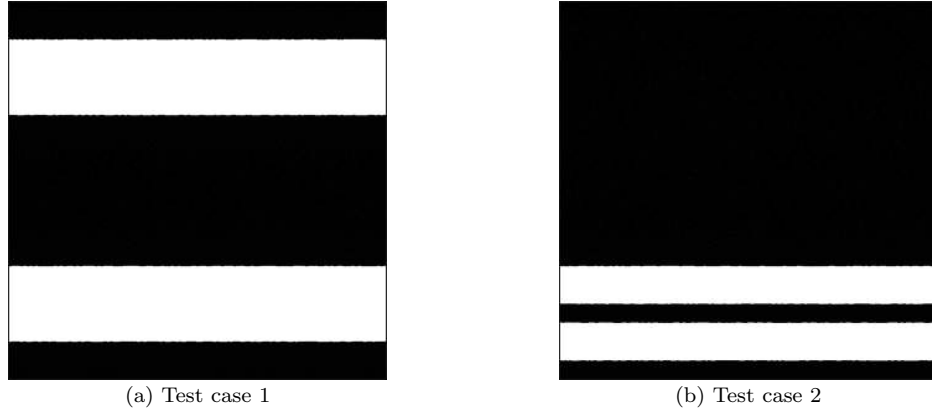


Figure 5.2: Initial distribution of the fluid domain  $\Omega_f^0$  (in white).

optimization. In the context of the method of Hadamard for shape optimization, the “hot” and “cold” domains at stake are explicitly defined at every iteration which makes it possible to formulate the non-mixing constraint in a geometric fashion.

The key idea is to consider the *two* signed distance functions  $d_{\Omega_{f,cold}}$  and  $d_{\Omega_{f,hot}}$  associated with the two fluid phases  $\Omega_{f,cold}$  and  $\Omega_{f,hot}$ . Following the introduction, this non-penetration condition is enforced by requiring the hot phase to remain at a minimum distance  $d_{\min}$  from the cold phase, which can mathematically be formulated as:

$$\forall x \in \Omega_{f,cold}, d_{\Omega_{f,hot}}(x) \geq d_{\min}, \quad (5.1.3)$$

or equivalently

$$\forall x \in \partial\Omega_{f,cold}, d_{\Omega_{f,hot}}(x) \geq d_{\min}. \quad (5.1.4)$$

Note that the roles of  $\Omega_{f,hot}$  and  $\Omega_{f,cold}$  in (5.1.3) and (5.1.4) are interchangeable. For our applications, we set  $d_{\min} = 0.15$  for the test case 1, and  $d_{\min} = 0.05$  for the test case 2. Following the strategy of Allaire, Jouve and Michailidis in [30, 234], we consider averaged penalty functionals in order to approximate the pointwise constraint (5.1.4). Such can be obtained by reformulating (5.1.4) with the infinity norm:

$$\left\| \frac{1}{d_{\Omega_{f,hot}}} \right\|_{L^\infty(\partial\Omega_{f,cold})} \leq \frac{1}{d_{\min}}. \quad (5.1.5)$$

Classically, we consider an approximation of the infinity norm by a  $L^p$  norm, which yields an averaged penalty functional  $P_{cold \rightarrow hot}(\Omega_f)$ :

$$P_{cold \rightarrow hot}(\Omega_f) := \left( \int_{\partial\Omega_{f,cold}} \left( \frac{1}{d_{\Omega_{f,hot}}} \right)^p dy \right)^{\frac{1}{p}} \simeq \left\| \frac{1}{d_{\Omega_{f,hot}}} \right\|_{L^\infty(\partial\Omega_{f,cold})}. \quad (5.1.6)$$

For our application, the parameter  $p$  involved is set to  $p = 4$ .

**Remark 5.1.** The formulation (5.1.4) involving  $\partial\Omega_{f,hot}$  is preferred to that (5.1.3) set on the whole domain  $\Omega_{f,hot}$ , because we expect that averaging on a smaller set in (5.1.6) yields in some sense a more accurate approximation of the infinity norm.

**Remark 5.2.** This constraint on the distance between the two connected components  $\Omega_{f,cold}$  and  $\Omega_{f,hot}$  is substantially different to the requirement of a minimum thickness for the solid phase  $\Omega_s$ . The use of the signed distance functions to  $\Omega_{f,cold}$  (or to  $\Omega_{f,hot}$ ) makes it possible to formulate an *actual* constraint, which is much more involved for a minimum thickness for  $\Omega_s$ , in comparison with the approach proposed by [30] and in section 4.3.4.

In order to balance the effect of the constraint (5.1.6) over both fluid phases: we introduce the symmetrized version of (5.1.6):

$$P_{hot \rightarrow cold}(\Omega_f) := \left( \int_{\partial\Omega_{f,hot}} \left( \frac{1}{d_{\Omega_{f,cold}}} \right)^p dy \right)^{\frac{1}{p}}. \quad (5.1.7)$$



It is then convenient to formulate the non-penetration constraint in terms of the harmonic mean of  $P_{hot\rightarrow cold}$  and  $P_{cold\rightarrow hot}$ , which yields an averaged measure of the distance  $d(\Omega_{f,cold}, \Omega_{f,hot})$  between both phases; the non penetration constraint (5.1.5) is approximated by the following one:

$$Q_{hot\leftrightarrow cold}(\Omega_f) := \frac{2}{P_{hot\rightarrow cold}(\Omega_f) + P_{cold\rightarrow hot}(\Omega_f)} \geq d_{\min}. \quad (5.1.8)$$

All in all, the considered optimization problem reads:

$$\begin{aligned} \min_{\Omega_f \subset D} \quad & J(\Omega_f) = - \left( \int_{\Omega_{f,cold}} \rho c_p \mathbf{v} \cdot \nabla T dx - \int_{\Omega_{f,hot}} \rho c_p \mathbf{v} \cdot \nabla T dx \right) \\ \text{s.t.} \quad & \begin{cases} DP(\Omega_f) = \int_{\partial\Omega_f^P} p ds - \int_{\partial\Omega_f^N} p ds \leq DP_0 \\ Q_{hot\leftrightarrow cold}(\Omega_f) \geq d_{\min}. \end{cases} \end{aligned} \quad (5.1.9)$$

This optimization problem is solved with the null space gradient flow described in chapter 3 and the level set based mesh evolution algorithm of chapter 1, section 1.4.2. The shape derivatives of  $J(\Omega_f)$  and  $DP(\Omega_f)$  are computed with the formulas given in volume form in chapter 2, proposition 2.3. The calculation of the shape derivative of the penalty functional  $Q_{hot\leftrightarrow cold}(\Omega_f)$  is described in the next section.

### 5.1.3 Shape derivative of the non-mixing constraint $Q_{hot\leftrightarrow cold}(\Omega_f)$

The expression of the shape derivative of  $Q_{hot\leftrightarrow cold}(\Omega_f)$  is easily obtained from those of  $P_{cold\rightarrow hot}$  and  $P_{hot\rightarrow cold}$ . Since the latter have similar expressions, we content ourselves with an outline of the calculation of the shape derivative of  $P_{cold\rightarrow hot}$ . A straightforward computation yields

$$\begin{aligned} DP_{cold\rightarrow hot}(\Omega_f)(\boldsymbol{\theta}) = \frac{1}{p} \left( \int_{\partial\Omega_{f,cold}} \frac{1}{|d_{\Omega_{f,hot}}|^p} dy \right)^{\frac{1}{p}-1} \left[ \int_{\partial\Omega_{f,cold}} \left( \frac{\partial}{\partial \mathbf{n}} + \kappa \right) \left( \frac{1}{|d_{\Omega_{f,hot}}|^p} \right) (\boldsymbol{\theta} \cdot \mathbf{n}) dy \right. \\ \left. - \int_{\partial\Omega_{f,cold}} \frac{p}{|d_{\Omega_{f,hot}}|^{p+1}} d'_{\Omega_{f,hot}}(\boldsymbol{\theta}) dy, \right] \end{aligned} \quad (5.1.10)$$

where we recall that  $d'_{\Omega_{f,hot}}(\boldsymbol{\theta})$  is the Eulerian derivative of the signed distance function  $d_{\Omega_{f,hot}}$  to  $\Omega_{f,hot}$  (see chapter 1, section 1.3.2). The first term of the above right-hand side has the usual structure of a boundary integral which can be easily numerically evaluated as a linear form in terms of  $\boldsymbol{\theta} \cdot \mathbf{n}$ . The last term is computed thanks to the variational method proposed in chapter 4. To be quite specific in this setting involving two distinct phases, we solve

$$\begin{aligned} \text{Find } u_{hot} \in V_\omega \text{ such that } \forall v \in V_\omega, \\ \int_{\partial\Omega_{f,hot}} u_{hot} v dy + \int_{D \setminus \bar{\Sigma}_{hot}} \omega (\nabla u_{hot} \cdot \nabla d_{\Omega_{f,hot}}) (\nabla v \cdot \nabla d_{\Omega_{f,hot}}) dx = \int_{\partial\Omega_{f,cold}} \frac{p}{|d_{\Omega_{f,hot}}|^{p+1}} v dy, \end{aligned} \quad (5.1.11)$$

where  $\bar{\Sigma}_{hot}$  denotes the skeleton set of  $\Omega_{f,hot}$  and the weight  $\omega$  is given by

$$\omega = 1/(1 + 100|d_{\Omega_{f,hot}} \Delta d_{\Omega_{f,hot}}|^{3.5}).$$

The reader is referred to chapter 4, section 4.2 regarding the definition of the space  $V_\omega$ , which is discretized with  $\mathbb{P}_1$  finite elements. Finally, the last integral term in (5.1.10) is readily obtained from the formula

$$\int_{\partial\Omega_{f,cold}} \frac{p}{|d_{\Omega_{f,hot}}|^{p+1}} d'_{\Omega_{f,hot}}(\boldsymbol{\theta}) dy = - \int_{\partial\Omega_{f,hot}} u_{hot} \boldsymbol{\theta} \cdot \mathbf{n} dy,$$

where the right-hand side is a boundary integral involving  $\boldsymbol{\theta} \cdot \mathbf{n}$  which can be easily discretized in our finite element setting.

### 5.1.4 Numerical results for two different configurations

Numerical results are depicted on Figs. 5.3 and 5.5 for the respective test cases 1 and 2. For each test case, we plot the initial and final designs, the temperature  $T$  and the velocity field  $\mathbf{v}$ , as well as some intermediate shapes obtained with our algorithm and the convergence histories for the objective and constraint functionals. We retrieve serpentine shapes as observed in [255].

This test case exhibits two striking features. At first, our approximation (5.1.6) of the  $L^\infty$  norm constraint by the averaged penalty functional  $P_{cold \rightarrow hot}(\Omega_f)$  with a  $L^p$  norm ( $p = 4$  in our case) works (surprisingly) very well: the ‘ideal’ pointwise constraint seems to be imposed at all intermediate iterations. On a different note, it is remarkable that our overall method is able to keep finding better shapes even after the saturation of the distant constraint, which happens very early in the optimization process. The reader will notice that some noise affect the convergence curves: we attribute it to inaccuracies related to the rather large Péclet and Reynolds number considered.

Finally, let us mention that we did not need a very fine mesh of the region of  $\Omega_s$  in between the cold and the hot domains  $\Omega_{f,cold}$  and  $\Omega_{f,hot}$  to handle the distance constraint (5.1.8). We plot on Figs. 5.4 and 5.6 the final meshes of the optimized shape for the respective test cases 1 and 2. For the second test case, it is visible that a resolution of about ten mesh elements (which means a skeleton located at approximately *five* mesh elements only) in between the two pipes allows for a satisfying approximation of the shape derivative of the distance constraint from the resolution of the variational problem (5.1.11).

## 5.2 TOPOLOGY OPTIMIZATION OF A 2-D AIR-OIL HEAT EXCHANGER

We now report on a study issued from a collaboration with Safran Aero Boosters. The problem at hand is concerned with the optimization of the cross section of an air-oil heat exchanger. Although the original problem is three-dimensional, a few assumptions allow to reduce it to 2 dimensions. The approach followed in this part is different to that of the previous section 5.1: it is not necessary to prescribe a non-mixing constraint for the air and oil phases because we consider a different physical model in which there is no “solid” phase  $\Omega_s$ ; the air phase occupies a domain  $\Omega_f$  to be optimized, and the oil phase is described by an isothermal boundary condition on  $\partial\Omega_f$ .

The general setting of the considered problem is described in section 5.2.1. Let us emphasize that the physical model considered in this setting is essentially the same as the fluid thermal model (5.1.1) and (5.1.2) of the previous section up to a change of boundary conditions on the optimized air-oil interface. The optimization problem is then mathematically formulated in section 5.2.2, where we define objective and constraint functions. A minimum thickness constraint is considered for the oil phase in order to improve the convergence towards manufacturable designs. Section 5.2.3 presents a variety of numerical results for different initializations and settings of the constraints. We mention the consideration of a different fluid model featuring a boundary condition on the stagnation pressure in section 5.2.4, which could be of direct interest for realistic industrial applications. The interesting point is the use of the rotational formulation of the Navier-Stokes equations [106]. However, our numerical results with this adaptation are not conclusive. We finally conclude in section 5.2.5 with a few general remarks about the extension of this approach to more realistic 3-d systems characterized by large Reynolds numbers.

### 5.2.1 Setting of the case study

The objective of the study is the optimization of the shape and topology of a 2-d air-oil heat exchanger. The optimization domain  $D$  contains an air phase which occupies a subdomain  $\Omega_f \subset D$  to be optimized. The complementary  $D \setminus \Omega_f$  is filled up with oil. The setting of the problem is depicted on the schematic of Figure 5.7. After a 3-d extrusion in the  $z$  direction, the system is interpreted as a heat exchanger featuring air flowing in the  $x$  direction and cooling down oil channels flowing in the transverse  $z$  direction.

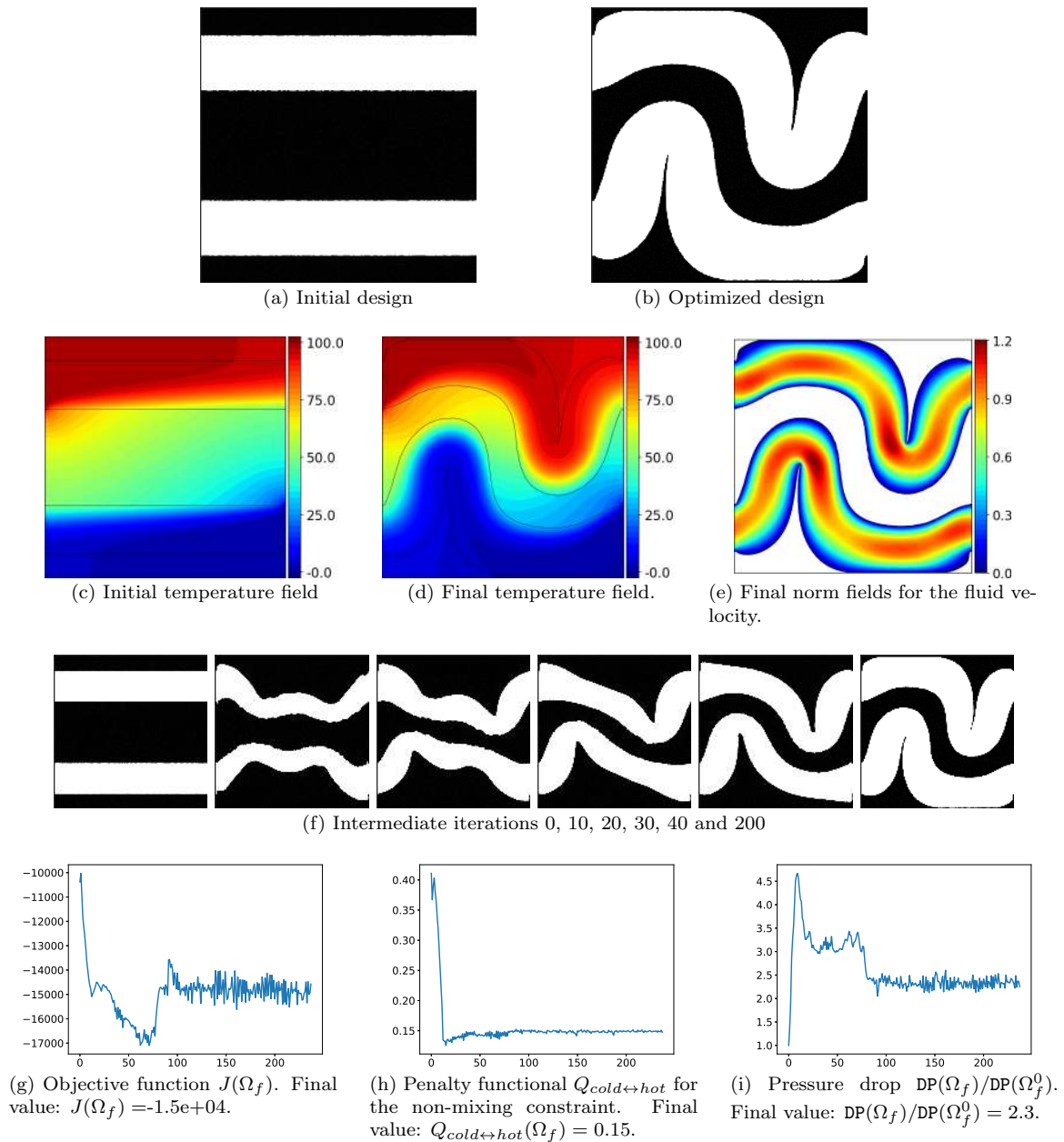


Figure 5.3: Optimization results for the test case 1 (cold and hot inlets entering  $D$  from opposite sides) of section 5.1.

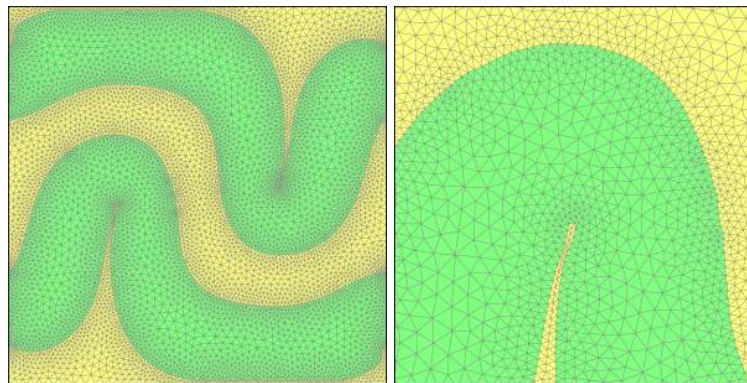


Figure 5.4: Zoom on a mesh of an intermediate optimization iteration.

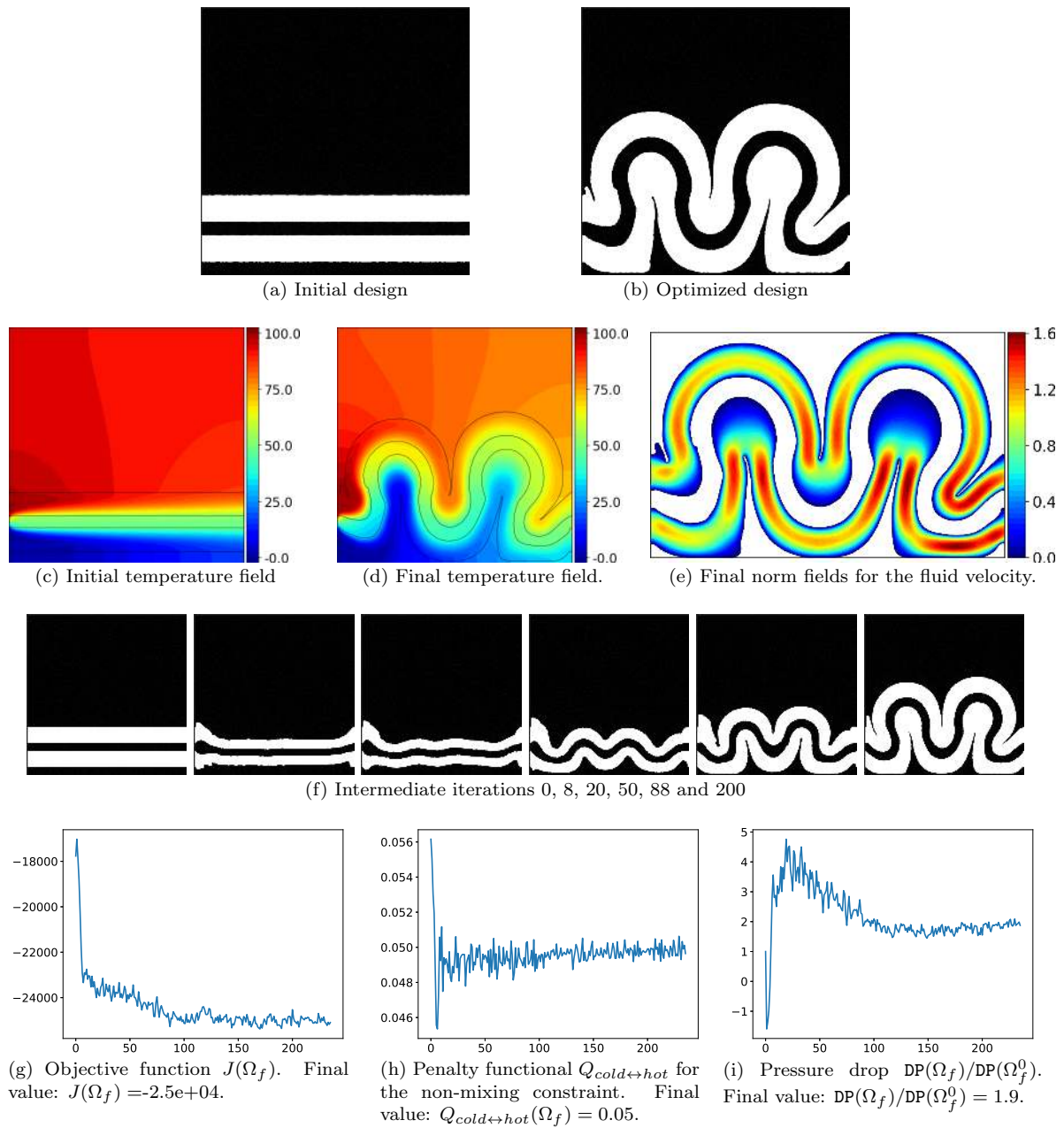


Figure 5.5: Optimization results for the test case 2 (cold and hot inlets entering  $D$  from the same side) of section 5.1.

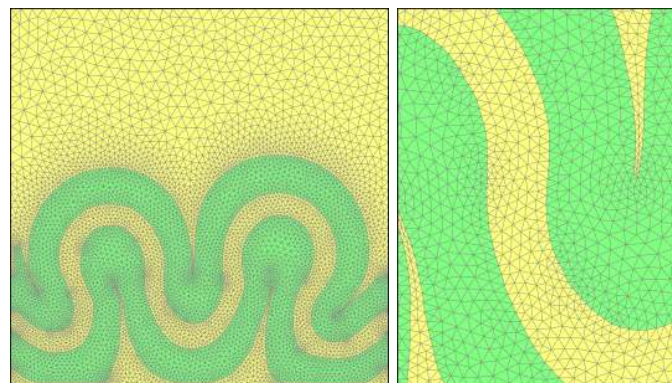


Figure 5.6: Zoom on a mesh of an intermediate optimization iteration.

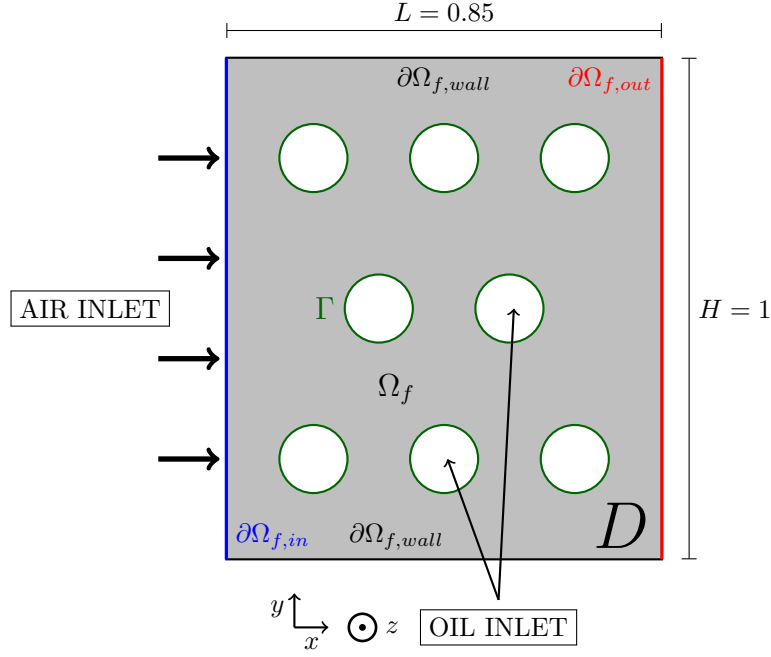


Figure 5.7: Setting of the air-oil heat exchanger of section 5.2.

The oil filling the phase  $D \setminus \Omega_f$  has a much higher thermal conductivity than air so that its temperature  $T$  can be considered to be uniform:  $T = T_{oil}$  in  $D \setminus \Omega_f$ . Assuming in addition the interface  $\Gamma$  to be made of a solid phase that is assumed to be a sufficiently thin and conductive material, it is sufficient to describe the system by the single physics of the air phase with a thermostatic temperature condition  $T = T_{oil}$  on the interface  $\Gamma$ .

The objective of the study is to demonstrate the feasibility of topology optimization with the method of Hadamard for the maximization of the heat exchanged between the two phases, under a maximum constraint on the outlet pressure drop.

### 5.2.2 Physical modeling and formulation of optimization problems

The air phase is characterized by fluid velocity and pressure  $(\mathbf{v}, p)$  solutions to the incompressible steady state Navier-Stokes equations:

$$\left\{ \begin{array}{ll} -\operatorname{div}(\sigma_f(\mathbf{v}, p)) + \rho \nabla \mathbf{v} \mathbf{v} = 0 & \text{in } \Omega_f \\ \operatorname{div}(\mathbf{v}) = 0 & \text{in } \Omega_f \\ \mathbf{v} = \mathbf{v}_0 & \text{on } \partial\Omega_{f,in} \\ \sigma_f(\mathbf{v}, p) \mathbf{n} = 0 & \text{on } \partial\Omega_{f,out} \\ \mathbf{v} = 0 & \text{on } \partial\Omega_{f,wall} \\ \mathbf{v} = 0 & \text{on } \Gamma, \end{array} \right. \quad (5.2.1)$$

where the fluid stress tensor is given by  $\sigma_f(\mathbf{v}, p) = 2\nu e(\mathbf{v}) - pI$  as in the previous section 5.1. The fluid velocity  $\mathbf{v}$  then determines a temperature field  $T$  through the following convection diffusion equation:

$$\left\{ \begin{array}{ll} -\operatorname{div}(k_f \nabla T) + \rho c_p \mathbf{v} \cdot \nabla T = 0 & \text{in } \Omega_f \\ T = T_{in} & \text{on } \partial\Omega_{f,in} \\ -k_f \frac{\partial T}{\partial \mathbf{n}} = 0 & \text{on } \partial\Omega_f \setminus \partial\Omega_{f,in} \\ T = T_{oil} & \text{on } \Gamma. \end{array} \right. \quad (5.2.2)$$

This model differs from the one considered in the previous section (equations (5.1.1) and (5.1.2)) due to the fact that a thermostatic boundary condition  $T = T_{oil}$  is imposed on the interface  $\Gamma$  rather than a diffusion model within the complementary subdomain  $D \setminus \Omega_f$ . Note that we assume, as in our previous studies, a zero normal temperature flux on the wall and outlet boundaries  $\partial\Omega_{f,wall}$  and  $\partial\Omega_{f,out}$ . The

	Input velocity $\ \mathbf{v}_0\ $	Pressure loss threshold $DP_0$	Re	Pe
Configuration 1	10	1300	40	475
Configuration 2	25	1030	100	1200
Configuration 3	40	475	160	1900

Table 5.1: Configurations considered for the input velocity and pressure constraint values.

fluid is entering the left boundary  $\partial\Omega_{f,in}$  with a parabolic velocity profile  $\mathbf{v}_0$  with maximum norm value  $\|\mathbf{v}_0\|_\infty$ . The “cold” air flow  $\Omega_f$  is entering with a temperature  $T_{in} := 310$ . The temperature of the “hot” oil phase  $D \setminus \Omega_f$  flowing in the opposite direction is constant and it is given by  $T_{oil} = 400$ . The capacity and density coefficients of the fluid are set to  $\rho = 1$ ,  $c_p = 1$ . Several values are considered for the viscosity  $\nu$ , the conductivity coefficient  $k_f$  and for the intensity  $\|\mathbf{v}_0\|_\infty$ , which are discussed below.

**Definition of objective and constraint functions** The goal is to maximize the exchanged heat between the air phase and the oil phase while imposing an upper bound on the static pressure drop:

$$\begin{aligned} \min_{\Omega_f \subset D} \quad & J(\Omega_f) := - \int_{\Omega_f} \rho c_p \mathbf{v} \cdot \nabla T dx \\ \text{s.t.} \quad & DP(\Omega_f) := \int_{\partial\Omega_{f,in}} p ds - \int_{\partial\Omega_{f,out}} p ds \leq DP_0. \end{aligned} \quad (5.2.3)$$

Upon integration by parts, the objective function  $J(\Omega_f)$  rewrites as the opposite of the heat transported from the inlet boundary  $\partial\Omega_{f,in}$  to the outlet boundary  $\partial\Omega_{f,out}$ :

$$J(\Omega_f) = - \left( \int_{\partial\Omega_{f,in}} \rho c_p T \mathbf{v} \cdot \mathbf{n} ds + \int_{\partial\Omega_{f,out}} \rho c_p T \mathbf{v} \cdot \mathbf{n} ds \right).$$

Although mathematically equivalent, we prefer to use the volume form in (5.2.3) which seems to us more accurate numerically (following the previous section and [chapter 2, section 2.5.7](#)).

The optimization problem (5.2.3) is solved for several values for the input velocity  $\|\mathbf{v}_0\|_\infty$  and the pressure loss threshold  $DP_0$ . The details of the considered three situations are provided in [Table 5.1](#) below.

**Reynolds, Péclet numbers and numerical values of the physical parameters** Since we do not rely on a turbulent model in the Navier-Stokes equations (5.2.1) for the determination of the fluid velocity and pressure  $(\mathbf{v}, p)$ , our study is restricted to moderate values of Reynolds and Péclet numbers. The viscosity and conductivity coefficients  $\nu$  and  $k_f$  are computed so as to fix the Reynolds and Péclet number

$$\text{Re} := \frac{\rho \|\mathbf{v}_0\|_\infty H}{\nu}, \quad \text{Pe} := \frac{\rho c_p \|\mathbf{v}_0\|_\infty H}{k_f}$$

to the values provided for each configuration in [Table 5.1](#).

**Minimum thickness constraint for the oil phase cross section** The resolution of the optimization problem (5.2.3) tends to produce very thin and elongated shapes for the oil cross section. These are favorable for the optimization problem (5.2.3), however they are numerically unstable. Indeed, small components of the design thinner than the prescribed mesh size tend to disappear due to numerical diffusion during the optimization process. We therefore consider an alternative version of (5.2.3) which tends to impose a minimum thickness  $d_{\min}$  of the oil channels. This constraint could also be of interest in the perspective of the manufacturing process. We follow the strategy proposed in [chapter 4, section 4.4.3](#): the performance of the system  $J(\Omega_f)$  is prescribed to be at least as good as a threshold value  $J_0$ , and

we minimize instead an energy  $E(\Omega_f)$  which favors areas of  $D \setminus \Omega_f$  thicker than  $d_{\min}$ :

$$\begin{aligned} \min_{\Omega_f \subset D} \quad & E(\Omega_f) := - \int_{D \setminus \Omega_f} d_{\Omega_f}^2 \max(-d_{\Omega_f} + d_{\min}/2, 0)^2 dx \\ \text{s.t.} \quad & \begin{cases} \text{DP}(\Omega_f) \leq \text{DP}_0 \\ J(\Omega_f) \leq J_0. \end{cases} \end{aligned} \quad (5.2.4)$$

The reader is referred to [chapter 4](#) for the computation of the shape derivative of  $E(\Omega_f)$ . The threshold value  $J_0$  measuring the performance of the heat exchange is determined empirically from the values obtained with the resolution of the original problem [\(5.2.3\)](#) which does not feature the minimum thickness constraint. For our application, we took  $J_0 = 10000$ .

**Shape derivatives for the heat exchanger problem** The physics considered for this test case is essentially the same as that of [\(5.1.1\)](#) and [\(5.1.2\)](#); the only difference is the Dirichlet boundary condition  $T = T_{oil}$  on the interface  $\Gamma$ . It is easily verified that the shape derivative formulas provided in [propositions 2.3](#) and [2.4](#) are still valid up to a change of boundary condition for the thermal adjoint equation: the adjoint problem [\(2.4.9\)](#) is solved with the boundary condition  $S = 0$  on  $\Gamma$ . Then, it can be shown that the shape derivative of an arbitrary functional  $J(\Omega_f, \mathbf{v}(\Omega_f), p(\Omega_f), T(\Omega_f))$  in volume and surface forms read respectively:

$$\begin{aligned} \frac{d}{d\boldsymbol{\theta}} \left[ J(\Gamma_{\boldsymbol{\theta}}, \mathbf{v}(\Gamma_{\boldsymbol{\theta}}), p(\Gamma_{\boldsymbol{\theta}}), T(\Gamma_{\boldsymbol{\theta}})) \right] (\boldsymbol{\theta}) &= \frac{\partial \mathfrak{J}}{\partial \boldsymbol{\theta}} (\boldsymbol{\theta}) + \int_{\Omega_f} [ -(\sigma_f(\mathbf{v}, p) : \nabla \mathbf{w} + \rho \mathbf{w} \cdot \nabla \mathbf{v} \cdot \mathbf{v}) \text{div}(\boldsymbol{\theta}) ] dx \\ &+ \int_{\Omega_f} [ \sigma_f(\mathbf{v}, p) : (\nabla \mathbf{w} \nabla \boldsymbol{\theta}) + \sigma_f(\mathbf{w}, q) : (\nabla \mathbf{v} \nabla \boldsymbol{\theta}) + \rho \mathbf{w} \cdot (\nabla \mathbf{v} \nabla \boldsymbol{\theta}) \cdot \mathbf{v} ] dx \\ &- \int_{\Omega_f} \text{div}(\boldsymbol{\theta}) (k_f \nabla T \cdot \nabla S + \rho c_p (\mathbf{v} \cdot \nabla T) S) dx \\ &+ \int_{\Omega_f} [ k_f (\nabla \boldsymbol{\theta} + \nabla \boldsymbol{\theta}^T) \nabla T \cdot \nabla S + \rho c_p \mathbf{v} \cdot (\nabla \boldsymbol{\theta}^T \nabla T) S ] dx, \end{aligned} \quad (5.2.5)$$

$$\frac{d}{d\boldsymbol{\theta}} \left[ J(\Gamma_{\boldsymbol{\theta}}, \mathbf{v}(\Gamma_{\boldsymbol{\theta}}), p(\Gamma_{\boldsymbol{\theta}}), T(\Gamma_{\boldsymbol{\theta}})) \right] (\boldsymbol{\theta}) = \frac{\overline{\partial \mathfrak{J}}}{\partial \boldsymbol{\theta}} (\boldsymbol{\theta}) + \int_{\Gamma} \left( 2\nu e(\mathbf{w}) : e(\mathbf{v}) + k_f \frac{\partial T}{\partial \mathbf{n}} \frac{\partial S}{\partial \mathbf{n}} \right) (\boldsymbol{\theta} \cdot \mathbf{n}) ds. \quad (5.2.6)$$

where the notation  $\mathfrak{J}$  refers to the modified functional of [\(2.4.7\)](#) introduced in [chapter 2](#).

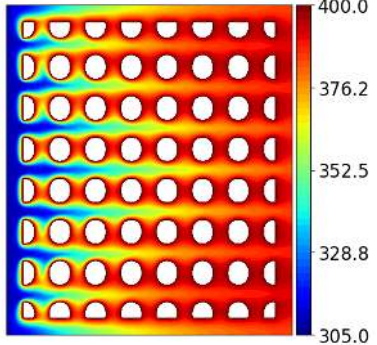
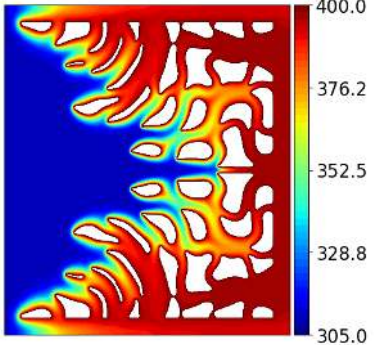
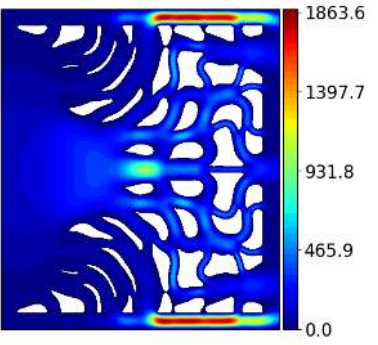
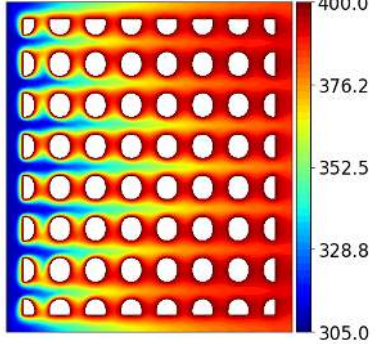
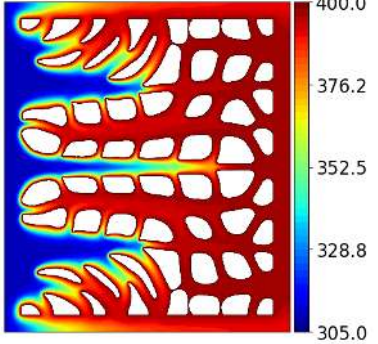
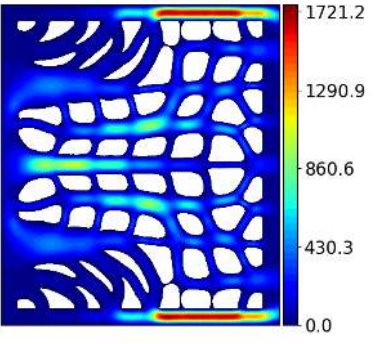
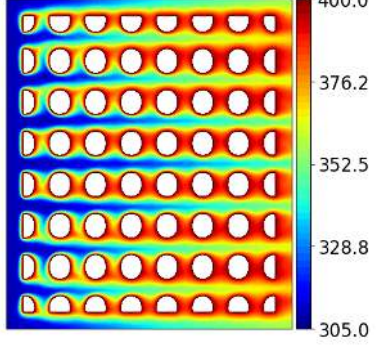
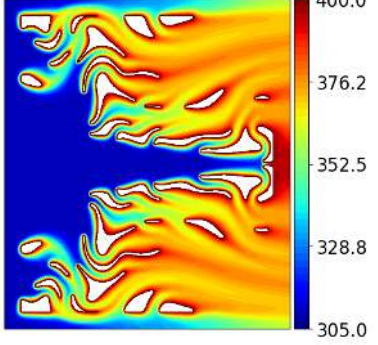
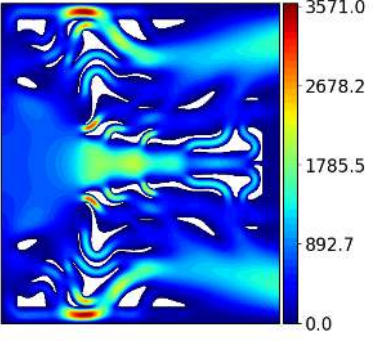
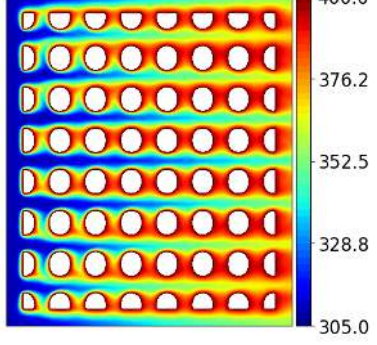
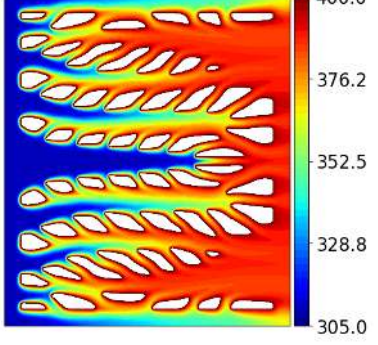
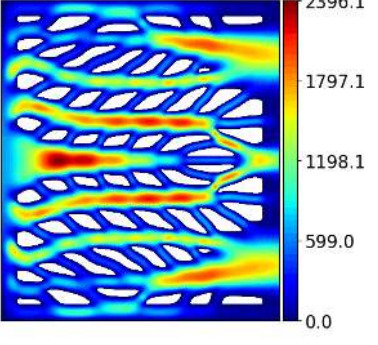
### 5.2.3 Numerical results

Both optimization problems [\(5.2.3\)](#) and [\(5.2.4\)](#) are solved for the three configurations considered in [Table 5.1](#). Since there is no unique solution, we propose a set of 24 results obtained by varying the shape chosen for the initialization. These are reported in the [Table 5.2](#) below where we display the final temperature and kinetic energy fields. We note that our optimization algorithm is able to (i) create recirculating fluid regions favorable to heat exchange and (ii) make the transverse oil channels assume an aerodynamic shape limiting the output pressure loss. We also observe that the effect of the minimum thickness constraint in [\(5.2.4\)](#) is most significant on configurations with maximal input velocity  $\|\mathbf{v}_0\|_{\infty} = 40$  which feature the most elongated structures. A few intermediate shapes are reported in [Figs. 5.8a](#) to [5.8c](#).

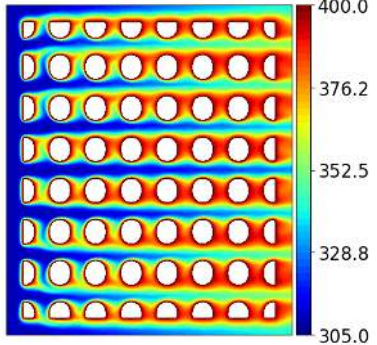
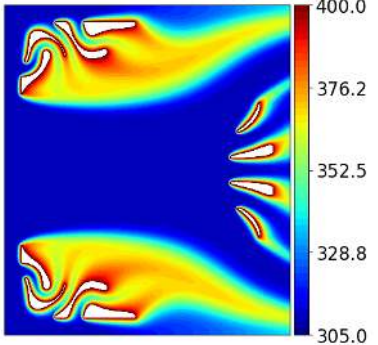
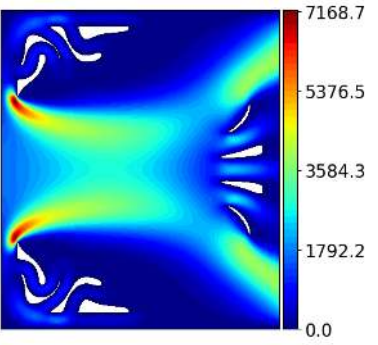
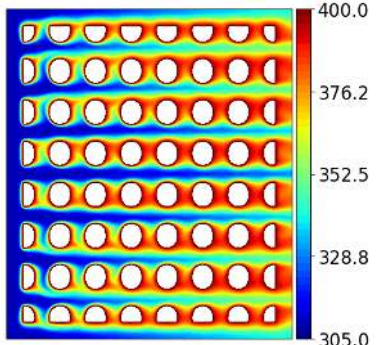
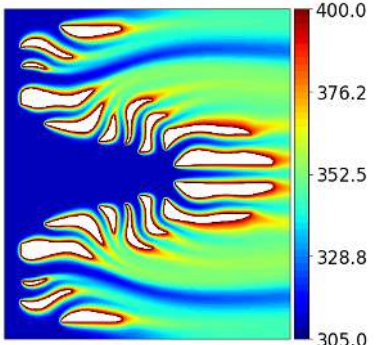
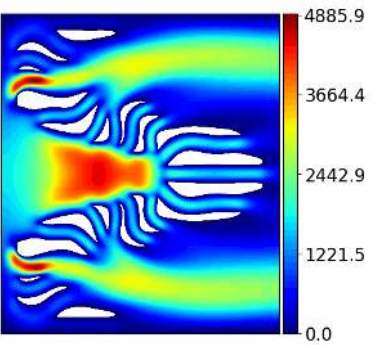
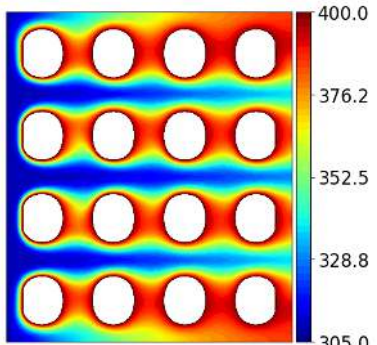
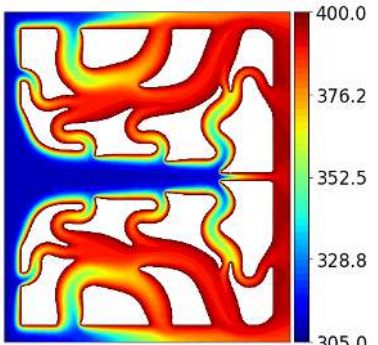
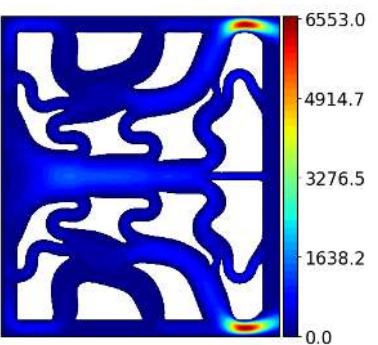
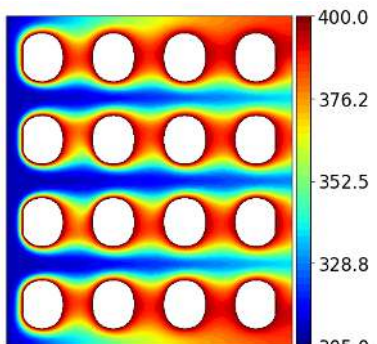
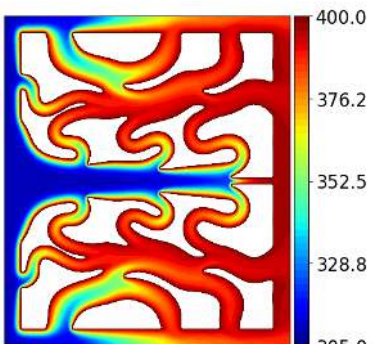
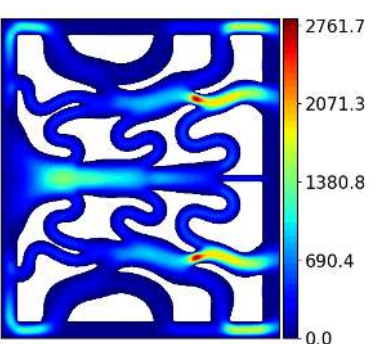
We relied on a rather fine mesh resolution (the minimum edge size was of the order of  $h_{\min} = 0.003$ ), which is illustrated on [Figure 5.9](#) where the mesh of the final shape is displayed for one of the test cases.

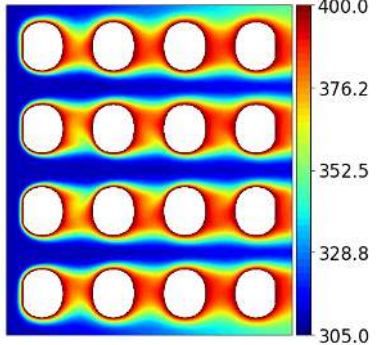
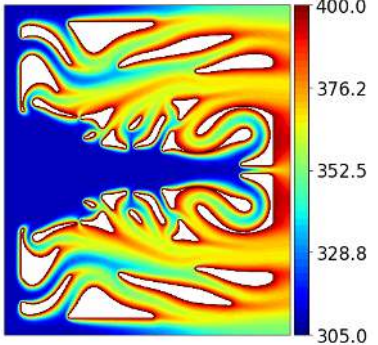
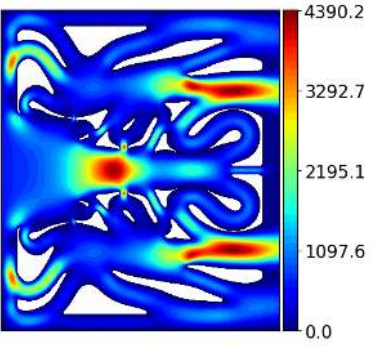
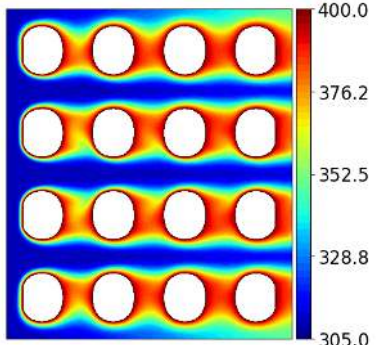
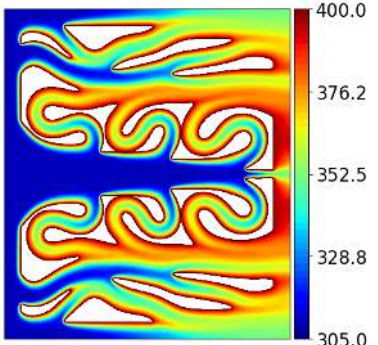
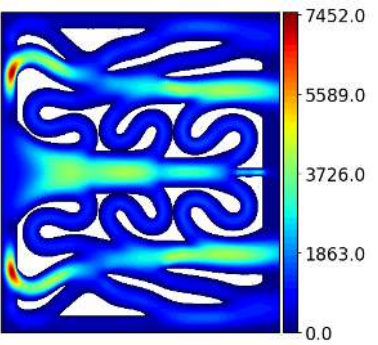
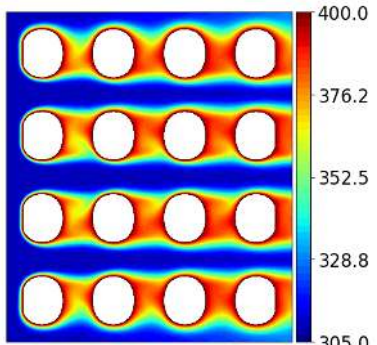
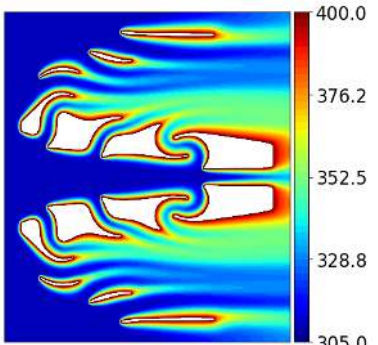
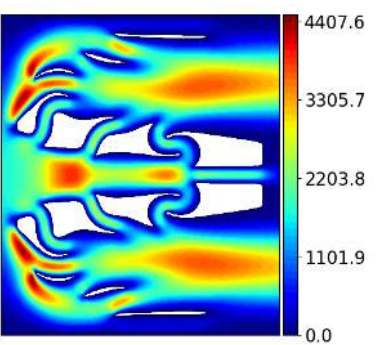
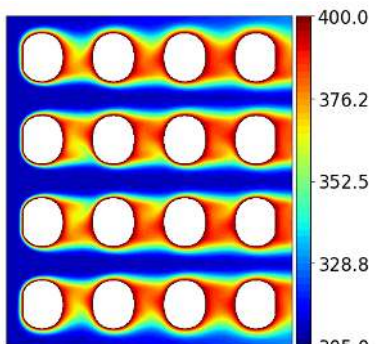
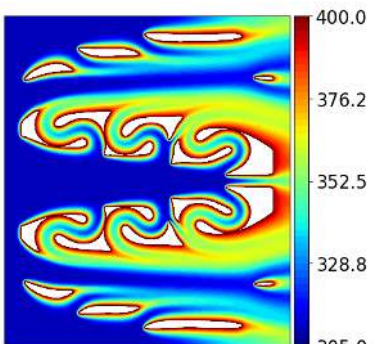
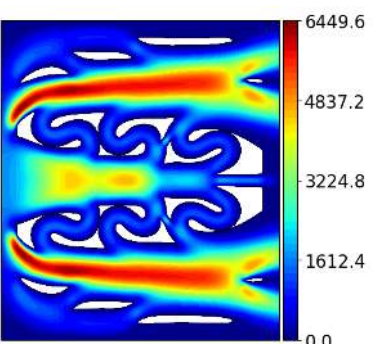
These results point out that a rather large variety of candidate shapes may be satisfactory solutions for the heat exchanger problem [\(5.2.3\)](#). Although the test case remains academic because we do not take turbulence into account, these rather unconventional designs further highlight longer term perspectives for topology optimization in the context of coupled fluid thermal industrial applications.

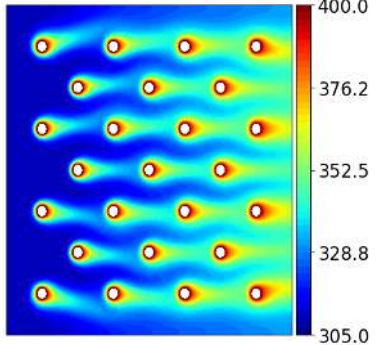
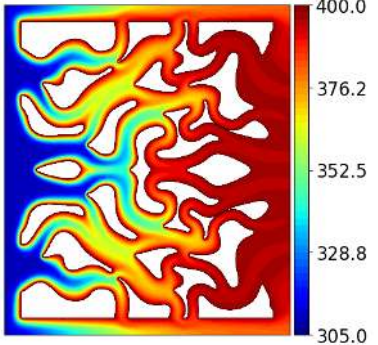
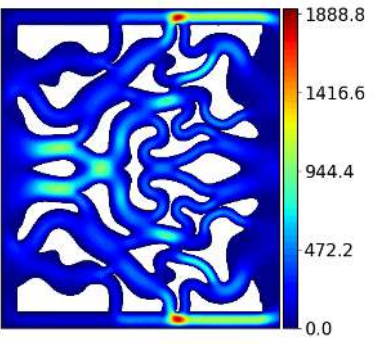
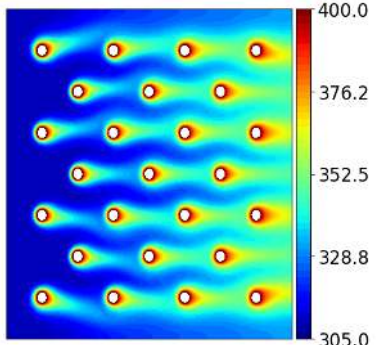
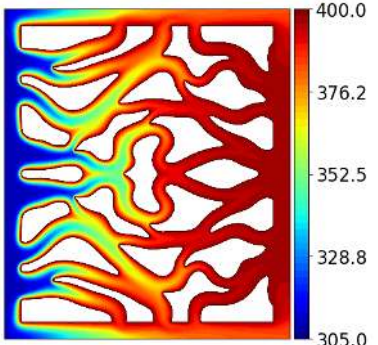
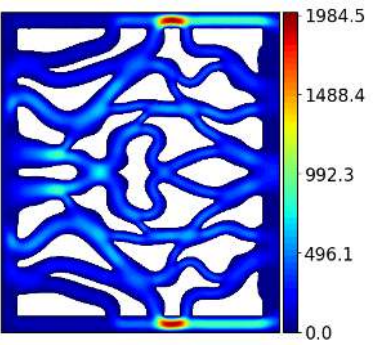
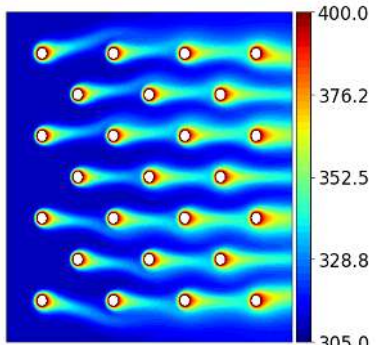
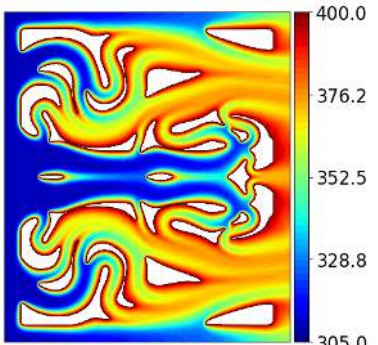
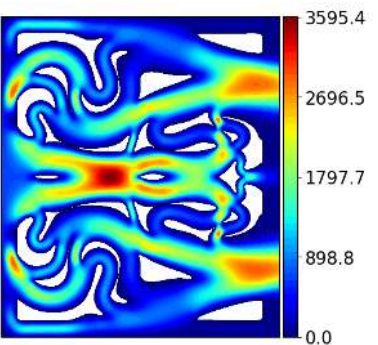
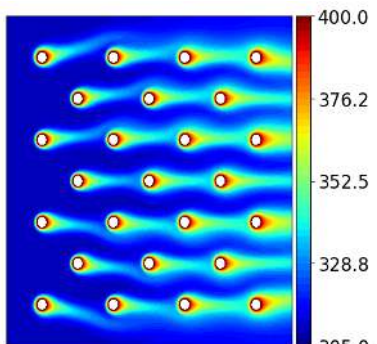
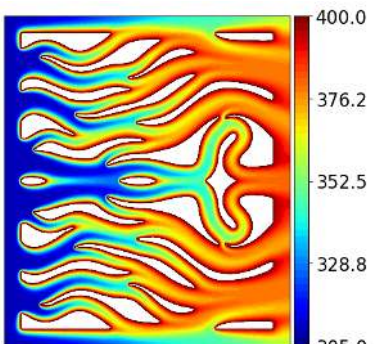
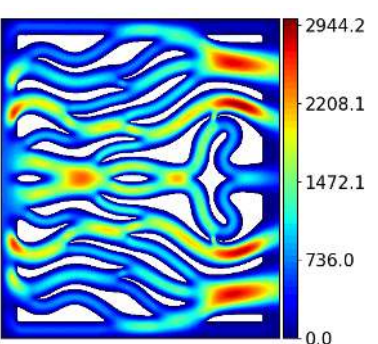
Test case	$T$ (Initial design)	$T$ (Optimized design)	$\ \mathbf{v}\ ^2$ (Optimized design)
-----------	----------------------	------------------------	---------------------------------------

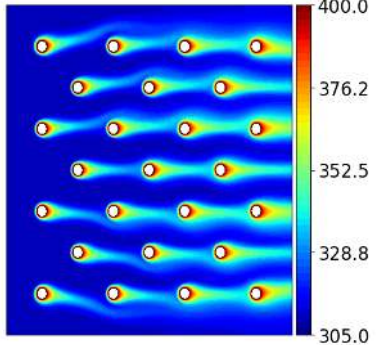
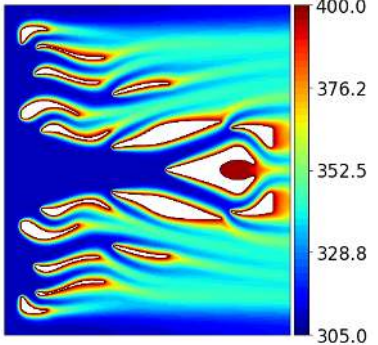
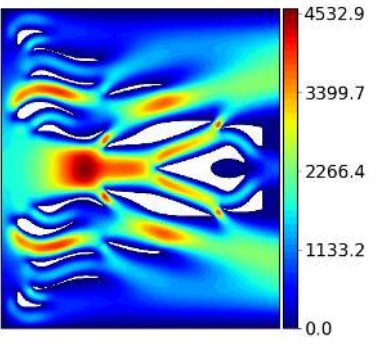
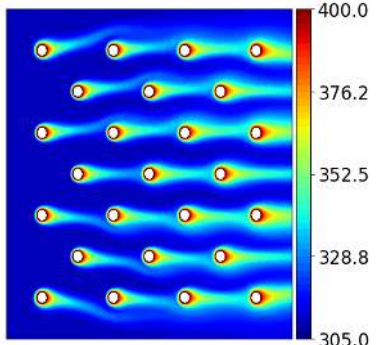
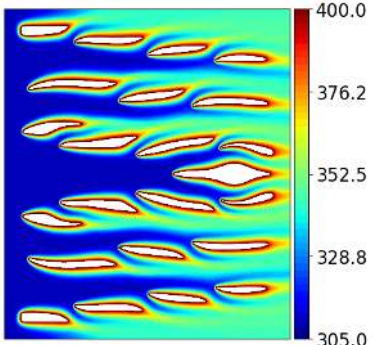
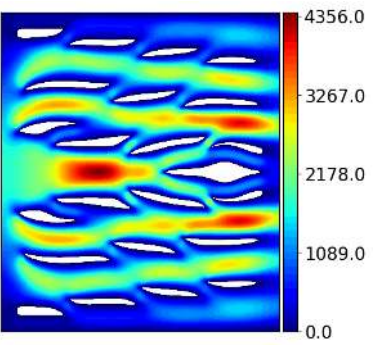
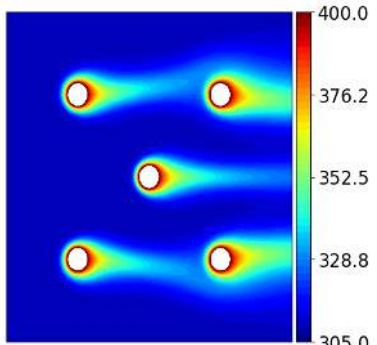
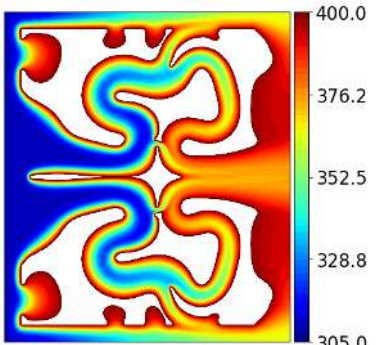
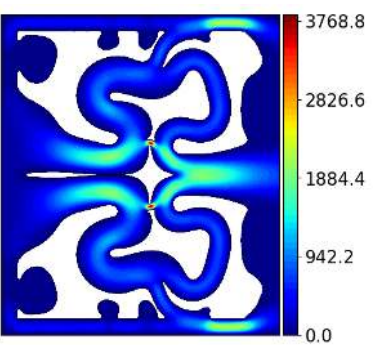
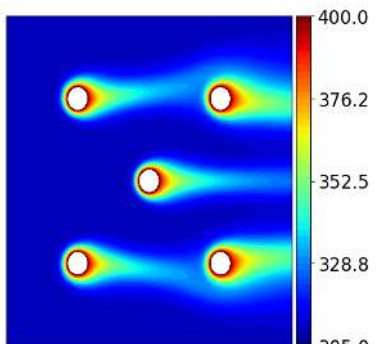
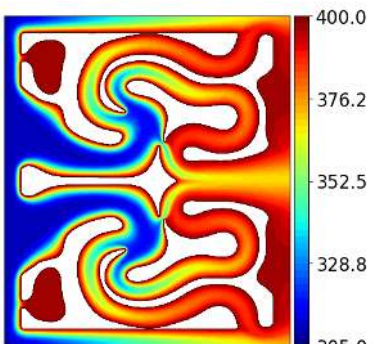
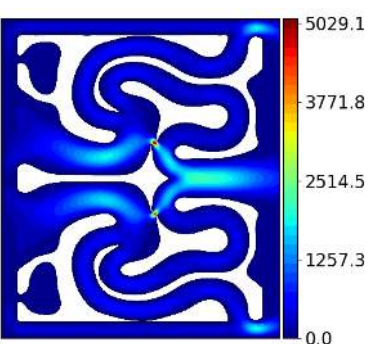
Test case	$T$ (Initial design)	$T$ (Optimized design)	$\ \mathbf{v}\ ^2$ (Optimized design)
<b>Test case 1</b> $\ \mathbf{v}_0\ _\infty=10$ $DP_0=1300$ $J_{final}=4350$ $DP_{final}=1113$ <i>Without</i> min. thickness constraint			
<b>Test case 2</b> $\ \mathbf{v}_0\ _\infty=10$ $DP_0=1300$ $J_{final}=4346$ $DP_{final}=1217$ <i>With</i> min. thickness constraint			
<b>Test case 3</b> $\ \mathbf{v}_0\ _\infty=25$ $DP_0=1030$ $J_{final}=8089$ $DP_{final}=983$ <i>Without</i> min. thickness constraint			
<b>Test case 4</b> $\ \mathbf{v}_0\ _\infty=25$ $DP_0=1030$ $J_{final}=9742$ $DP_{final}=1030$ <i>With</i> min. thickness constraint			



Test case	$T$ (Initial design)	$T$ (Optimized design)	$\ v\ ^2$ (Optimized design)
<p><b>Test case 5</b>  <math>\ v_0\ _\infty=40</math>  <math>DP_0=475</math>  <math>J_{final}=3472</math>  <math>DP_{final}=392</math>  <i>Without</i> min. thickness constraint</p>			
<p><b>Test case 6</b>  <math>\ v_0\ _\infty=40</math>  <math>DP_0=475</math>  <math>J_{final}=7285</math>  <math>DP_{final}=520</math>  <i>With</i> min. thickness constraint</p>			
<p><b>Test case 7</b>  <math>\ v_0\ _\infty=10</math>  <math>DP_0=1300</math>  <math>J_{final}=4086</math>  <math>DP_{final}=1308</math>  <i>Without</i> min. thickness constraint</p>			
<p><b>Test case 8</b>  <math>\ v_0\ _\infty=10</math>  <math>DP_0=1300</math>  <math>J_{final}=4168</math>  <math>DP_{final}=1188</math>  <i>With</i> min. thickness constraint</p>			

Test case	$T$ (Initial design)	$T$ (Optimized design)	$\ v\ ^2$ (Optimized design)
<b>Test case 9</b> $\ v_0\ _\infty=25$ $DP_0=1030$ $J_{final}=7667$ $DP_{final}=968$ <i>Without</i> min. thickness constraint			
<b>Test case 10</b> $\ v_0\ _\infty=25$ $DP_0=1030$ $J_{final}=7508$ $DP_{final}=1112$ <i>With</i> min. thickness constraint			
<b>Test case 11</b> $\ v_0\ _\infty=40$ $DP_0=475$ $J_{final}=5731$ $DP_{final}=479$ <i>Without</i> min. thickness constraint			
<b>Test case 12</b> $\ v_0\ _\infty=40$ $DP_0=475$ $J_{final}=6847$ $DP_{final}=524$ <i>With</i> min. thickness constraint			

Test case	$T$ (Initial design)	$T$ (Optimized design)	$\ \mathbf{v}\ ^2$ (Optimized design)
<p><b>Test case 13</b>  <math>\ \mathbf{v}_0\ _\infty=10</math>  <math>DP_0=1300</math>  <math>J_{final}=4208</math>  <math>DP_{final}=1140</math>  <i>Without</i> min. thickness constraint</p>			
<p><b>Test case 14</b>  <math>\ \mathbf{v}_0\ _\infty=10</math>  <math>DP_0=1300</math>  <math>J_{final}=4252</math>  <math>DP_{final}=1157</math>  <i>With</i> min. thickness constraint</p>			
<p><b>Test case 15</b>  <math>\ \mathbf{v}_0\ _\infty=25</math>  <math>DP_0=1030</math>  <math>J_{final}=7785</math>  <math>DP_{final}=1022</math>  <i>Without</i> min. thickness constraint</p>			
<p><b>Test case 16</b>  <math>\ \mathbf{v}_0\ _\infty=25</math>  <math>DP_0=1030</math>  <math>J_{final}=8711</math>  <math>DP_{final}=1106</math>  <i>With</i> min. thickness constraint</p>			

Test case	$T$ (Initial design)	$T$ (Optimized design)	$\ \mathbf{v}\ ^2$ (Optimized design)
<b>Test case 17</b> $\ \mathbf{v}_0\ _\infty=40$ $DP_0=475$ $J_{final}=6236$ $DP_{final}=470$ <i>Without</i> min. thickness constraint			
<b>Test case 18</b> $\ \mathbf{v}_0\ _\infty=40$ $DP_0=475$ $J_{final}=7822$ $DP_{final}=498$ <i>With</i> min. thickness constraint			
<b>Test case 19</b> $\ \mathbf{v}_0\ _\infty=10$ $DP_0=1300$ $J_{final}=3361$ $DP_{final}=1149$ <i>Without</i> min. thickness constraint			
<b>Test case 20</b> $\ \mathbf{v}_0\ _\infty=10$ $DP_0=1300$ $J_{final}=3582$ $DP_{final}=1064$ <i>With</i> min. thickness constraint			

Test case	$T$ (Initial design)	$T$ (Optimized design)	$\ \mathbf{v}\ ^2$ (Optimized design)
<b>Test case 21</b> $\ \mathbf{v}_0\ _\infty=25$ $DP_0=1030$ $J_{final}=2972$ $DP_{final}=986$ Without min. thickness constraint			
<b>Test case 22</b> $\ \mathbf{v}_0\ _\infty=25$ $DP_0=1030$ $J_{final}=5330$ $DP_{final}=1589$ With min. thickness constraint			
<b>Test case 23</b> $\ \mathbf{v}_0\ _\infty=40$ $DP_0=475$ $J_{final}=2847$ $DP_{final}=476$ Without min. thickness constraint			
<b>Test case 24</b> $\ \mathbf{v}_0\ _\infty=40$ $DP_0=475$ $J_{final}=4925$ $DP_{final}=1051$ With min. thickness constraint			

Table 5.2: Topology optimization results for the air-oil heat exchanger case study.

### 5.2.4 An alternative model featuring a stagnation pressure boundary condition

Industrial specifications of heat exchangers often impose prescribed input pressure values  $P_{in}$  and  $P_{out}$  at the inlet and the outlet of the system rather than the upper bound condition  $DP(\Omega_f) \leq DP_0$  on the pressure loss. Furthermore,  $P_{in}$  and  $P_{out}$  correspond to the “stagnation” pressure  $p + \rho v^2/2$  ( $v = \|\mathbf{v}\|^2$ )

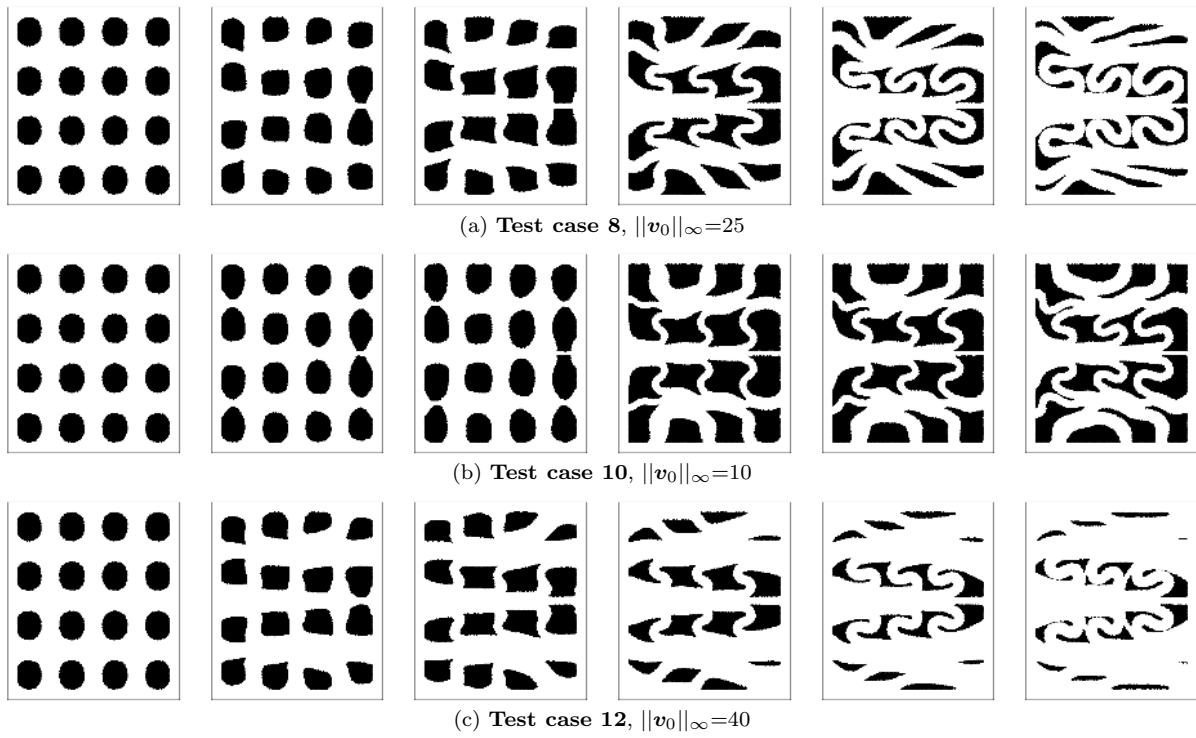


Figure 5.8: Iterations 0, 10, 20, 100, 200, 400 for several test case configurations.



Figure 5.9: Mesh of the final shape of a test case featuring  $\|v_0\| = 10$ .

in our case) rather than the static pressure  $p$  (the “stagnation pressure” can be directly experimentally measured with pitot tubes). These considerations suggest to consider a variant of the optimization problem (5.2.3) where the pressure loss constraint would be enforced as a boundary condition involving the pressures  $P_{in}$  and  $P_{out}$  in the Navier-Stokes equations (5.1.1).

It turns out that it is indeed possible to impose the value of the stagnation pressure upon a judicious rewriting of the Navier-Stokes equations in their classical “rotational form”.

### Rewriting of the physical model and shape derivatives

The Navier-Stokes system in rotational form classically reads

$$\left\{ \begin{array}{ll} -\nu\Delta\mathbf{v} + \nabla\left(p + \frac{1}{2}\|\mathbf{v}\|^2\right) + \rho(\nabla \times \mathbf{v}) \times \mathbf{v} = 0 & \text{in } \Omega_f \\ \operatorname{div}(\mathbf{v}) = 0 & \text{in } \Omega_f \\ p + \frac{1}{2}\rho\|\mathbf{v}\|^2 = P_{in} & \text{on } \partial\Omega_{f,in} \\ \mathbf{v} \times \mathbf{n} = 0 & \text{on } \partial\Omega_{f,in} \\ p + \frac{1}{2}\rho\|\mathbf{v}\|^2 = P_{out} & \text{on } \partial\Omega_{f,out} \\ \mathbf{v} \times \mathbf{n} = 0 & \text{on } \partial\Omega_{f,out} \\ \mathbf{v} = 0 & \text{on } \partial\Omega_{f,wall} \\ \mathbf{v} = 0 & \text{on } \Gamma. \end{array} \right. \quad (5.2.7)$$

The associated variational formulation then reads (see the works of Conca et. al. [105, 106]):

Find  $(\mathbf{v}, q) \in V_{\mathbf{v},q}$  such that  $\forall(\mathbf{w}, r) \in V_{\mathbf{v},q}$ ,

$$\begin{aligned} \int_{\Omega_f} [\nu(\nabla \times \mathbf{v}) \cdot (\nabla \times \mathbf{w}) + \rho(\nabla \times \mathbf{v}) \times \mathbf{v} \cdot \mathbf{w} - r\operatorname{div}(\mathbf{v}) - q\operatorname{div}(\mathbf{w})]dx \\ = - \int_{\partial\Omega_{f,in}} P_{in}\mathbf{w} \cdot \mathbf{n}ds - \int_{\partial\Omega_{f,out}} P_{out}\mathbf{w} \cdot \mathbf{n}ds \end{aligned} \quad (5.2.8)$$

where  $V_{\mathbf{v},q}$  is the functional space

$$V_{\mathbf{v},q} = \{(\mathbf{v}, q) \in H^1(\Omega, \mathbb{R}^d) \times L^2(\Omega)/\mathbb{R} \mid \mathbf{v} = 0 \text{ on } \Gamma \cup \partial\Omega_{f,wall} \text{ and } \mathbf{v} \times \mathbf{n} = 0 \text{ on } \partial\Omega_{f,in} \cup \partial\Omega_{f,out}\}.$$

We then consider the minimization problem (5.2.3) without the pressure loss constraint:

$$\min_{\Omega_f \subset D} J(\Omega_f). \quad (5.2.9)$$

Since the variational formulation (5.2.8) is different to that associated with the more standard Navier-Stokes equations (5.2.1), the shape derivative of  $J(\Omega_f)$  needs to be recomputed. Very briefly, the adjoint fluid variable, denoted  $(\mathbf{w}, r)$ , is obtained by solving the following variational problem:

Find  $(\mathbf{w}, r) \in V_{\mathbf{v},q}$  such that  $\forall(\mathbf{w}', r') \in V_{\mathbf{v},q}$ ,

$$\begin{aligned} \int_{\Omega_f} \nu(\nabla \times \mathbf{w}) \cdot (\nabla \times \mathbf{w}') + \rho(\nabla \times \mathbf{w}') \times \mathbf{v} \cdot \mathbf{w} + \rho(\nabla \times \mathbf{v}) \times \mathbf{w}' \cdot \mathbf{w} - r\operatorname{div}(\mathbf{w}') - r'\operatorname{div}(\mathbf{w})dx \\ = - \int_{\Omega_f} \rho c_p \mathbf{w}' \cdot \nabla T S dx + \frac{\partial \mathfrak{J}}{\partial(\mathbf{v}, q)} \cdot (\mathbf{w}', r'). \end{aligned} \quad (5.2.10)$$

The thermal adjoint variable  $S$  is computed by the previous equation of chapter 2, (2.4.9). It can then be shown that the shape derivative of an arbitrary functional  $J(\Gamma, \mathbf{v}, q, T)$  reads in surface form:

$$\frac{d}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=0} [J(\Omega_{\boldsymbol{\theta}}, \mathbf{v}_{\boldsymbol{\theta}}, q_{\boldsymbol{\theta}}, T_{\boldsymbol{\theta}})] \cdot \boldsymbol{\theta} = \frac{\partial \mathfrak{J}}{\partial \boldsymbol{\theta}} \cdot \boldsymbol{\theta} + \int_{\partial\Omega_f} \left( \nu(\nabla \times \mathbf{w}) \cdot (\nabla \times \mathbf{v}) + k_f \frac{\partial T}{\partial n} \frac{\partial S}{\partial n} \right) (\boldsymbol{\theta} \cdot \mathbf{n}) ds. \quad (5.2.11)$$

**Remark 5.3.** This rotational formulation does not allow to prescribe the stagnation pressure  $P + \rho\|\mathbf{v}\|^2/2$  on some parts of the boundary  $\partial\Omega_f$ , and the normal stress flux  $\sigma_f(\mathbf{v}, p) \cdot \mathbf{n}$  at some other parts, which could be also of industrial interest. However, this might be possible by resorting to a different variational treatment accounting for these boundary conditions, see [68].

## Numerical results

Our numerical results do not seem conclusive; the main reason lies in that the boundary condition on the stagnation pressure does not allow to bound sufficiently the values of the velocity in the whole domain. Indeed, the prescribed input and output pressure tends to make the flow velocity very large between close obstacles: the Reynolds number is not well controlled, which makes the numerical treatment difficult.

However when the optimization succeeds, the optimized shapes look similar to those obtained previously with the pressure loss imposed as a constraint. An instance of such shapes is illustrated on Figure 5.10 below. For these reasons, it seems to us preferable to stick to the first approach (5.2.3) (or

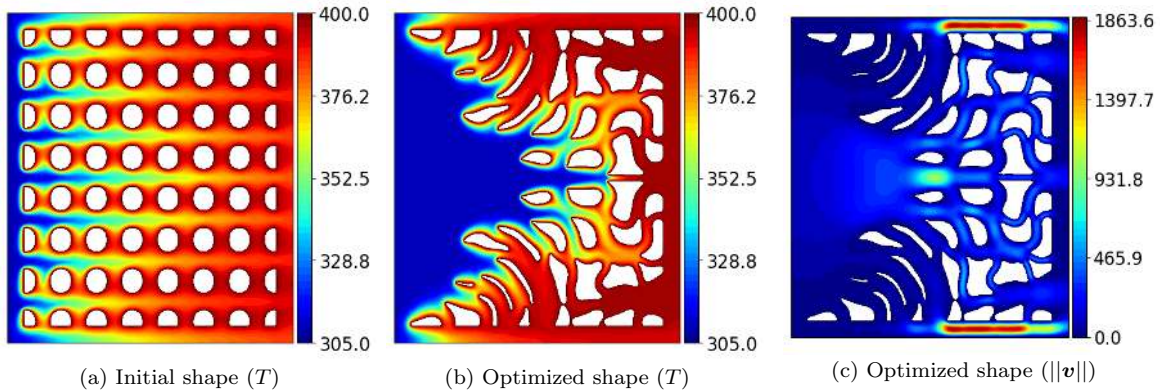


Figure 5.10: An optimization result for the shape optimization of the problem (5.2.9) featuring the Navier-Stokes equations (5.2.7) with a boundary condition on the stagnation pressure  $p + \rho v^2/2$ .

its reformulation (5.2.4) accounting for the minimum thickness constraint), which is also that followed by the majority of academics involving heat exchangers, e.g. [116, 255].

### 5.2.5 Conclusions

In this study, we have demonstrated the relevance of shape and topology optimization for generating unconventional heat exchanger designs. The method could be extended to other choices of physical models up to the use of different formulas for the shape derivatives: future work could take into account much larger, turbulent, Reynolds and Péclet numbers by considering turbulence models. Once again, we emphasize the ability of the method to take into account arbitrary and multiple optimization constraints. The extension to true 3-d test cases at low Reynolds number would not add additional theoretical difficulties in principle; however the implementation is delicate and still the object of ongoing work.

Let us mention, however, that several difficulties are still to be expected regarding the applications of such shape optimization methods to realistic industrial applications featuring large Reynolds and Péclet numbers (typically of the order of 30,000 or more).

First, when the Reynolds number is large, the assumption of stationary velocity and pressure fields does not make sense because it is well known that stationary solutions are unstable; the time-dependent solution varying rather on some attractor [304].

Second, even assuming the optimization to be numerically tractable, current heat exchanger designs suggest that optimal topologies would be characterized by a microstructure featuring very small details periodically repeated. Such designs would be very difficult to obtain with the method of Hadamard because very costly to mesh explicitly. It seems to us that a homogenization approach would be relevant for tackling such design problems. The recent works of [254, 166, 27] provide some perspectives regarding the automatic generation of optimized microstructures from homogenized models in the context of linear elasticity. In the final chapter 7, we provide a few preliminary theoretical contributions towards the application of such methods for fluid models.



## CHAPTER 6

# TOWARDS 3-D AND INDUSTRIAL APPLICATIONS: IMPLEMENTATION RECIPES FOR A VARIETY OF NUMERICAL TEST CASES

### Contents

---

<b>6.1 Implementation recipes for 3-d constrained topology optimization of multiphysics system</b> . . . . .	<b>209</b>
6.1.1 Abstract programming paradigm for constrained shape optimization . . . . .	210
6.1.2 Interfacing <code>python</code> and <code>FreeFEM</code> . . . . .	211
6.1.3 Processing operations on level set functions: generation of initial designs, taking into account non-optimizable regions, symmetrization and regularization . . . . .	213
6.1.4 Domain decomposition and preconditioning for 3-d variational problems . . . . .	216
<b>6.2 A few (moderately) large-scale three dimensional multiphysics applications</b> . . . . .	<b>219</b>
6.2.1 Cantilever beam subject to traction or torsion loads . . . . .	221
6.2.2 Optimal design for pure thermal heat conduction . . . . .	222
6.2.3 Lift–Drag topology optimization for aerodynamic design . . . . .	226
6.2.4 A 3-d fluid-structure interaction test case . . . . .	233

---

This chapter attempts to demonstrate the ability of multiphysics shape optimization by the method of Hadamard to deal with 3-d problems approaching industrial test cases.

From the numerical point of view, the extension of the shape optimization algorithm outlined in [chapter 1, section 1.4](#) from 2-d to 3-d is delicate. Several additional ingredients are required in order to achieve a satisfactory efficiency, or even to be able to achieve a single optimization run (for fluid applications). The first [section 6.1](#) outlines the most important features of our implementation: we namely discuss the use of domain decomposition techniques and preconditioning for the resolution of the 3-d multiphysics state equations involved in our shape optimization test cases. We also provide beforehand a brief presentation of our `python/FreeFEM` implementation, and some details about various technical operations (rarely described in the literature) applied in the course of optimization iterations in order to enforce non-optimizable regions, symmetry and non-degeneracy of discretized shapes.

In the next [section 6.2](#), we present a variety of 3-d test cases solved thanks to our implementation. Four examples are considered: the first three of them are single physics applications in either linear elasticity, heat conduction, or fluid mechanics. The fourth one is a true multiphysics example; it is concerned with fluid-structure interaction. These test cases are described as *moderately* large-scale, in the sense that the problems considered make extensive use of parallel computing and preconditioning in order to be run in reasonable CPU time, however their size remains rather small (our largest test case features up to 1.7 millions degrees of freedom) when compared to that of industrial problems (reaching about the billion of degrees of freedom). Our results, however, are promising and are preliminary to more challenging applications.

### 6.1 IMPLEMENTATION RECIPES FOR 3-D CONSTRAINED TOPOLOGY OPTIMIZATION OF MULTIPHYSICS SYSTEM

In this section, we discuss several aspects of our implementation in `python` (version 3.6) and `FreeFEM` [183] for 3-d shape and topology optimization of coupled thermal-fluid elastic systems. The first two [sections 6.1.1](#) and [6.1.2](#) provide information regarding several choices of programming paradigms. [Section 6.1.3](#) describes several “hidden” operations applied to shapes in the course of optimization iterations, in order to account for non-optimizable regions, symmetry, or to avoid mesh degeneracy. Finally, our parallel implementation of finite element operations using domain decomposition and preconditioning is discussed in [section 6.1.4](#).

### 6.1.1 Abstract programming paradigm for constrained shape optimization

A single implementation of the null space gradient flow algorithm as described in [chapter 3, algorithm 3.1](#) is used for both 2-d and 3-d problems, which could in principle be used for any optimization problem set on a *manifold* as soon as a generic minimal set of ingredients (described below) is provided.

In a few words, the algorithm was implemented in `python` via a `nullspace` function whose prototype reads

```
def nullspace(problem: Optimizable, params=None, results=None):
    """
    Solve the optimization problem
    min      J(x)
    x in V
    under the constraints
    g_i(x)=0 for all i=1..p
    h_j(x)<=0 for all j=1..q

    problem: an instance of the class Optimizable
    params: a dictionary of optimization parameters
    results: a previous output of the nullspace function
            (the algorithm will restart from the last iteration)
    """
```

The most important argument of the function `nullspace` is the variable `problem` which instantiates an abstract class `Optimizable`. An `Optimizable` object encodes a generic optimization problem

$$\begin{aligned} \min_{x \in \mathcal{X}} \quad & J(x) \\ \text{s.t.} \quad & \begin{cases} g_i(x) = 0 \text{ for } 1 \leq i \leq p \\ h_j(x) \leq 0 \text{ for } 1 \leq j \leq q. \end{cases} \end{aligned} \tag{6.1.1}$$

The set  $\mathcal{X}$  is assumed to be a manifold equipped with the structure outlined in [chapter 3, section 3.6.1](#) whose representative illustration is reproduced on [Figure 6.1](#) below with the current notation. In order

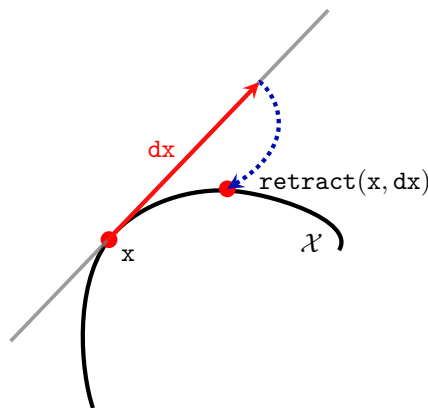


Figure 6.1: Tangential motion and retraction on an abstract manifold  $\mathcal{X}$ .

to solve (6.1.1), it is sufficient that the following information be provided (e.g. as object methods) by the supplied instance `problem` of an `Optimizable` object:

- $J, g_i, h_j : \mathcal{X} \rightarrow \mathbb{R}$ : objective functions and constraints;
- $DJ, Dg_i, Dh_j : \mathcal{X} \rightarrow \mathbb{R}^n$ : Fréchet derivatives of objective and constraints as functions. Here,  $n$  is thought of as the dimension of the tangent space of  $\mathcal{X}$  to  $x \in \mathcal{X}$ . For such a given  $x$ ,  $DJ(x)^T dx$  is the variation of the objective function  $J$  at  $x$  along the tangent direction  $dx \in \mathbb{R}^n$ ;
- $A : \mathcal{X} \rightarrow \mathbb{R}^{n \times n}$ : local inner product needed for the computation of gradients

$$\nabla J(x) := A(x)^{-1}DJ(x), \quad \nabla g_i(x) := A(x)^{-1}Dg_i(x), \quad \nabla h_j(x) := A(x)^{-1}Dh_j(x) \in \mathbb{R}^n.$$

Of course,  $A(\mathbf{x})$  must be a symmetric positive definite matrix.

- **retract** :  $\mathcal{X} \times \mathbb{R}^n \rightarrow \mathcal{X}$ : a *retraction* that convert a current point  $\mathbf{x}$  and a tangent vector  $\mathbf{dx} \in \mathbb{R}^n$  into a new point  $\text{retract}(\mathbf{x}, \mathbf{dx}) \in \mathcal{X}$  on the manifold  $\mathcal{X}$ . This mapping must be compatible with the previous definitions of the derivatives  $DJ, Dg_i$  and  $Dh_j$  in the following sense:

$$J(\text{retract}(\mathbf{x}, \Delta t \mathbf{dx})) = J(\mathbf{x}) + \Delta t DJ(\mathbf{x})^T \mathbf{dx} + o(\Delta t) \text{ as } \Delta t \rightarrow 0.$$

- **accept**: an optional function that is called by the optimization algorithm when the next point  $\mathbf{x} \in \mathcal{X}$  is accepted, which serves e.g. for saving current available information before proceeding to the next iteration.

The above ingredients constitute all the necessary information required by any first order optimization algorithm acting on a manifold (equipped with a retraction), and in particular by the null space algorithm outlined in [chapter 3, algorithm 3.1](#). In our implementation for shape optimization, the current guess  $\mathbf{x}$  contains the path to the current mesh file of the optimized domain, and  $\mathbf{dx}$  is a finite element vector discretizing a deformation field. The function **A** returns the finite element matrix associated with the regularization, while **retract** encodes the advection and remeshing step. When passing from 2-d to 3-d, or if we were to rely on other numerical representations of shapes (e.g. using the level set method on a fixed), it is sufficient to update the above functions **A** and **retract** accordingly. This programming paradigm allows us to implement separately 2-d and 3-d test case (which required different solvers for the physical state equations), while calling the same null space optimization algorithm.

### 6.1.2 Interfacing python and FreeFEM

Our treatment of 2-d and 3-d shape optimization test cases with the method of Hadamard is implemented in both programming languages `python` and `FreeFEM` [183]:

1. on the one hand, the `nullspace` optimization routine of the previous section is conveniently implemented in `python`, which is a very user-friendly language and which allows for easy debugging. For instance, it is possible to pause a running instance of the code at any step of the optimization process (such as calling `mmg` for remeshing, computing the signed distance function with `mshdist`, etc...);
2. on the other hand, `FreeFEM` is used for all finite element related operations: the assembly of sparse matrices discretizing variational forms, integration on meshed subdomains or boundaries, the resolution of linear systems, the use of domain decomposition methods, mesh interpolation, etc... In particular, the language relies on a `C++` kernel which allows to run these operations very efficiently. Furthermore, the syntax of the language is very close to mathematics and it allows to easily implement the expressions of the shape derivatives of [chapter 2, propositions 2.3 and 2.4](#) (see also [34] for detailed examples of the use of `FreeFEM` in topology optimization).

The interface between both languages is realized thanks to a preprocessing meta-language: special instructions are added to `FreeFEM` (non executable) source files which are then parsed by a `python` routine and converted into a proper `.edp` file to be executed by `FreeFEM`. An example of such an augmented `.edp` code taken from our 2-d implementation is provided below:

```
SET (FLUID_DOMAIN,"3")
DEFAULT (THERMIC_ENABLED,"0")
IF THERMIC_ENABLED
  solve thermic(T,S)=
    int2d(Th,$FLUID_DOMAIN) (kf*grad(T) '*grad(S)
  IF FLUID_ENABLED
    \\\+rho*cp*grad(T) '*[vx,vy]*S
  ENDIF
  \\\)
  +int2d(Th,$SOLID_DOMAIN) (ks*grad(T) '*grad(S))
IFDEF Qf
  -int2d(Thf) ($Qf*S)
ENDIF
```

```

IFDEF Qs
  -int2d(Ths)($Qs*S)
ENDIF
IFDEF h
  +int1d(Th,$BCTN)($h*S)
ENDIF
IFDEF TO
  +on($BCTD,T=$TO)
ENDIF
\\;
ENDIF

```

This piece of code implements the variational form of [chapter 2, \(2.4.3\)](#):

$$\begin{aligned}
 \text{thermic}(T,s) &= \int_{\Omega_s} k_s \nabla T \cdot \nabla S dx + \int_{\Omega_f} (k_f \nabla T \cdot \nabla S + \rho c_p S \mathbf{v} \cdot \nabla T) dx \\
 &= \int_{\Omega_s} Q_s S dx + \int_{\Omega_f} Q_f S dx + \int_{\partial\Omega_T^N} h S ds.
 \end{aligned}$$

The blue keywords SET, DEFAULT, IF, etc. . . correspond to preprocessing instructions parsed by `python`: they determine sections of the code to be included in the final executable depending on the values of the preprocessing variables prefixed by the dollar ‘\$’ symbol. For instance, the term  $\rho c_p S \mathbf{v} \cdot \nabla T$  will not be included in the final code if `$FLUID_ENABLED` is not set to 1. The dollar prefixed variables can be easily accessed from the `python` implementation, or assigned directly in the FreeFem code with the instructions SET or DEFAULT. The double backslash symbols ‘\’ indicates that the previous line carriage return must be removed. For instance, if `$THERMIC_ENABLED=1`, `$FLUID_ENABLED=1`, `$TO=100` and `$BCTD=1`, then the above source code is converted by our `python` parser as the following “standard” FreeFEM source code

```

solve thermic(T,S)=
  int2d(Th,3)(kf*grad(T)*grad(S)+rho*cp*grad(T)*[vx,vy]*S)
  +int2d(Th,2)(ks*grad(T)*grad(S))
  +on(1,T=100);

```

Such a preprocessing language proves to be very convenient in order (i) to use a single source code for all our test cases whatever the number of considered physical equations and (ii) to easily change test case parameter values from either `python` or FreeFEM thanks to the dollar prefixed variables and (iii) maintain a good readability of the implementation.

FreeFEM source files augmented with preprocessing instructions are called by the various functions `J,DJ, A, retract`, etc. . . declared in the implementation of an `Optimizable` object (see the previous [section 6.1.1](#)). Below is reported a typical output of a running instance of our 2-d code, which offers a clear picture of all the elementary steps of the shape optimization algorithm outlined in [chapter 1, algorithm 1.1](#).

```

0. J=0.1079 G=[0.302] H=[]
FreeFem++ 01_cantilever/01_run/scripts/sensitivities.edp -nw (0.44)
FreeFem++ 01_cantilever/01_run/scripts/scalar_product.edp -nw (0.29)
advect -nocfl -dt 1.0 01_cantilever/01_run/meshes/Th_0000.mesh -c
  01_cantilever/01_run/scalars_dir/phi_0000.chi.sol
  -s 01_cantilever/01_run/tmp/theta.sol
  -o 01_cantilever/01_run/tmp/phi.o.sol (0.08)
FreeFem++ 01_cantilever/01_run/scripts/symmetrize.edp -nw (0.27)
mmg2d_03 -nr -hmin 0.02 -hmax 0.1 -hgrad 1.3 -hausd 1e-4 -ls
  -sol 01_cantilever/01_run/tmp/phi_tmp.sol
  -out 01_cantilever/01_run/tmp/Th.o.mesh
  01_cantilever/01_run/meshes/Th_0000.mesh (0.16)
mv 01_cantilever/01_run/tmp/Th.o.mesh 01_cantilever/01_run/tmp/Th_tmp.mesh (0.00)
mshdist -v 0 -ncpu 1 01_cantilever/01_run/tmp/Th_tmp.mesh -dom -fmm (0.03)

```

```
FreeFem++ 01_cantilever/01_run/scripts/solve_state.edp -nw (0.43)
[...]
```

### 6.1.3 Processing operations on level set functions: generation of initial designs, taking into account non-optimizable regions, symmetrization and regularization

We now describe a few very classical operations that are implemented in most level-set based topology optimization codes but which are not often detailed in published works. The context and notation assumed are that of [chapter 2](#): a hold-all domain  $D \subset \mathbb{R}^d$  ( $d = 2$  or  $3$ ) is considered, and the goal is to optimize the shape of the interface  $\Gamma = \partial\Omega_s \cap \partial\Omega_f$  between a solid subdomain  $\Omega_s \subset D$  and a fluid subdomain  $\Omega_f = D \setminus \overline{\Omega_s}$ .

#### Generation of initial meshes and initial designs

Our implementation relies on a fully explicit mesh discretization of the hold-all domain  $D$  featuring the design shape  $\Omega_f \subset D$  as a submesh. This is very convenient for solving finite element problems on  $\Omega_f$  or  $\Omega_s$ , since boundary conditions can then be applied in a straightforward manner on the boundary  $\partial\Omega_f$  or  $\partial\Omega_s$ .

It is worth mentioning that generating an initial mesh featuring a correct topology respecting user-defined labels on the boundary triangles is an art by itself: there exist several software programs in order to do this (such as [TetGen](#) [283] or [gmsht](#) [167]). For our purposes, we found very convenient to do it in a level-set fashion assisted with the help of the library [mmg](#) [108]:

1. a first mesh  $\mathcal{T}_0$  for the computational domain  $D$  is generated. Usually,  $D$  is a box so that this step is very easily achieved, for instance using a [FreeFEM](#) command of the kind

```
// Generate a mesh for the box domain [0,2]*[0,1]*[0,1]
mesh3 Th = cube(60, 40, 40, [x*2,y,z]);
```

2. a level set function  $\phi$  is generated on the domain  $D$  such that the negative subdomain of  $\phi$  on  $\partial D$ ,

$$\{x \in \partial D \mid \phi_0(x) < 0\}$$

delimits a set of connected components on the boundary of  $D$  associated with various boundary conditions. The level set of  $\phi_0$  is explicitly discretized within  $\mathcal{T}_0$  and remeshed with [mmg](#) so that a new mesh  $\mathcal{T}_1$  where these boundary patches are explicitly meshed is obtained;

3. the obtained boundary patches are tagged with specific labels according to the desired boundary conditions (see e.g. the red patches on [Figure 2.3](#) below);
4. an initial design  $D = \Omega_s \cup \Omega_f$  is proposed by the user under the form of a level-set function  $\phi_1$  (satisfying  $\phi_1 \leq 0$  in  $\Omega_f$ ). An initial mesh  $\mathcal{T}$  of  $D$  featuring the initially proposed design for  $\Omega_f$  discretized as a submesh and the correct boundary labels is obtained by discretizing the negative subdomain of  $\phi_1$  in  $\mathcal{T}_1$ . Note that boundary labels do not need to be rewritten neither at this step nor at subsequent stages of the optimization because the remeshing library [mmg](#) preserves meshed interfaces between any two regions labeled with distinct tags.

Classically, complex initial domains can be generated by applying min-max operations to level-set functions associated with elementary shapes (such as sphere or half-spaces), since minimum and maximum of two level-set functions amount to considering to respectively the intersection and the reunion of the corresponding subdomains [247].

#### Non optimizable subdomains

It is very customary to impose that some non-optimizable subdomains  $\omega$  belong either to the fluid or the solid part  $\Omega_f$  or  $\Omega_s$ . A particular attention must be given to the treatment of this constraint during the computation step of a descent direction  $\theta \in W^{1,\infty}(D, \mathbb{R}^d)$ : this deformation  $\theta$  should satisfy

$$\theta \cdot n_\omega \geq 0 \text{ on } \partial\omega \tag{6.1.2}$$

where  $\mathbf{n}_\omega$  denotes the outward normal to  $\omega$ . The constraint (6.1.2) ensures that the boundary of the optimized shape does not penetrate the non optimizable region  $\omega$ . In principle, this requirement could be imposed as an additional constraint in the mathematical program determining the current descent direction from the knowledge of the shape derivatives. For instance, in the context of the unconstrained minimization of an objective functional  $J(\Gamma)$ , a “best” descent direction  $\boldsymbol{\theta}$  taking into account the non optimizable region  $\omega$  could be obtained by solving the quadratic minimization problem

$$\begin{aligned} \min_{\boldsymbol{\theta} \in V} \quad & DJ(\Gamma)(\boldsymbol{\theta}) \\ \text{s.t.} \quad & \begin{cases} \boldsymbol{\theta} \cdot \mathbf{n}_\omega \geq 0 \\ \|\boldsymbol{\theta}\|_V \leq 1, \end{cases} \end{aligned} \tag{6.1.3}$$

where  $\|\cdot\|_V$  is the regularizing norm considered for the identification of  $DJ(\Gamma)$  with a gradient (see chapter 1, section 1.4.1).

In practice, it is easier to solve (6.1.3) by replacing the inequality constraint  $\boldsymbol{\theta} \cdot \mathbf{n}_\omega \geq 0$  with an equality constraint  $\boldsymbol{\theta} \cdot \mathbf{n}_\omega = 0$ , or even by simply imposing that  $\boldsymbol{\theta} = 0$  on  $\omega$ . This is precisely what we do in all our numerical examples with satisfactory results. The descent direction  $\boldsymbol{\theta}$  solving the optimization problem (6.1.3) is then used as an advection field driving the evolution of a level-set function  $\phi$  for the optimized shape  $\Omega_f$  ( see the step 2 of the algorithm algorithm 1.2 recalled in chapter 1). In principle, imposing  $\boldsymbol{\theta} \cdot \mathbf{n}_\omega \geq 0$  on  $\partial\omega$  should guarantee at the continuous level that the interface of the optimized shape do not penetrate  $\partial\omega$ . Unfortunately, due to numerical inaccuracies, too large time-steps, or incorrect initial designs, a post-treatment of the advected function is used in order to correct small violations of the non-optimizable region. Let us introduce

- a level-set function  $\Phi_\omega$  associated with the non-optimizable region  $\omega$ , i.e.

$$\forall x \in \omega, \Phi_\omega(x) < 0 \text{ and } \forall x \in D \setminus \bar{\omega}, \Phi_\omega(x) \geq 0;$$

- a level-set function  $\Psi$  associated with the distribution  $\mathcal{X} \subset \omega$  of material desired in the non-optimizable region of  $\omega$ , i.e. it should always hold

$$\{x \in \omega \mid \Psi(x) < 0\} \subset \Omega_f \text{ and } \{x \in \omega \mid \Psi(x) > 0\} \subset \Omega_s.$$

Then denoting by  $\phi$  the obtained level set function for the next iteration after the advection step, we perform the operation (see Figure 6.2 below)

$$\phi \leftarrow \min(\max(\phi, -\Phi_\omega), \max(\Psi, \Phi_\omega)), \tag{6.1.4}$$

which corresponds to the domain update

$$\Omega_f \leftarrow (\Omega_f \cap (D \setminus \bar{\omega})) \cup \mathcal{X}. \tag{6.1.5}$$

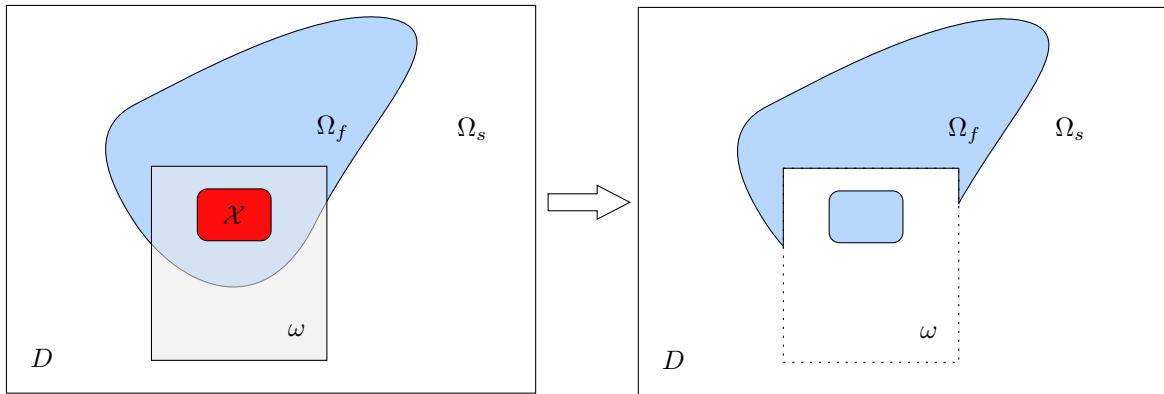


Figure 6.2: Illustration of the level-set operation (6.1.4) enforcing non-optimizable regions  $\omega$ : the distribution of material  $\Omega$  inside the domain  $\omega$  should match exactly the red set  $\mathcal{X}$ . After the operation,  $\Omega$  is the new blue color domain on the right.

Note that the post-treatment (6.1.4) is only a “projection” of  $\Omega_f$  onto the set of shapes fulfilling the requirement that  $\Omega_f \cap \omega = \mathcal{X}$ . To be consistent with this operation, i.e. to guarantee that it does not affect the decreasing property (6.1.3) of  $\theta$ , it is very important to impose the condition (6.1.2), or at least  $\theta \cdot \mathbf{n}_\omega = 0$  on  $\partial\omega$ .

**Remark 6.1.** In the context where a fixed computational domain  $D$  containing the optimized shape  $\Omega$  is used, the boundaries of  $D$  should be set as non-optimizable as well, because deformations making  $\Omega$  larger than  $D$  should not be allowed. A slight difference of implementation was considered regarding this point whether we use either the surface or the volume expression of the shape derivative:

- if the surface expression (2.4.14) of the shape derivative is used, then the considered space of deformation fields is

$$V = \{v \nabla d_{\Omega_f} \mid v \in H^1(D)\},$$

that is the set of vector fields aligned with the constant extension of the normal to  $\partial\Omega_f$ . Since shape derivatives are written as surface integrals over  $\Gamma \subset \partial\Omega_f$ , and we make sure numerically that  $\Gamma$  never becomes a subset of  $\partial D$ , no constraint is required on the value of  $v$  regarding the fact that  $\partial D$  is a non-optimizable boundary.

- if the volume expression (2.4.14) of the shape derivative is used, then

$$V = H^1(D, \mathbb{R}^d)$$

and it is necessary to enforce that  $\theta \cdot \mathbf{n}_D \geq 0$  on  $\partial D$  where  $\mathbf{n}_D$  is the outward normal to  $D$ . In our implementation, we rather enforce  $\theta \cdot \mathbf{n}_D = 0$  in the identification problem (6.1.3).

### Enforcing symmetries of the iterated domains

Many optimization test cases feature inherent domain symmetries; it is then expected that locally optimal designs for these problems should also be symmetric. In many contributions of the literature, symmetry is used in order to reduce the size of the computational domain: the overall optimization is performed on only one part of the domain  $D$ , and the final design on  $D$  as a whole is inferred by symmetry.

However, it is sometimes desirable to perform the computations involved in the resolution of the shape optimization problem on the whole domain  $D$  and to symmetrize shapes as a post-treatment. Such a post-treatment is needed because in general, numerical inaccuracies tend to quickly make optimized shapes nonsymmetric when the computational mesh is not symmetric. In our implementation where each intermediate design is obtained with a level set function  $\phi$ , the symmetry of  $\phi$  with respect to a given symmetry transformation  $S_1$  (satisfying  $S_1 \circ S_1 = I$ ) is enforced by the operation

$$\phi \leftarrow \frac{1}{2}(\phi + \phi \circ S_1) \quad (\text{for one symmetry } S_1).$$

Likewise, invariance of the design with respect to two symmetry transformations  $S_1, S_2$  (satisfying  $S_1 \circ S_1 = S_2 \circ S_2 = I$ ) is imposed via

$$\phi \leftarrow \frac{1}{5}(\phi + \phi \circ S_1 + \phi \circ S_2 + \phi \circ S_1 \circ S_2 + \phi \circ S_2 \circ S_1) \quad (\text{for two symmetries } S_1, S_2).$$

The compositions  $\phi \circ S_1, \phi \circ S_2$  are computed by interpolation (which is an automatic feature in **FreeFEM**). These transformations ensure that the updated level set function satisfies  $\phi \circ S_1 = \phi$  and  $\phi \circ S_2 = \phi$ .

### Regularization of the level-set function to avoid mesh degeneracy

A very important ingredient for handling topological changes in the context of the mesh evolution method of [24] is a regularization of the level-set function. After the step 2 of **algorithm 1.1**, a level-set function  $\phi$  is obtained at the vertices of a computational mesh of  $D$  whose negative subdomain corresponds to the new shape. A new mesh of  $D$  is then created in which this updated shape exists as a submesh (step 3 in **algorithm 1.1**). Before performing this remeshing step,  $\phi$  is regularized into a new level set function  $\tilde{\phi}$  by solving the elliptic problem

$$\text{Find } \tilde{\phi} \in H^1(D), \text{ such that } \forall \psi \in H^1(D), \int_D (\gamma^2 \nabla \tilde{\phi} \cdot \nabla \psi + \tilde{\phi} \psi) dx = \int_D \phi \psi dx. \quad (6.1.6)$$

The regularizing length scale is typically set to  $\gamma = 0.01\mathbf{hmin}$  where  $\mathbf{hmin}$  is the desired *minimum* mesh edge size. It may appear surprising that one regularizes  $\phi$  at a scale lower than the mesh size. In fact, the minimum mesh size  $\mathbf{hmin}$  provided to `mmg` might be violated if the input topology of the mesh prescribes it: for instance, `mmg` would never remove a bubble of arbitrary size in an input mesh because the topology of the computed output mesh would then change (see [108]). If the Hadamard's shape derivative dictates to continuously decrease the size of these bubbles at every optimization step, the resulting successive meshes quickly degenerate and the finite element resolution becomes impossible. Removing these bubbles (or any other topological features much smaller than  $\mathbf{hmin}$ ) by the pretreatment (6.1.6) allows to avoid these issues.

### Handling connectivities of the meshes for finite element related operations.

A particular care must be paid in the implementation of the variational formulations (2.4.2) to (2.4.4) for the state equations, or in the calculation of the shape derivatives (2.4.13) and (2.4.14). Indeed, the state variables at play  $(\mathbf{v}, p), T, \mathbf{u}$  (as well as the adjoint variables) are living on different meshes: the fluid variables  $(\mathbf{v}, p)$  and the elastic variable  $\mathbf{u}$  are for instance solved on submeshes of the global meshed domain  $D$  corresponding respectively to the fluid or to the solid subdomains  $\Omega_f, \Omega_s$ . As a result, interpolation matrices are extensively used in order to transfer nodal information from one mesh to another (e.g. for transferring the value of  $\sigma_f(\mathbf{v}, p)$  on the *fluid* interface  $\Gamma$  seen as a boundary of the meshed domain  $\Omega_f$ , to the *solid* interface  $\Gamma$  seen as a boundary of the meshed domain  $\Omega_s$ ). This is achieved in `FreeFEM` thanks to the instruction `interpolate`;

#### 6.1.4 Domain decomposition and preconditioning for 3-d variational problems

Passing from the 2-d implementation to its 3-d counterpart is *theoretically* without difficulty, but it requires in practice a substantial amount of effort. We shall not discuss the (quite important) differences between 2-d and 3-d regarding *remeshing* issues (the reader is again referred to [108]); we shall focus instead on the difficulties related to the resolution of variational problems by the finite element method.

The cornerstone of the passage from 2-d to 3-d lies in the assembly and inversion of large sparse linear systems obtained from the discretization of the physical equations of chapter 2, (2.4.2) to (2.4.4). Generally, linear systems resulting from 2-d applications are sufficiently small so that a direct factorization based method can be used [305]. In 3-d, it is possible to use direct methods only for very low resolution problems: indeed, modern direct solvers based on LU factorizations such as `MUMPS` [38] have a complexity of order  $O(N^2)$  ( $N$  is the number of degrees of freedom of the finite element approximation) which becomes quickly too expensive in terms of both CPU time and required memory. In this context, one classically resorts to *iterative* methods merely based on matrix-vector products, which are relatively inexpensive to compute due to the sparsity of the matrices involved in the context of finite element problems. The most popular iterative methods are the conjugate gradient method (CG) for symmetric positive definite problems, and GMRES for the general non symmetric case.

The treatment of large sparse systems as usually encountered in large scale 3-d applications involves two additional ingredients:

1. *preconditioning*: iterative methods may take many iterations to converge in reasonable CPU time for ill-conditioned linear systems [173]. It is often possible to accelerate the resolution of such linear systems

$$Ax = b, \tag{6.1.7}$$

by left or right multiplying the (large) square matrix  $A$  and the vector  $b$  with a *preconditioner*  $M$ :

$$MAx = Mb \text{ or } AMy = b. \tag{6.1.8}$$

A good preconditioner  $M$  is a square matrix approximating well the inverse of  $A$ . In that case, it is expected that iterative methods applied to (6.1.8) will converge in much less iterations than when applied to the original problem (6.1.7).

2. *domain decomposition*: the computational domain  $D$  is divided into a number  $m$  of subdomains:  $D = D_1 \cup \dots \cup D_m$ . Very roughly, the inverses of the "restrictions"  $A_1, \dots, A_m$  of the finite element



matrix  $A$  to these subdomains are used to build a block preconditioner

$$M := \begin{bmatrix} A_1^{-1} & & \\ & \ddots & \\ & & A_m^{-1} \end{bmatrix}.$$

This allows to distribute the resolution of the linear system (6.1.7) on multiple cpus, because all operations involving the restriction inverses  $A_1^{-1}, \dots, A_m^{-1}$  (or their approximation using a limited number of GMRES iterations) can be performed in parallel.

For our applications, we rely on the PETSc library [54, 55, 56] and its interface in FreeFEM developed by Jolivet [199, 132] which allows to solve finite element problems with a large library of state-of-the-art preconditioned iterative methods. The domain decomposition step was achieved thanks to a macro `buildMinimalist` (see [198]):

```
// Partition the mesh Th for the finite element space Pk and compute
// associated partition of unity D and connectivity arrays arrayIntersection
// and restrictionIntersection
buildMinimalist(Th, D, arrayIntersection, restrictionIntersection, Pk);
//Now, Th is one of the submeshes of the domain decomposition [...]
```

Importantly, the use of these techniques requires a significant amount of effort, because the whole implementation needs to be thought parallel for scalability; this includes operations ranging from the finite element matrix assembly, the evaluation of volume or surface integrals, up to the numerical assembly and regularization of shape derivatives (2.4.13) and (2.4.14).

We now detail the choice (physics dependent) preconditioners for the weakly coupled fluid thermal mechanical system (2.2.1) to (2.2.3). For these matters, we have been very much assisted by Pierre Jolivet and his tutorial on FreeFEM [198].

**Linear elasticity:** the linear elasticity system (2.2.3) is an elliptic vector system. The Geometric Algebraic Multigrid preconditioner (GAMG) is known to be very efficient for solving linear elasticity problems [8] (among other smoothed aggregation methods). The idea of multigrid methods is to construct a preconditioner by using the (iterative) inverses of the elasticity operator restricted to coarser meshes; these are much cheaper to compute and allow to obtain good approximations of the low frequency component of the solution.

We used the GAMG preconditioner in our FreeFEM implementation by calling

```
// set GAMG preconditioner for elasticity matrix AElasticity
string petsc_options_elasticity = "-pc_type gamg -ksp_type gmres -ksp_max_it 200"
    + " -pc_gamg_threshold 0.01";
set(AElasticity, sparams = petsc_options_elasticity, nearnullspace = Rb);
```

The variable `Rb` contains the near null space of the elasticity operator (i.e. rigid body modes on the whole computational domain), it is required to help the preconditioner to achieve good performance.

**Heat conduction-convection** For the convective heat problem (2.2.2), we rely on the pointwise aggregation multigrid preconditioner (BoomerAMG) implemented in the `hypr` library [149] which is supposed to be slightly more efficient than `gamg` for scalar problems (in fact, for our moderately large applications, the difference was not obvious). This preconditioner was used in our implementation by calling

```
// set BoomerAMG preconditioner for the thermic system matrix AThermic
string petsc_options_thermic = "-pc_type hypr";
set(AThermic, sparams = petsc_options_thermic, nearnullspace = Rb);
```

This time, the near null space variable `Rb` corresponds to the constant functions. Note that for these first two considered physics, we observed that the use of the standard conjugate gradient method converges well without the need for a preconditioner (the latter allows nevertheless to significantly speed

up this convergence). This situation is very particular to the fact that the partial differential operators corresponding to these physics are elliptic (if no convection is involved).

**Navier-Stokes equations** Solving the 3-d steady state Navier-Stokes equations is much more challenging than the previous elasticity and heat conduction problems, because unpreconditioned iterative solvers fail to converge in a reasonable amount of time.

Our implementation of the resolution of the steady state incompressible Navier-Stokes equations relies on the Augmented Lagrangian Preconditionner recently described in [238, 279] and for which the source code is available in `FreeFEM`. The main ingredient of this method is the addition of a penalization term for the divergence constraint in the variational formulation (2.4.2) associated with the nonlinear Navier-Stokes problem (2.2.1):

$$\text{Find } (\mathbf{v}, p) \in V_{\mathbf{v},p}(\Gamma) \text{ such that } \forall (\mathbf{w}, q) \in V_{\mathbf{v},p}(\Gamma),$$

$$\int_{\Omega_f} [\sigma_f(\mathbf{v}, p) : \nabla \mathbf{w} + \rho \mathbf{w} \cdot \nabla \mathbf{v} \cdot \mathbf{v} - q \operatorname{div}(\mathbf{v}) - p \operatorname{div}(\mathbf{w}) + \gamma \operatorname{div}(\mathbf{v}) \operatorname{div}(\mathbf{w})] dx = \int_{\Omega_f} \mathbf{f}_f \cdot \mathbf{w} dx. \quad (6.1.9)$$

The coefficient  $\gamma > 0$  penalizes the constraint  $\operatorname{div}(\mathbf{v}) = 0$ . The resolution of the nonlinear problem (6.1.9) is performed with the Newton method: at each step  $k \geq 0$  of the process, an increment  $(\delta \mathbf{v}_k, \delta p_k)$  is computed by solving the linearization of (6.1.9) around  $(\mathbf{v}_k, p_k)$ :

$$\text{Find } (\delta \mathbf{v}_k, \delta p_k) \in V_{\mathbf{v},p}(\Gamma) \text{ such that } \forall (\mathbf{w}, q) \in V_{\mathbf{v},p}(\Gamma),$$

$$\int_{\Omega_f} [\sigma_f(\delta \mathbf{v}_k, \delta p_k) : \nabla \mathbf{w} + \rho \mathbf{w} \cdot \nabla \mathbf{v}_k \cdot \delta \mathbf{v}_k + \rho \mathbf{w} \cdot \nabla(\delta \mathbf{v}_k) \cdot \mathbf{v}_k - q \operatorname{div}(\delta \mathbf{v}_k) - \delta p_k \operatorname{div}(\mathbf{w})] dx$$

$$+ \int_{\Omega_f} \gamma \operatorname{div}(\delta \mathbf{v}_k) \operatorname{div}(\mathbf{w}) dx = \int_{\Omega_f} \mathbf{f}_f \cdot \mathbf{w} dx. \quad (6.1.10)$$

The next iterate  $(\mathbf{v}_{k+1}, p_{k+1})$  is then obtained by setting

$$\mathbf{v}_{k+1} := \mathbf{v}_k + \delta \mathbf{v}_k, \quad p_{k+1} := p_k + \delta p_k.$$

In our implementation, the initial guess  $(\mathbf{v}_0, p_0)$  is the solution of the Stokes counterpart problem to (6.1.9) (obtained with  $\rho = 0$ ). The difficult part of the method is the resolution of the so-called Oseen problem (6.1.10) which turns to be very poorly ill-conditioned. For our application we rely on the Augmented Lagrangian preconditioner of [238]. It relies on the use of a suitable preconditioner for the block matrix discrete operator associated to the problem (6.1.10), which is of the form:

$$\mathbf{A0seen} = \begin{bmatrix} A & B^T \\ B & 0 \end{bmatrix} \quad (6.1.11)$$

where  $A$  is the matrix discretizing the bilinear form

$$(\delta \mathbf{v}_k, \mathbf{w}) \mapsto \int_{\Omega_f} [\sigma_f(\delta \mathbf{v}_k, \delta p_k) : \nabla \mathbf{w} + \rho \mathbf{w} \cdot \nabla \mathbf{v}_k \cdot \delta \mathbf{v}_k + \rho \mathbf{w} \cdot \nabla(\delta \mathbf{v}_k) \cdot \mathbf{v}_k + \gamma \operatorname{div}(\delta \mathbf{v}_k) \operatorname{div}(\mathbf{w})] dx$$

and  $B$  is the discretization of the divergence operator  $(\mathbf{w}, q) \mapsto \int_{\Omega_f} q \operatorname{div} \mathbf{w} dx$ . The linear solver GMRES is very slow to converge if it is called directly on (6.1.11). A preconditioner is obtained by using the inverse of  $A$  (or an approximation of it); following [238], we use the (right) preconditioner `M0seen` given by

$$\mathbf{M0seen} := \begin{bmatrix} A^{-1} & 0 \\ B & S, \end{bmatrix} \quad (6.1.12)$$

where the inverse  $A^{-1}$  is estimated thanks to an additive Schwarz method (ASM; see [313]). The matrix  $S$  is an approximation of the pressure Schur complement  $S = -BA^{-1}B^T$ , it is estimated from its inverse  $S^{-1} = -(\gamma + 1/\operatorname{Re})M_p$  where  $M_p$  is the mass matrix associated with  $(p, q) \mapsto \int_{\Omega_f} pq dx$  and  $\operatorname{Re} = \rho/\nu$  is the Reynolds number (see [238] for the details). Many variants of this technique using  $A^{-1}$  to obtain a preconditioner exist, which can be related to Uzawa-type methods [52]. Note that setting  $\gamma > 0$  in (6.1.9) is required for the (flexible) GMRES solver to converge.

To be very specific, this preconditioner is implemented in `FreeFEM` thanks to the following lines of code:

```
// set Augmented Lagrangian preconditioner for the Oseen matrix AOseen
string paramsOseen = " -ksp_type fgmres -pc_type fieldsplit "
+ " -pc_fieldsplit_type multiplicative -ksp_atol 1e-6";
+ " -fieldsplit_velocity_pc_type asm -fieldsplit_velocity_pc_asm_overlap 1 "
+ " -fieldsplit_velocity_sub_pc_type lu "
+ " -fieldsplit_velocity_sub_pc_factor_mat_solver_type mumps";
+ " -fieldsplit_velocity_ksp_type gmres -fieldsplit_velocity_ksp_rtol 1e-1 "
+ " -fieldsplit_velocity_ksp_pc_side right"
+ " -fieldsplit_velocity_ksp_restart 50";
+ " -fieldsplit_pressure_pc_type jacobi -fieldsplit_pressure_ksp_type cg "
+ " -fieldsplit_pressure_ksp_max_it 5";
set(AOseen, sparams=paramsOseen, fields=fields[], names=names,
    schurPreconditioner = S, schurList = listX[]);
```

The pressure block term  $S$  is inverted by a CG method. The parameter `listX` refers to the matrix blocks with the “velocity” or “pressure” labels.

Note that we also use this preconditioner for the resolution of the fluid adjoint problem (6.2.10) which involves the transpose of the Oseen matrix. Note that some adaptation of this method would be needed if the objective function were to depend on the pressure, because the divergence  $\text{div}(\mathbf{w})$  of the fluid adjoint variable  $\mathbf{w}$  would not equal zero.

### Comparative performance of various stages of the parallel finite element implementation

We report in Table 6.1 various running times corresponding to the main operations performed during one iteration of our optimization algorithm. For simplicity, we considered a situation where only the fluid physics is involved and modeled by a Stokes system. More precisely, the following operations of our overall implementation are accelerated by parallel computing thanks to the use of domain decomposition:

1. all steps required for solving the state equations (2.2.1) to (2.2.3);
2. the computation of the values of the objective and constraint functions;
3. all steps required for computing the shape derivative of these functionals, including the resolution of adjoint systems;
4. the resolution of identification problems of the form (1.4.4) which feature the linear inversion of an elliptic problem.

We note, in Table 6.1, that the most computationally expensive tasks consisted in the resolution of the Stokes problem and its adjoint, as well as the various finite element matrix or vector assembly steps. All things considered, the scaling obtained is quite satisfactory and allowed us to run fluid shape optimization test cases in reasonable CPU time.

## 6.2 A FEW (MODERATELY) LARGE-SCALE THREE DIMENSIONAL MULTIPHYSICS APPLICATIONS

In this whole section, we go back to the three physics setting of chapter 2: a computational domain  $D = \Omega_s \cup \Omega_f \subset \mathbb{R}^3$  is given, which is the disjoint union of solid and fluid phases  $\Omega_s$  and  $\Omega_f$ . The behavior of the fluid-solid system is described by the weakly coupled system of partial differential equations (2.2.1) to (2.2.3), which determines the fluid velocity and pressure fields  $(\mathbf{v}, p)$  in  $\Omega_f$ , the temperature field  $T$  in  $D$  and the elastic displacement  $\mathbf{u}$  for  $\Omega_s$ . The ultimate goal is to optimize the shape of the interface  $\Gamma = \Omega_s \cap \Omega_f$  (a three dimensional surface) in order to solve constrained optimization problems of the form

$$\begin{aligned} \min_{\Gamma} \quad & J(\Gamma, \mathbf{v}(\Gamma), p(\Gamma), T(\Gamma), \mathbf{u}(\Gamma)) \\ \text{s.t.} \quad & \begin{cases} g_i(\Gamma, \mathbf{v}(\Gamma), p(\Gamma), T(\Gamma), \mathbf{u}(\Gamma)) = 0, & 1 \leq i \leq p, \\ h_j(\Gamma, \mathbf{v}(\Gamma), p(\Gamma), T(\Gamma), \mathbf{u}(\Gamma)) \leq 0, & 1 \leq j \leq q, \end{cases} \end{aligned} \quad (6.2.1)$$

where the notation convention is that introduced in chapter 2. In what follows, we treat four instances of the program (6.2.1). Our first three sets of examples involve only one physics at a time. In section 6.2.1, we examine very classical problems in linear elasticity, by considering the optimization of mechanical

Action	1 cpu	4 cpus	8 cpus
Loaded meshes	0.48	0.48	0.49
Partitioning	0.05	7.0	6.5
Saved partitioning data	0.82	0.3	0.21
Build Stokes Matrix and Rhs	10.73	2.7	1.38
Solved Stokes problem	56.42	18.72	10.98
Built interpolate Fhflvp matrix	0.38	0.16	0.09
Computed objective J	6.96	2.12	1.17
Read partitioning data	1.53	0.71	0.51
Adjoint fluid matrix assembly	18.8	5.51	2.59
Connectivity operation for non optimizable subdomains	3.97	1.27	0.69
Riesz matrix built	4.14	1.31	0.71
Assemble adjoint fluid RHS	6.25	1.8	0.89
Adjoint fluid problem resolution	103.84	38.04	23.24
Assemble volumetric shape derivative	72.68	20.9	10.92
Identification of the shape derivative of J to a gradient	1.67	0.66	0.5
<b>Total FreeFEM running CPU time</b>	<b>308.14</b>	<b>112.53</b>	<b>71.55</b>

Table 6.1: Running cpu times (2.60 GHz) for all finite element operations processed at every optimization iteration. The test case considered is the drag minimization around an obstacle for a 3-d Stokes problem featuring approx 22,000 vertices.

structures subject to either traction or torsion loads. The originality of our work is the treatment of these problems with the mesh evolution technique of [24] at much higher resolutions than in the seminal paper. We also discuss qualitative differences between the use of the surface expression (2.4.14) of the shape derivative and that of the volume expression (2.4.13).

The next section 6.2.2 considers an optimal design problem in pure heat conduction, which is a 3-d extension of the 2-d case treated in section 2.5.7. Although the physics at play is the least complicated to solve among the four considered situations (it involves only a scalar elliptic problem), the obtained final design is very intricate and illustrates well the efficiency of our mesh evolution method and the remeshing library `mmg`.

The third context of interest is that of shape and topology optimization in fluid mechanics; we tackle in section 6.2.3 the very classical problem of finding optimal aerodynamic designs with respect to the induced lift and drag forces. This problem has been the object of much effort in the literature, see e.g. [191, 260]; however, these contributions most often consider industrial contexts featuring very high Reynolds numbers and where the shape to optimize is parameterized by a small number of parameters (which make sense because very small design update can lead to a substantial increase of performance). Often, automatic differentiation is used rather than Hadamard's shape derivatives in order to obtain the sensitivity to these parameters in the context of the resolution of the physics with industrial codes. The novelty of our work is the application of our *topology optimization* method, relying on analytic Hadamard's shape derivatives, which allows to compute optimal aerodynamic designs (at low Reynolds number) without resorting to any parameterization of the shape.

Finally, our fourth and last test case of section 6.2.4 features two weakly coupled physics: a vertical plate is pushed down by a fluid; the problem at hand is to find a distribution of solid material around the plate in order to make the whole structure the least compliant as possible. This is the most computationally involved test case considered in this thesis: discretization meshes at play contained up to 250,000 vertices in the fluid domain (and the same in the solid domain), which required the resolution of linear systems featuring more than  $2 \times 10^6$  degrees of freedom.

All these results can still be considered as preliminary: future work will seek to find optimal designs for convective heat transfer problems where both fluid mechanics and thermal conduction interact. Many

improvements could be considered in order to reach much larger, industrial size problems: one of the limiting factor which prevented us to obtain such results lies in the fact that many important steps of the optimization algorithm are still *sequential*, including for instance the remeshing step or the computation of the signed distance function.

### 6.2.1 Cantilever beam subject to traction or torsion loads

We start by reproducing the classical benchmark test case of a 3-d cantilever beam subject to either a flexural or a torsion load. The computational domain  $D$  is a box of dimensions  $2 \times 1 \times 1$ . The solid structure is fixed at four squares of size  $0.3 \times 0.3$  located on the left-hand side of the boundary  $\partial D$  as depicted on [Figure 6.3](#). A force  $\mathbf{g}$  is applied on a disk-shaped region at the center of the right-hand side of  $\partial D$ ; two cases are considered:

- *traction load*: we set  $\mathbf{g} := -\mathbf{e}_y$  where  $(\mathbf{e}_x, \mathbf{e}_y, \mathbf{e}_z)$  denotes the canonical basis of  $\mathbb{R}^3$ , i.e. the load  $\mathbf{g}$  is vertical, pointing downward;
- *torsion load*: we set

$$\mathbf{g} = \begin{pmatrix} 0 \\ \frac{(z-0.5)}{\sqrt{(y-0.5)^2+(z-0.5)^2}} \\ -\frac{(y-0.5)}{\sqrt{(y-0.5)^2+(z-0.5)^2}} \end{pmatrix}$$

which corresponds to a torsion force field.

The goal is to minimize the compliance of the structure  $\Omega_s$  under a volume constraint:

$$\begin{aligned} \min \quad & J(\Omega_s, \mathbf{u}(\Omega_s)) := \int_{\Omega_s} A\mathbf{e}(\mathbf{u}) : \mathbf{e}(\mathbf{u}) dx \\ \text{s.t.} \quad & \text{Vol}(\Omega_s) := \int_{\Omega_s} dx = V_{\text{target}}. \end{aligned} \tag{6.2.2}$$

where  $V_{\text{target}}$  is a target volume, set to 0.15 in this example.

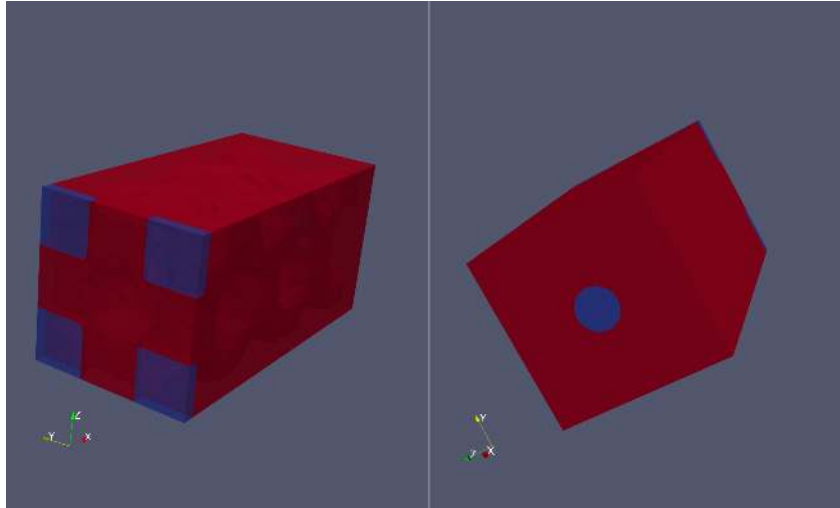


Figure 6.3: Computational domain for the 3-d cantilever test case of [section 6.2.1](#) subjected to a bending load. The blue regions correspond to four square-shaped Dirichlet boundaries to which the whole structure is fixed, and a disk-shaped Neumann boundary to which either a traction or a torsion force field is applied.

The former situation, involving traction load has already been treated with the level set mesh evolution algorithm in [\[24\]](#), with a much smaller resolution however (18,081 vertices in [\[24\]](#) vs. 108,605 in our case for the first mesh).

Note that for the four numerical examples presented below, convergence was not fully attained as the objective function was still decreasing at the end of the performed iterations. However it is expected that more iterations would not substantially change the physical outline of the shape.

**Traction test case** Two numerical results are shown for the traction test case on Figs. 6.4 and 6.5 corresponding to the use of respectively the surface and volume expressions of the shape derivative. The optimization paths are quite different. Qualitatively, the use of the surface expression of the shape derivative (identified with a gradient via the inner product of chapter 1, (1.4.8)) seems to favor the creation of holes while less topological changes occur when the volume expression is used (identified with a gradient via the inner product of chapter 1, (1.4.7)) which favors the occurrence of walls rather than bars on this example. These plots further illustrate the use of remeshing: the level set mode of `mmg` (option `-1s`) makes it possible to obtain high quality meshes of the shape at every optimization iteration.

The convergence histories for the objective and constraint functions in both cases are provided in Figure 6.6. The final designs obtained by using either the surface or volume expression of the shape derivative show very similar performance.

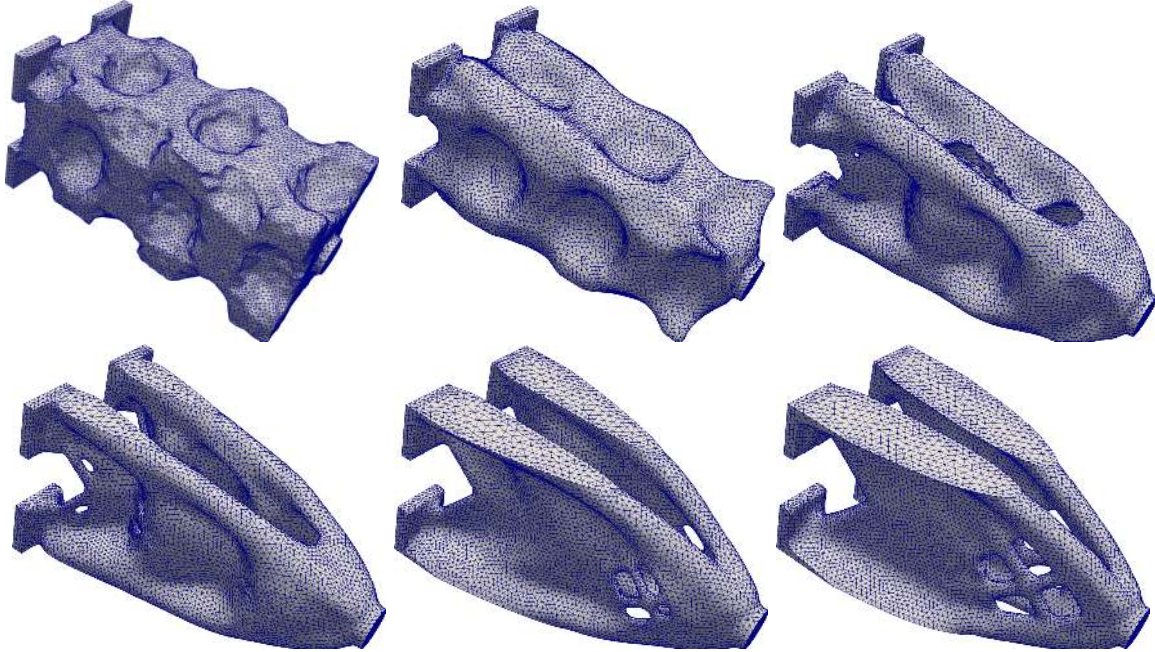


Figure 6.4: From top left to bottom right: iterations 0, 10, 20, 40, 100 and 400 of the optimization of a 3-d cantilever under *traction* load (section 6.2.1) using the *surface* expression of the shape derivative.

**Torsion test case** Optimization results for the *torsion* test case are shown on Figs. 6.7 and 6.8 corresponding to the use of either the surface or volume expressions of the shape derivative. The convergence histories for the objective and constraint functions are plotted on Figure 6.9. For this test case, a clear difference is observed between both computations the use of the volume expression favors a shell-shaped structure without bars at the upper end at the end of the iterations, which seems to be more efficient at least for the first iterations (note again that convergence is not fully attained).

### 6.2.2 Optimal design for pure thermal heat conduction

We now focus on the optimization of a pure heat conduction test case (only the thermal equation of chapter 2, (2.2.2) is solved). The problem considered is the direct extension to 3-d of the test case previously treated in chapter 2, section 2.5.3. The setting is represented on Figure 6.10: the hold-all domain  $D$  is a box with size  $1 \times 1 \times 1$ . It is divided between two phases  $\Omega_s$  with (low) conductivity  $k_s = 1$  and  $\Omega_f$  with (high) conductivity  $k_f = 100$ . A Dirichlet boundary condition is imposed on a small square of size  $0.4 \times 0.4$  at the bottom face of  $\partial D$  where the temperature is prescribed to  $T = 0$ . All other external boundaries of the cube  $D$  are adiabatic ( $\partial T / \partial \mathbf{n} = 0$ ). The whole domain is heated with a source  $Q_s = Q_f = 10^4$  and the goal is to find the shape of the interface  $\Gamma = \partial\Omega_s \cap \partial\Omega_f$  of the two materials which minimizes the average temperature over  $D$  subject to a volume constraint:

$$\begin{aligned} \min_{\Gamma} \quad & J(\Gamma, T(\Gamma)) = \int_D T dx \\ \text{s.t.} \quad & \text{Vol}(\Omega_f) = V_{\text{target}}. \end{aligned} \tag{6.2.3}$$

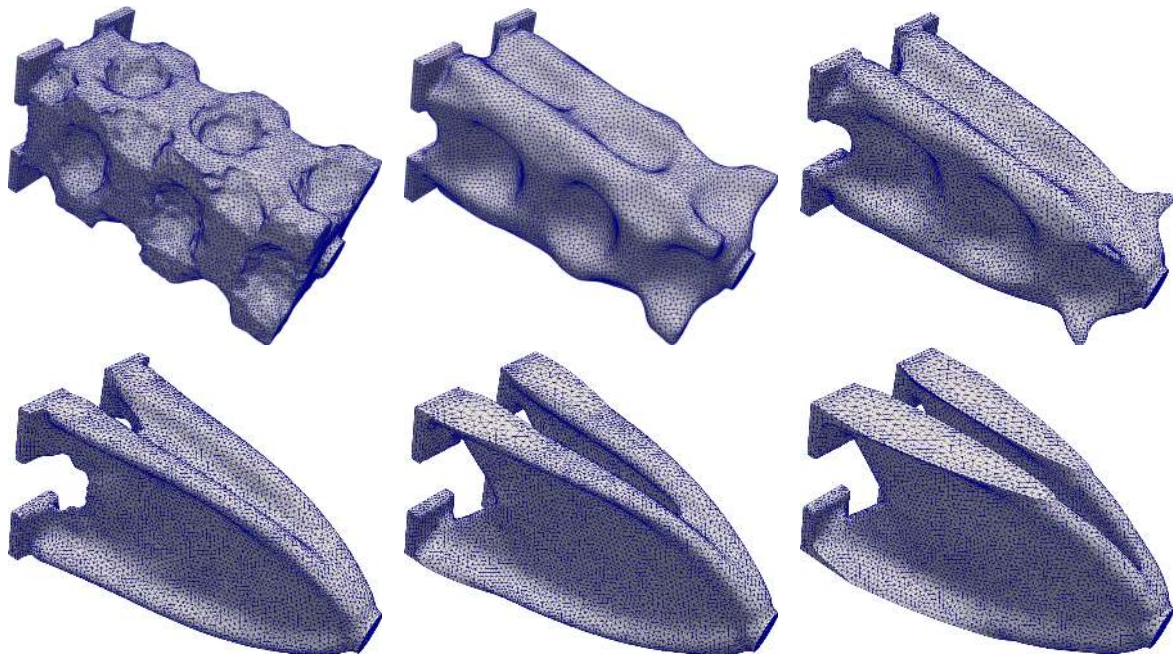


Figure 6.5: From top left to bottom right: iterations 0, 10, 20, 40, 100 and 280 of the optimization of a 3-d cantilever subjected to a *traction* load (section 6.2.1) using the *volume* expression of the shape derivative. Note the angle view is not the same as in Figure 6.4 for better visualization.

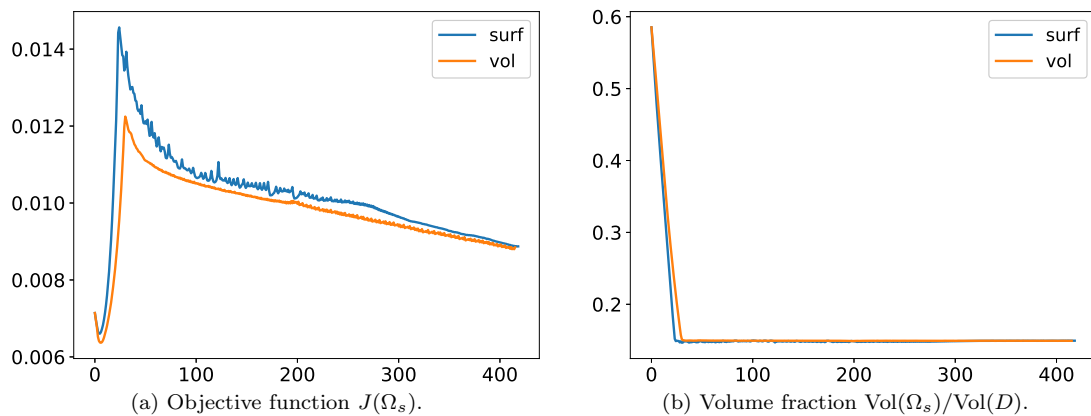


Figure 6.6: Convergence histories for the 3-d cantilever test cases subjected to a *traction* load of section 6.2.1 using either the *surface* or *volume* expression of the shape derivative.

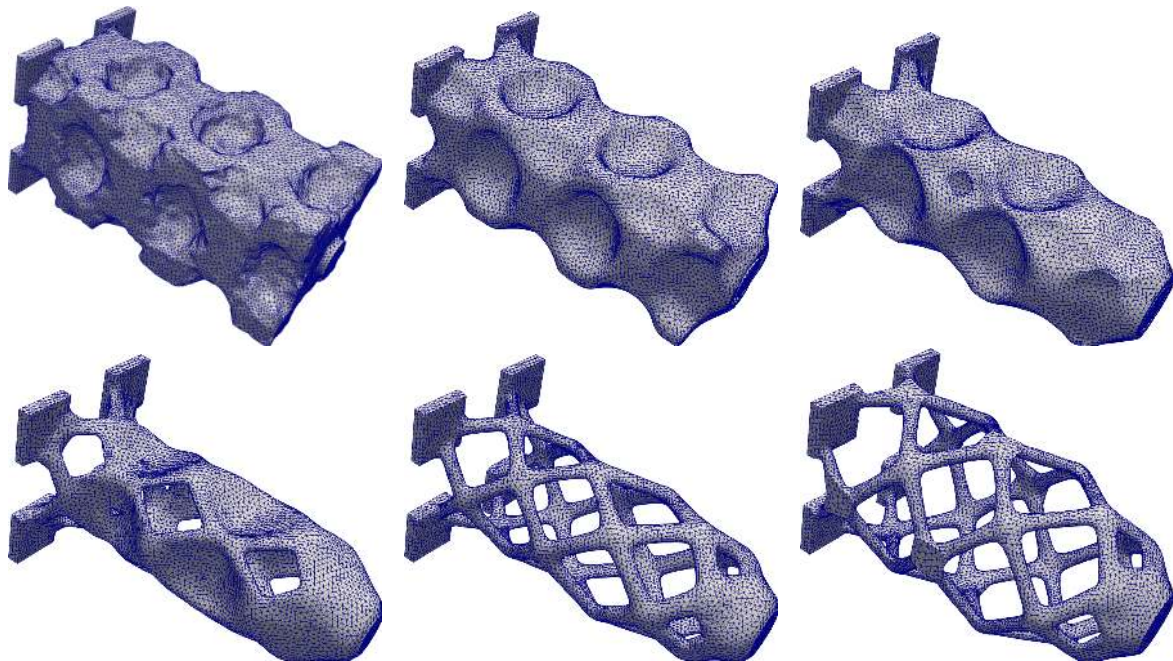


Figure 6.7: From top left to bottom right: iterations 0, 10, 20, 40, 100 and 400 of the optimization of a 3-d cantilever subjected to a *torsion* load (section 6.2.1) using the *surface* expression of the shape derivative.

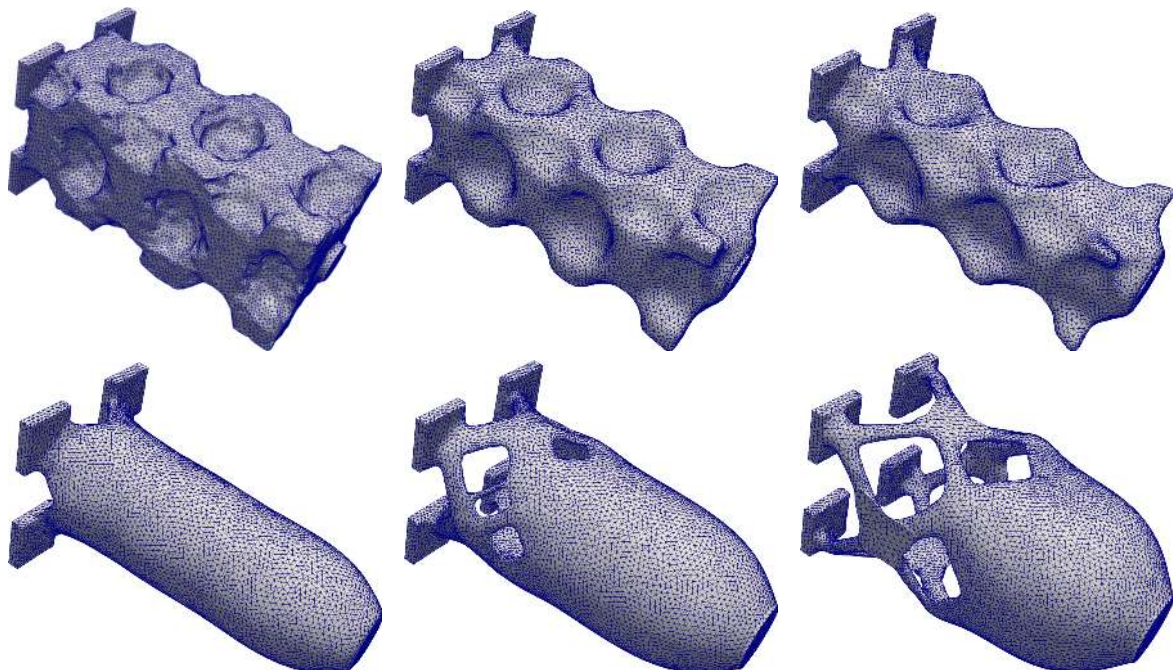


Figure 6.8: From top left to bottom right: iterations 0, 10, 20, 40, 100 and 350 of the optimization of a 3-d cantilever subjected to a *torsion* load (section 6.2.1) using the *surface* expression of the shape derivative.



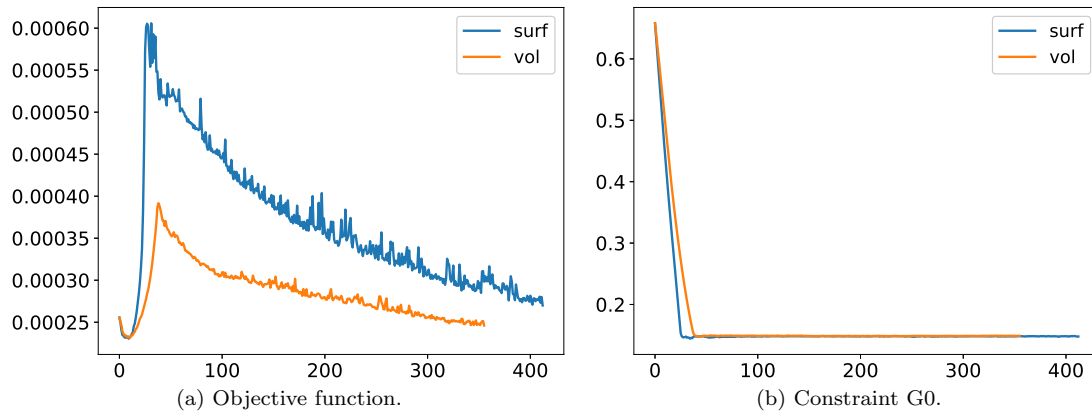


Figure 6.9: Convergence histories for the 3-d cantilever test case (subjected to a *torsion* load) of [section 6.2.1](#) using either the *surface* or *volume* expression of the shape derivative.

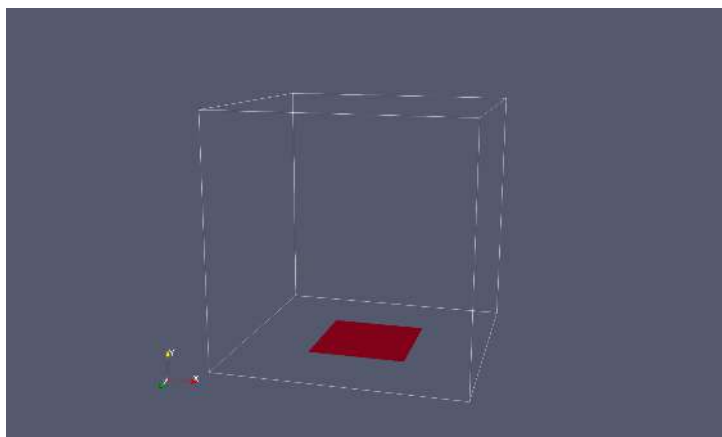


Figure 6.10: Setting for the pure conduction test case.

In this three-dimensional context, the volume constraint is set to  $V_{target} := 0.05$ .

Results are reported on Figs. 6.11 and 6.12, relying on the *volume* expression for the shape derivative. The resolution of the mesh varies from 63,761 vertices for the initial design to 206,464 for the final design. Remarkably, the remeshing software *mmg* and our optimization method are able to capture sheet like structures not thicker than one or two mesh element size. For this numerical example, the difference with the use of the surface expression of the shape derivative is very striking (we report a design obtained with the surface expression on Figure 6.13): the optimized shape presents much less details in the latter case and achieves a worse performance, as illustrated on Figure 6.14. Naturally, a lot of other factors come into play which could also account for the differences (such as the level of regularization of the shape derivative), for which we do not have a definitive explanation.

### 6.2.3 Lift–Drag topology optimization for aerodynamic design

This section investigates 3-d shape and topology optimization for lift-drag problems in aerodynamic design. The lift functional is the vertical force generated by a flow around an obstacle; it is commonly defined as a *surface integral* involving the normal stress tensor. The drag is the energy dissipated by the fluid around the obstacle. Lift-Drag shape optimization is a very classical problem which has been the object of a very large amount of contributions, see e.g. [259, 236, 191, 192, 142, 211, 175]. However, these references have considered situations very close to realistic applications where

- (i) the physics is more challenging than in our case, featuring for instance compressible fluids or characterized by much larger Reynolds numbers;
- (ii) the shape design is usually described by means of CAD parameters to optimize;
- (iii) very small updates of the shape may lead to substantial gains of efficiency. In this context, it makes sense to seek for improved geometries by means of very small deformations of the proposed CAD design.

Few works have actually tried to apply shape and topology optimization techniques, were the design shape is allowed to deform freely, to lift-drag problems: we are essentially aware of [100, 205, 164, 307]. In what follows, we treat a lift-drag optimization problem with the method of Hadamard and our topology optimization framework on 2-d and 3-d examples featuring a very small Reynolds number ( $Re = 200$ ). We are not aware, to the best of our knowledge, of analogous results in the 3-d setting.

A first part of this section is devoted to the computation of the shape derivative of this functional: although it has already been considered in several works of the literature, the calculation and the numerical implementation of the resulting formulas do not seem completely standard to us; we propose a special treatment based on a classical idea of [63, 137].

In a second part, we present 2-d and 3-d numerical designs for an instance of the lift-drag optimization problem.

#### Shape derivatives of the lift functional

Let  $D = \mathbb{R}^d$  be a computational domain featuring a liquid phase  $\Omega_f$  flowing around a solid obstacle  $\Omega_s \subset\subset D$ . The notation convention is again that assumed in chapter 2: the flow is entering the domain from a Dirichlet boundary  $\partial\Omega_f^D$  with a given velocity  $\mathbf{v} = \mathbf{v}_0$  and exits the domain with a zero normal stress boundary condition  $\sigma_f(\mathbf{v}, p) \cdot \mathbf{n} = 0$  on  $\partial\Omega_f^N$ . The remaining part of the fluid boundary is the solid interface  $\Gamma = \partial\Omega_f \setminus (\partial\Omega_f^D \cup \partial\Omega_f^N)$  which is to be optimized.

The *lift* generated by the obstacle is the total force exerted by the fluid on the solid interface  $\Gamma = \Omega_f \cap \Omega_s = \partial\Omega_s$  in the vertical,  $y$ -direction:

$$\text{Lift}(\Gamma, \mathbf{v}(\Gamma), p(\Gamma)) := - \int_{\Gamma} \mathbf{e}_y \cdot \sigma_f(\mathbf{v}, p) \cdot \mathbf{n} ds, \tag{6.2.4}$$

where we recall that the fluid stress tensor is given by  $\sigma_f(\mathbf{v}, p) = 2\nu e(\mathbf{v}) - pI$ , and the notation  $\mathbf{a} \cdot M \cdot \mathbf{b} := \mathbf{a}^T M \mathbf{b}$  for any vectors  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$  and matrix  $M \in \mathbb{R}^{d \times d}$ . Note that the minus sign accounts for our convention of chapter 2 whereby the normal  $\mathbf{n}$  is pointing outward the fluid domain. Several authors have considered the optimization of the lift functional, e.g. in [292], [214] based on surrogate models, [206] based on control points, [205] with the SIMP method and a different objective functional, and [262, 246] in a time varying setting, [219]. Several difficulties are involved in the evaluation of the shape derivative of the lift functional:

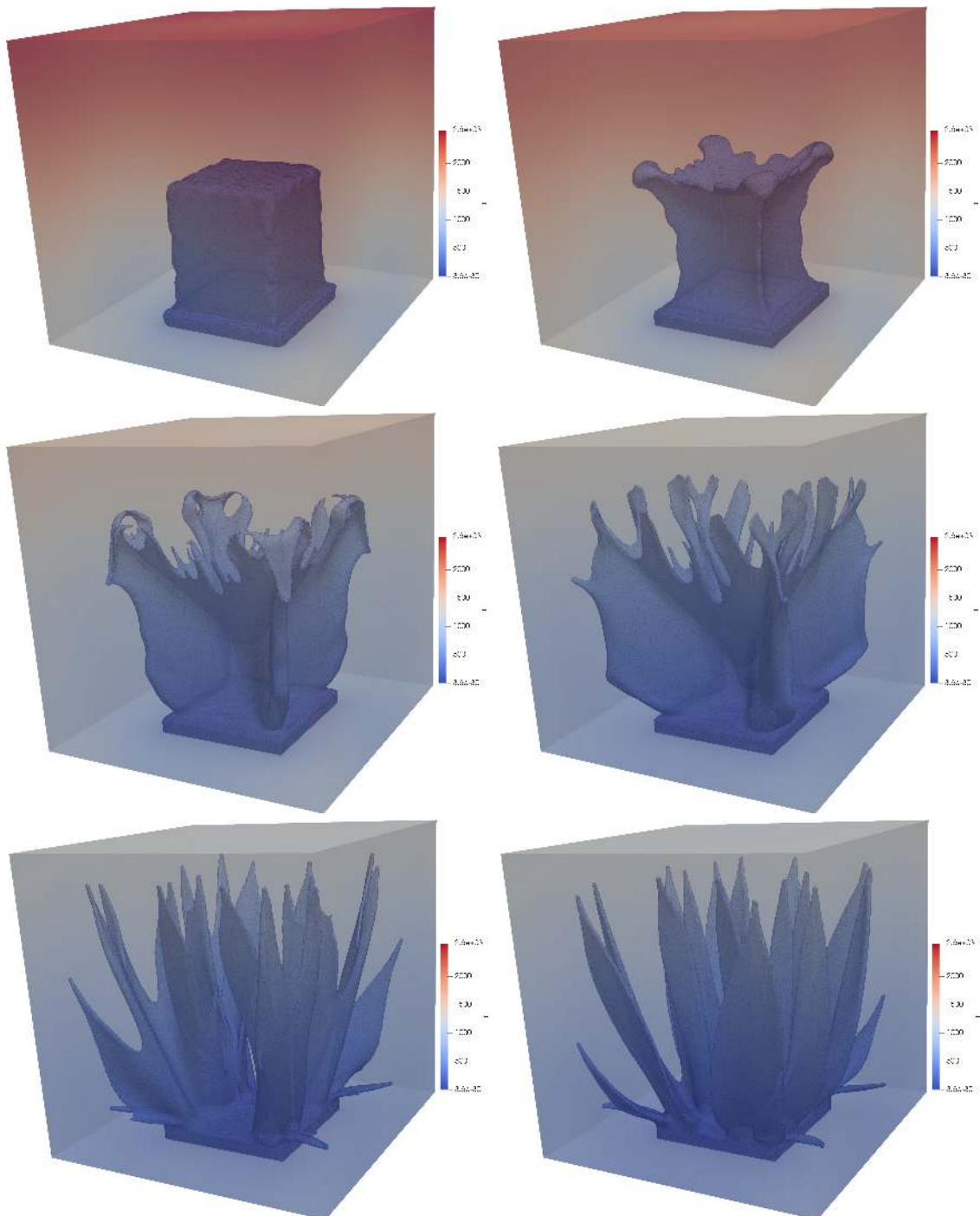


Figure 6.11: From top left to bottom right: iterations 0, 5, 15, 30, 100 and 258 of the optimization of a 3-d design for heat conduction using the *volume* expression of the shape derivative (test case of [section 6.2.2](#)).

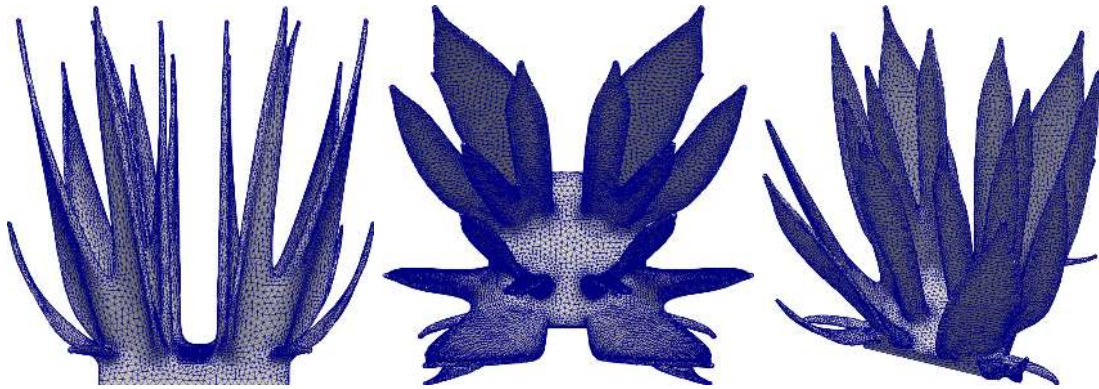


Figure 6.12: Different 3-d views of the optimized design for the heat conduction test case of section 6.2.2.

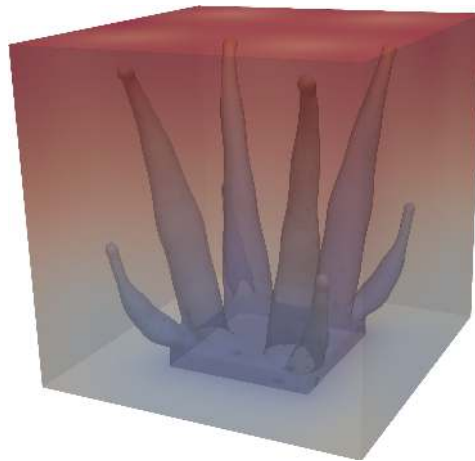


Figure 6.13: Optimized 3-d design for heat conduction obtained with the *surface* expression of the shape derivative (test case of section 6.2.2).

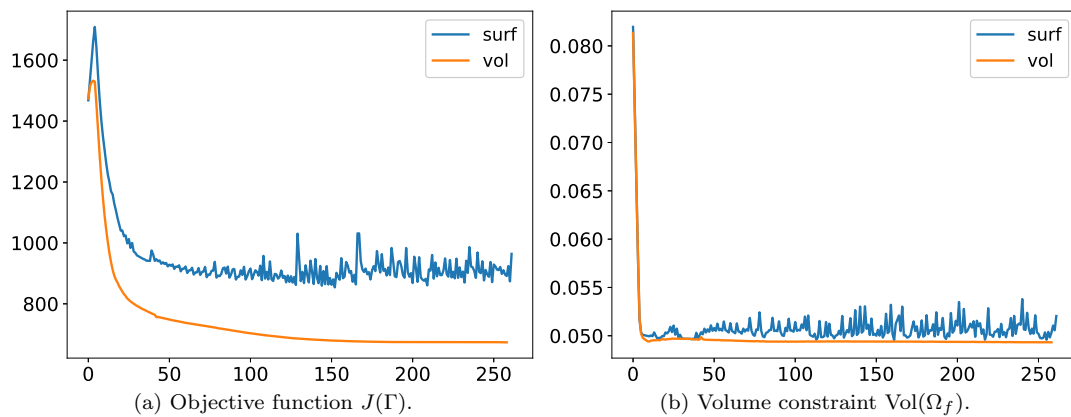


Figure 6.14: Convergence histories for the 3-d heat conduction test case of section 6.2.2 using either the volume or the surface expression of the shape derivative.

1. from the theoretical point of view, the partial derivative  $(\mathbf{v}, p) \mapsto \partial \text{Lift} / \partial (\mathbf{v}, p)$  is not a continuous linear form of  $H^1(\Omega_f) \times L^2(\Omega_f)$ : the fact that the normal derivative  $\sigma_f(\mathbf{v}, p) \cdot \mathbf{n}$  has a trace on  $\Gamma$  (as an element of  $H^{-1/2}(\Omega_f, \mathbb{R}^d)$ ) is related to the fact that  $(\mathbf{v}, p)$  solves the Navier-Stokes system (see [chapter 2, section 2.4](#)) but it may not be true for arbitrary variations  $(\mathbf{w}, q) \in H^1(\Omega_f) \times L^2(\Omega_f)$ ;
2. even assuming a different variational setting where these partial derivatives would make sense (see [\[296\]](#) for an example of such setting), the functional is posed on the boundary. Introducing as in [chapter 2](#) the transported function

$$\mathfrak{Lift}(\boldsymbol{\theta}, \widehat{\mathbf{v}}, \widehat{p}) := \text{Lift}(\Gamma_{\boldsymbol{\theta}}, \widehat{\mathbf{v}} \circ (I + \boldsymbol{\theta})^{-1}, \widehat{p} \circ (I + \boldsymbol{\theta})^{-1}), \quad (6.2.5)$$

the computation of the partial derivative

$$\frac{\partial \mathfrak{Lift}}{\partial \boldsymbol{\theta}}$$

is not straightforward since (i) it involves the gradient of  $v$ , and (ii) the variations of the normal  $\mathbf{n}$  come into play. This approach has been followed by [\[292\]](#) for the compressible Navier-Stokes equations.

Here, we propose an alternative method, which relies on a reformulation of the lift functional [\(6.2.4\)](#) as a volume integral: let  $\mathcal{X} \in H^1(\Omega_f)$  be any extension of the constant function 1 on  $\Gamma$  (i.e.  $\mathcal{X} = 1$  on  $\Gamma$ ) and vanishing on the complementary boundary (i.e.  $\mathcal{X} = 0$  on  $\partial\Omega_f \setminus \Gamma$ ). In our particular implementation,  $\mathcal{X}$  is obtained as the solution to the following Poisson problem:

$$\begin{cases} -\Delta \mathcal{X} = 0 & \text{in } \Omega_f \\ \mathcal{X} = 1 & \text{on } \Gamma \\ \mathcal{X} = 0 & \text{on } \partial\Omega_f \setminus \Gamma \end{cases} \quad (6.2.6)$$

The function  $\mathcal{X} \mathbf{e}_y$  yields then an extension of the vector field  $\mathbf{e}_y$  vanishing on  $\partial\Omega_f \setminus \Gamma$ , which enables to rewrite  $\text{Lift}(\Gamma)$  as a volume integral:

$$\begin{aligned} \text{Lift}(\Gamma) &= - \int_{\Gamma} \mathcal{X} \mathbf{e}_y \cdot \sigma_f(\mathbf{v}, p) \cdot \mathbf{n} ds \\ &= - \int_{\Omega_f} \text{div}(\mathcal{X} \sigma_f(\mathbf{v}, p) \cdot \mathbf{e}_y) dx \\ &= - \int_{\Omega_f} (\nabla \mathcal{X} \cdot \sigma_f(\mathbf{v}, p) \cdot \mathbf{e}_y + \mathcal{X} \text{div}(\sigma_f(\mathbf{v}, p)) \cdot \mathbf{e}_y) dx \\ &= \int_{\Omega_f} (\mathcal{X} \mathbf{f}_f \cdot \mathbf{e}_y - \rho \mathcal{X} \mathbf{e}_y \cdot \nabla \mathbf{v} \cdot \mathbf{v} - \nabla \mathcal{X} \cdot \sigma_f(\mathbf{v}, p) \cdot \mathbf{e}_y) dx \end{aligned} \quad (6.2.7)$$

where we have used the state equation  $-\text{div}(\sigma_f(\mathbf{v}, p)) + \rho \nabla \mathbf{v} \mathbf{v} = \mathbf{f}_f$  to obtain the last line. This rewriting is rather classical and is often used in numerical applications, since it is known to yield a more accurate evaluation of the lift functional [\[63, 137\]](#).

**Remark 6.2.** The last equality of [\(6.2.7\)](#) can be considered as a *definition* for the meaning of the normal derivative  $\sigma_f(\mathbf{v}, p) \cdot \mathbf{n}$  in [\(6.2.4\)](#), since it does not depend on the chosen extension  $\mathcal{X} \in H^1(\Omega_f)$  satisfying  $\mathcal{X} = 1$  on  $\Gamma$  and  $\mathcal{X} = 0$  on  $\partial\Omega_f \setminus \Gamma$  (from the variational formulation [\(2.4.2\)](#)).

Equation [\(6.2.7\)](#) can now be easily differentiated with respect to the shape. Let us introduce the functional spaces  $V_{\mathbf{v}, p}(\Gamma)$  and  $V_{\mathcal{X}}(\Gamma)$  associated with the solutions  $(\mathbf{v}, p)$  and  $\mathcal{X}$ :

$$V_{\mathbf{v}, p}(\Gamma) := \{(\mathbf{w}, q) \in H^1(\Omega_f, \mathbb{R}^d) \times L^2(\Omega_f) \mathbb{R} \mid \mathbf{w} = 0 \text{ on } \partial\Omega_f\},$$

$$V_{\mathcal{X}}(\Gamma) := \{\Psi \in H^1(\Omega_f) \mid \Psi = 0 \text{ on } \partial\Omega_f.\}$$

We also denote by  $\mathbf{v}_0 + V_{\mathbf{v}, p}(\Gamma)$  and  $1 + V_{\mathcal{X}}(\Gamma)$  the affine spaces

$$\mathbf{v}_0 + V_{\mathbf{v}, p}(\Gamma) := \{(\mathbf{v}, p) \in H^1(\Omega_f, \mathbb{R}^d) \times L^2(\Omega_f) \mathbb{R} \mid \mathbf{v} = 0 \text{ on } \partial\Omega_f \setminus \partial\Omega_f^D \text{ and } \mathbf{v} = \mathbf{v}_0 \text{ on } \partial\Omega_f^D\},$$

$$1 + V_{\mathcal{X}}(\Gamma) := \{\mathcal{X} \in H^1(\Omega_f) \mid \mathcal{X} = 0 \text{ on } \partial\Omega_f \setminus \Gamma \text{ and } \mathcal{X} = 1 \text{ on } \Gamma.\}$$

In what follows, we assume that the extension  $\mathcal{X}$  of the constant function 1 on  $\Gamma$  is given by the solution to the Poisson problem (6.2.6). We denote, for any  $\boldsymbol{\theta} \in W^{1,\infty}(D, \mathbb{R}^d)$ ,  $(\widehat{\mathbf{v}}, \widehat{p}) \in \mathbf{v}_0 + V_{\mathbf{v},p}(\Gamma)$  and  $\widehat{\mathcal{X}} \in 1 + V_{\mathcal{X}}$ ,

$$\mathfrak{Lift}(\boldsymbol{\theta}, \widehat{\mathbf{v}}, \widehat{p}, \widehat{\mathcal{X}}) := \text{Lift}(\Gamma_{\boldsymbol{\theta}}, \widehat{\mathbf{v}} \circ (I + \boldsymbol{\theta})^{-1}, \widehat{p} \circ (I + \boldsymbol{\theta})^{-1}, \widehat{\mathcal{X}} \circ (I + \boldsymbol{\theta})^{-1}) \quad (6.2.8)$$

the transported functional on the reference situation. Note that since  $\mathfrak{Lift}$  does not depend on the choice of the extension  $\mathcal{X}$  satisfying  $\mathcal{X} = 1$  on  $\Gamma$ , (6.2.8) and (6.2.5) coincide for such  $\mathcal{X}$ . Then, with the notation of chapter 2:

$$\begin{aligned} \forall \boldsymbol{\theta} \in W^{1,\infty}(D, \mathbb{R}^d), \quad \frac{\partial \mathfrak{Lift}}{\partial \boldsymbol{\theta}}(\boldsymbol{\theta}) &= \int_{\Omega_f} (\mathcal{X} \mathbf{f}_f \cdot \mathbf{e}_y - \rho \mathcal{X} \mathbf{e}_y \cdot \nabla \mathbf{v} \cdot \mathbf{v} - \nabla \mathcal{X} \cdot \sigma_f(\mathbf{v}, p) \cdot \mathbf{e}_y) \text{div}(\boldsymbol{\theta}) dx \\ &+ \int_{\Omega_f} (\rho \mathcal{X} \mathbf{e}_y \cdot \nabla \mathbf{v} \nabla \boldsymbol{\theta} \cdot \mathbf{v} + (\nabla \boldsymbol{\theta}^T \nabla \mathcal{X}) \cdot \sigma_f(\mathbf{v}, p) \cdot \mathbf{e}_y + \nu \nabla \mathcal{X} \cdot (\nabla \mathbf{v} \nabla \boldsymbol{\theta} + \nabla \boldsymbol{\theta}^T \nabla \mathbf{v}^T) \cdot \mathbf{e}_y) dx, \\ \forall \boldsymbol{\theta} \in W^{1,\infty}(D, \mathbb{R}^d), \quad \frac{\overline{\partial \mathfrak{Lift}}}{\partial \boldsymbol{\theta}}(\boldsymbol{\theta}) &= \int_{\Gamma} (\mathcal{X} \mathbf{f}_f \cdot \mathbf{e}_y - \rho \mathcal{X} \mathbf{e}_y \cdot \nabla \mathbf{v} \cdot \mathbf{v} - \nabla \mathcal{X} \cdot \sigma_f(\mathbf{v}, p) \cdot \mathbf{e}_y) \boldsymbol{\theta} \cdot \mathbf{n} ds \\ &+ \int_{\Gamma} ((\nabla \mathcal{X} \cdot \mathbf{n})(\mathbf{n} \cdot \sigma_f(\mathbf{v}, p) \cdot \mathbf{e}_y) + \nu (\nabla \mathcal{X} \cdot \nabla \mathbf{v} \cdot \mathbf{n})(\mathbf{e}_y \cdot \mathbf{n}) + \nu (\mathbf{e}_y \cdot \nabla \mathbf{v} \cdot \mathbf{n})(\nabla \mathcal{X} \cdot \mathbf{n})) \boldsymbol{\theta} \cdot \mathbf{n} ds \quad (6.2.9) \\ &= \int_{\Gamma} [\mathbf{f}_f \cdot \mathbf{e}_y + 2\nu e(\mathbf{v}) : (\nabla \mathcal{X} \otimes \mathbf{e}_y)] \boldsymbol{\theta} \cdot \mathbf{n} ds, \end{aligned}$$

$$\forall (\mathbf{w}', q') \in V_{\mathbf{v},p}(\Gamma) \quad \frac{\partial \mathfrak{Lift}}{\partial (\mathbf{v}, p)}(\mathbf{w}', q') = - \int_{\Omega_f} (\rho \mathcal{X} \mathbf{e}_y \cdot \nabla \mathbf{w}' \cdot \mathbf{v} + \rho \mathcal{X} \mathbf{e}_y \cdot \nabla \mathbf{v} \cdot \mathbf{w}' + \nabla \mathcal{X} \cdot \sigma_f(\mathbf{w}', q') \cdot \mathbf{e}_y) dx,$$

The simplifications in the second equality of (6.2.9) are classical consequences of the tangential gradients  $\nabla_{\Gamma} \mathbf{v} = 0$  and  $\nabla_{\Gamma} \mathcal{X}$  being 0 on  $\Gamma$ .

Finally, the partial derivative with respect to  $\mathcal{X}$  is zero, because an easy integration by part implies that

$$\forall \Psi' \in V_{\mathcal{X}}(\Gamma), \quad \frac{\partial \mathfrak{Lift}}{\partial \mathcal{X}}(\Psi') = - \int_{\Gamma} \Psi' \mathbf{e}_y \cdot \sigma_f(\mathbf{v}, p) \cdot \mathbf{n} ds = 0.$$

Therefore, there is no contribution of the variations of the extension  $\mathcal{X}$  defined in (6.2.6) to the shape derivative of  $\mathfrak{Lift}$ .

All in all, applying the result of chapter 2, propositions 2.3 and 2.4 yields expressions for the shape derivative of  $\text{Lift}$  (eqn. (6.2.4)) in surface and volume forms:

**Proposition 6.1.** *Let  $(\mathbf{w}, q) \in V_{\mathbf{v},p}(\Gamma)$  the adjoint fluid variables solutions to*

$$\begin{aligned} \forall (\mathbf{w}', q') \in V_{\mathbf{v},p}(\Gamma), \\ \int_{\Omega_f} \left( \sigma_f(\mathbf{w}, q) : \nabla \mathbf{w}' + \rho \mathbf{w} \cdot \nabla \mathbf{w}' \cdot \mathbf{v} + \rho \mathbf{w} \cdot \nabla \mathbf{v} \cdot \mathbf{w}' - q' \text{div}(\mathbf{w}') \right) dx = \frac{\partial \mathfrak{Lift}}{\partial (\mathbf{v}, p)}(\mathbf{w}', q'). \quad (6.2.10) \end{aligned}$$

The lift functional  $\Gamma \mapsto \text{Lift}(\Gamma, \mathbf{v}(\Gamma), p(\Gamma))$  is differentiable with respect to the shape and the shape derivative reads (in volumetric form):

$$\begin{aligned} \forall \boldsymbol{\theta} \in W^{1,\infty}(D, \mathbb{R}^d), \quad \frac{d}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=0} \left[ \text{Lift}(\Gamma_{\boldsymbol{\theta}}, \mathbf{v}(\Gamma_{\boldsymbol{\theta}}), p(\Gamma_{\boldsymbol{\theta}}), T(\Gamma_{\boldsymbol{\theta}}), \mathbf{u}(\Gamma_{\boldsymbol{\theta}})) \right] (\boldsymbol{\theta}) \\ = \frac{\partial \mathfrak{Lift}}{\partial \boldsymbol{\theta}}(\boldsymbol{\theta}) + \int_{\Omega_f} [\mathbf{w} \cdot \text{div}(\mathbf{f}_f \otimes \boldsymbol{\theta}) - (\sigma_f(\mathbf{v}, p) : \nabla \mathbf{w} + \rho \mathbf{w} \cdot \nabla \mathbf{v} \cdot \mathbf{v}) \text{div}(\boldsymbol{\theta})] dx \\ + \int_{\Omega_f} [\sigma_f(\mathbf{v}, p) : (\nabla \mathbf{w} \nabla \boldsymbol{\theta}) + \sigma_f(\mathbf{w}, q) : (\nabla \mathbf{v} \nabla \boldsymbol{\theta}) + \rho \mathbf{w} \cdot (\nabla \mathbf{v} \nabla \boldsymbol{\theta}) \cdot \mathbf{v}] dx \quad (6.2.11) \end{aligned}$$

or (in surface form):

$$\frac{d}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=0} \left[ \text{Lift}(\Gamma_{\boldsymbol{\theta}}, \mathbf{v}(\Gamma_{\boldsymbol{\theta}}), p(\Gamma_{\boldsymbol{\theta}}), T(\Gamma_{\boldsymbol{\theta}}), \mathbf{u}(\Gamma_{\boldsymbol{\theta}})) \right] (\boldsymbol{\theta}) = \frac{\overline{\partial \mathfrak{Lift}}}{\partial \boldsymbol{\theta}}(\boldsymbol{\theta}) + \int_{\Gamma} (\mathbf{f}_f \cdot \mathbf{w} + \sigma_f(\mathbf{w}, q) : \nabla \mathbf{v}) ds. \quad (6.2.12)$$

The above formulas are easily implemented in our shape and topology optimization framework; they are at the basis of the resolution of the test cases presented in the next paragraph.

### Numerical test cases

We consider the problem of maximizing the Lift generated by a flow obstacle  $\Omega_s \subset\subset D$  subject to an upper bound constraint on the drag (so that the obstacle remains aerodynamic). In addition, the volume occupied by the obstacle and the location of its center of mass are prescribed:

$$\begin{aligned} \min \quad & -\text{Lift}(\Gamma, \mathbf{v}(\Gamma), p(\Gamma)) \\ \text{s.t.} \quad & \begin{cases} \text{Drag}(\Gamma, \mathbf{v}(\Gamma), p(\Gamma)) \leq \text{DRAG}_0 \\ \text{Vol}(\Omega_f) = V_0 \\ \mathbf{X}(\Omega_s) := \frac{1}{|\Omega_s|} \int_{\Omega_s} \mathbf{x} dx = \mathbf{x}_0. \end{cases} \end{aligned} \quad (6.2.13)$$

The drag functional has been considered in [chapter 2, section 2.5.4](#); it is defined by

$$\text{Drag}(\Gamma, \mathbf{v}(\Gamma), p(\Gamma)) := \int_{\Omega_f} \sigma_f(\mathbf{v}, p) : \nabla \mathbf{v} dx = \int_{\Omega_f} 2\nu e(\mathbf{v}) : e(\mathbf{v}) dx.$$

The constant  $\text{DRAG}_0$  is set to  $\alpha \text{DRAG}^*$  ( $\alpha = 1.1$  in 2-d and  $\alpha = 1.5$  in 3-d) where  $\text{DRAG}^*$  is the value of the minimum drag problem subject to the same volume and position constraints.

We start by solving the problem in 2-d; the setting is identical to that of [section 2.5.4](#):  $D = [0, 1] \times [0, 1]$ ,  $V_{\text{target}} = 0.03$ ,  $x_0 = (0.5, 0.5)$ ,  $\|\mathbf{v}_0\| = 1$ ,  $\rho = 1$ ,  $\nu = 1/200$  and  $\text{Re} = 200$ . The flow is entering the left-hand boundary with velocity  $\mathbf{v}_0 = \mathbf{e}_x$  and it exits the right-hand boundary with zero normal stress ( $\sigma_f(\mathbf{v}, p) \cdot \mathbf{n} = 0$ ). A slip boundary condition  $\mathbf{v} \cdot \mathbf{n} = 0$  is prescribed at other boundaries. The minimum drag value computed for this example is  $\text{DRAG}^* \simeq 0.11$ , so that the upper bound constraint is given by  $\text{DRAG}_0 = \text{DRAG}^* \times 1.1 = 0.121$ . The initial and optimized shapes with the corresponding velocity fields are shown on respectively [Figure 6.15](#). A few intermediate iterations are shown on [Figure 6.16](#), and the optimization histories are reported on [Figure 6.18](#). Finally, the mesh discretizing the final shape is plotted on [Figure 6.17](#); both the boundary of the obstacle and the outlet were refined in order to numerically capture boundary layers and vortex patterns.

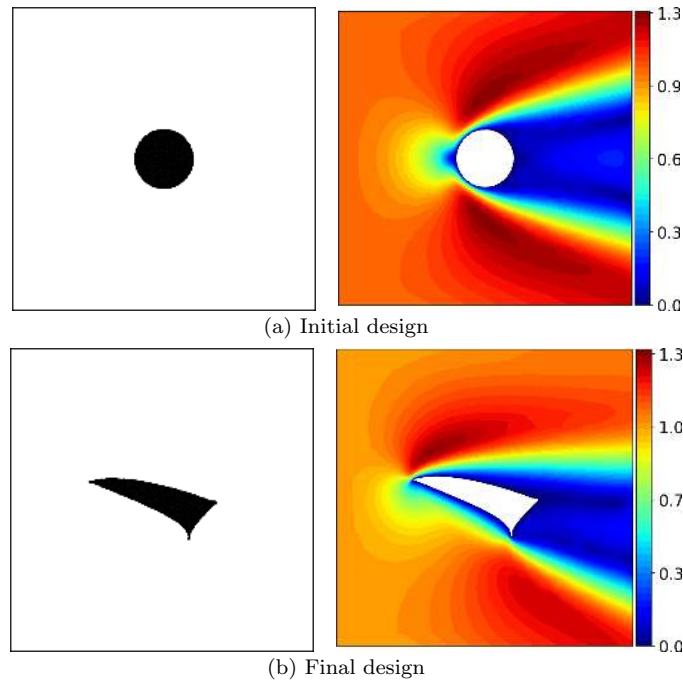


Figure 6.15: Optimization results for the 2-d lift-drag optimization problem of [section 6.2.3](#). The norm of the velocity field  $\mathbf{v}$  is plotted on the right.

Our design is slightly similar to those obtained in the work [\[205\]](#) with a completely different method: a density based approach was used and a different expression was used to estimate the lift functional.

We then solve the same problem in 3-d: the hold-all domain is the box  $D = [0, 1] \times [0, 1] \times [0, 1]$ . A flow is entering with a velocity  $\mathbf{v}_0 = \mathbf{e}_x$  on the left-hand side of the domain. A slip boundary condition

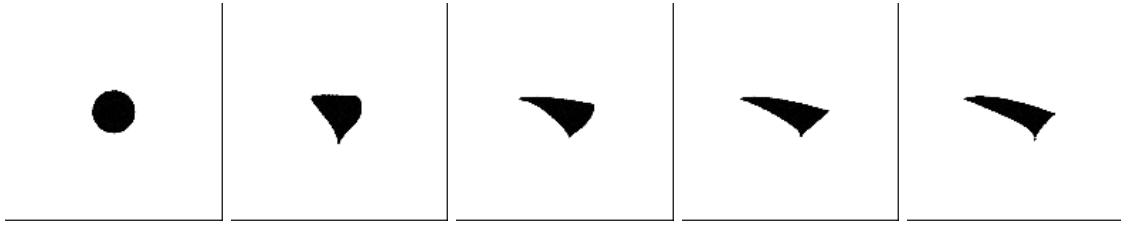


Figure 6.16: Intermediate iterations 0, 8, 30, 80 and 200 for the 2-d lift-drag optimization problem of section 6.2.3

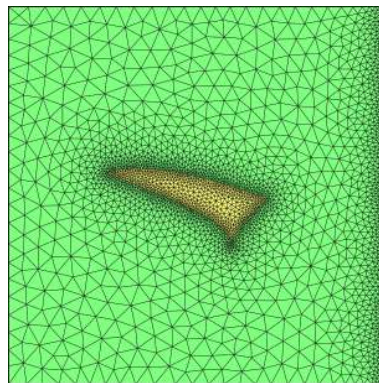


Figure 6.17: Mesh of the final shape for the 2-d lift-drag optimization problem of section 6.2.3.

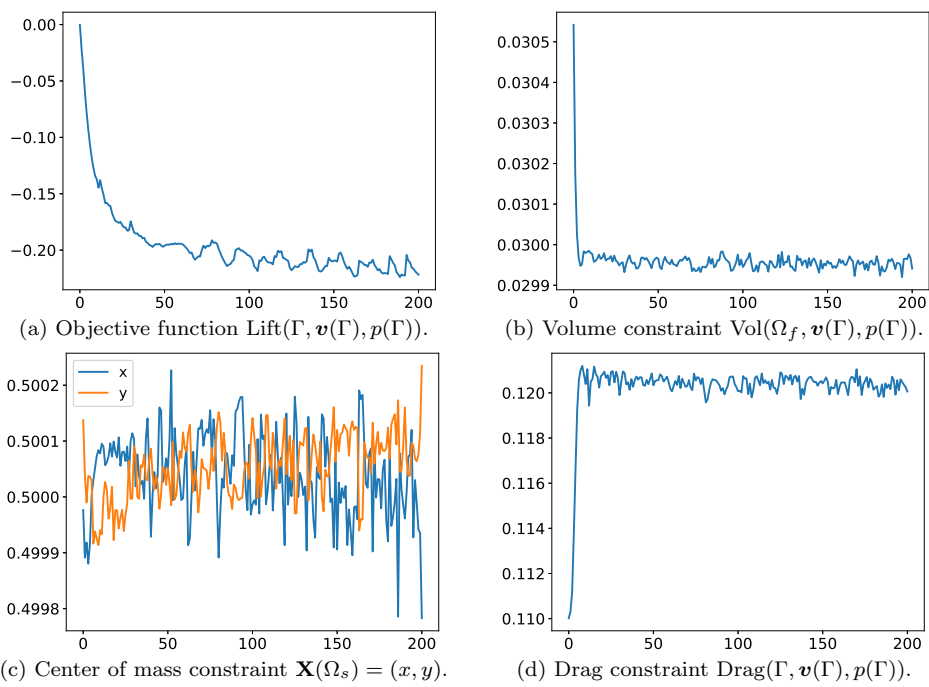


Figure 6.18: Convergence histories for the 2-d lift-drag optimization test case of section 6.2.3.



$\mathbf{v} \cdot \mathbf{n} = 0$  is imposed at other boundaries of the cube. The physical parameters are still set to  $\nu = 1/200$  and  $\rho = 1$  corresponding to  $\text{Re} = \frac{\rho h v_{\max}}{\nu} = 200$  (the characteristic length is  $h = 1$ ). The volume fraction target is set to  $V_{\text{target}} = 0.01$ . For this configuration, the computed value for the minimum drag problem is  $\text{DRAG}^* = 0.0304$  which yields an upper bound  $\text{DRAG}_0 = 1.5 \times \text{DRAG}^* = 0.0456$ .

The optimized shape with associated velocity profile is shown on [Figure 6.19](#). The final fluid mesh features 54,299 vertices. The finite element problems associated with the Navier-Stokes and adjoint equations are solved by using the domain decomposition technique described in [section 6.1.4](#) with 12 cpus (2.60GHz). Convergence histories are shown on [Figure 6.18](#). A single resolution of the state equations including domain decomposition and the Newton loop with the Newton method takes approximately 2 minutes. The computation of the shape derivative, including adjoining system resolutions takes approximately the same time. Every remeshing step (not performed in parallel) is achieved within approximately one minute.

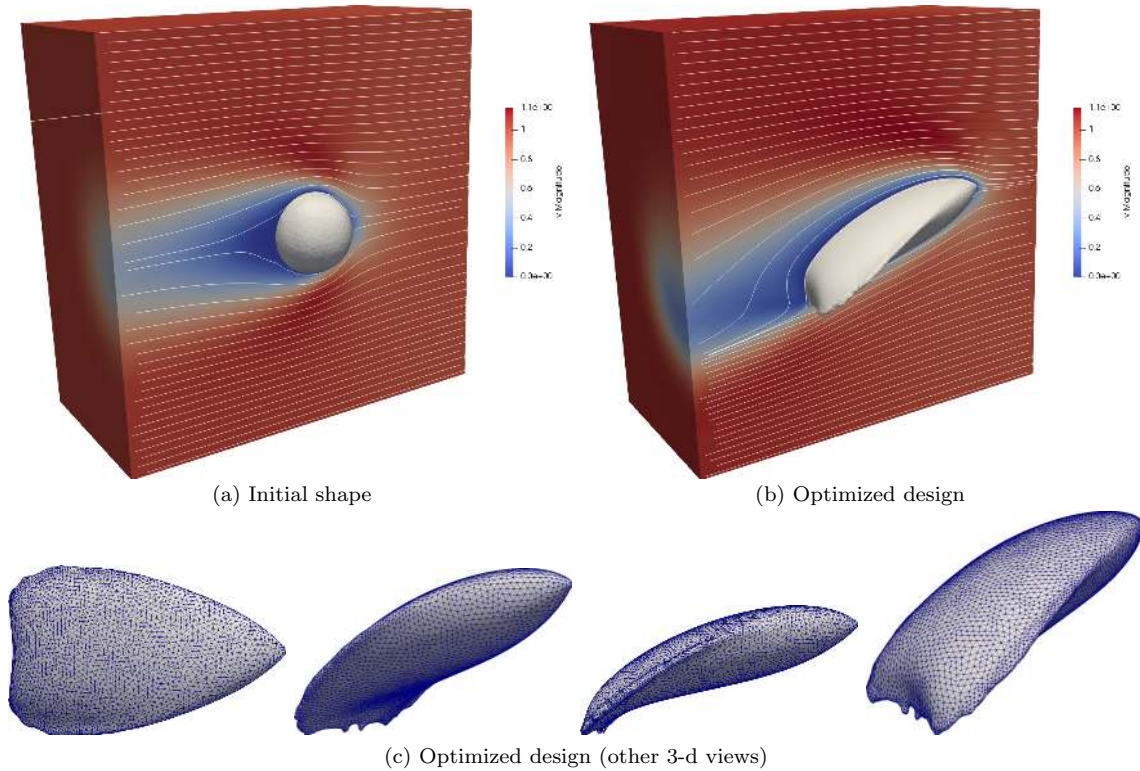


Figure 6.19: Optimized shape for the 3-d lift-drag maximization problem of [section 6.2.3](#).

### 6.2.4 A 3-d fluid-structure interaction test case

Our last test case is concerned with a rather large-scale fluid-structure interaction problem. A fluid is entering the box  $D = [0, 4] \times [0, 1] \times [0, 1]$  with an input velocity  $\mathbf{v}_0 = y\mathbf{e}_x$  on the left-hand boundary. A no-slip boundary condition  $\mathbf{v} = 0$  is prescribed on the bottom face of the domain. The top and side faces assume a slip boundary condition  $\mathbf{v} \cdot \mathbf{n} = 0$ . The flow exits the domain with a zero normal stress boundary condition  $\sigma_f(\mathbf{v}, p) \cdot \mathbf{n} = 0$ . A mechanical structure  $\Omega_s \subset D$  is fixed ( $\mathbf{u} = 0$ ) on a square patch of the bottom face and is subjected to the stress induced by the fluid (namely,  $\mathbf{u}$  is the solution to [\(2.2.3\)](#)).

A vertical plate is set as a non optimizable part of the mechanical structure (as well as a small layer above the bottom Dirichlet boundary): the setting is made visible on [Figure 6.21](#). The goal is to find how to distribute additional material in order to make the structure  $\Omega_s$  as rigid as possible. The problems features of course a volume constraint on the mechanical structure, so that it reads

$$\begin{aligned} \min_{\Gamma} \quad & J(\Gamma, \mathbf{u}(\Gamma)) = \int_{\Omega_s} A\mathbf{e}(\mathbf{u}) : \mathbf{e}(\mathbf{u}) dx \\ \text{s.t.} \quad & \text{Vol}(\Omega_s) = V_{\text{target}}. \end{aligned} \tag{6.2.14}$$

The Reynolds number, fluid density and viscosity are respectively set to  $\text{Re} = 60$ ,  $\rho = 1$  and  $\nu = \rho h \|\mathbf{v}_0\|_{\infty} / \text{Re} = 0.012$  (the characteristic length is  $h = 0.7$  corresponding to the height of the non-

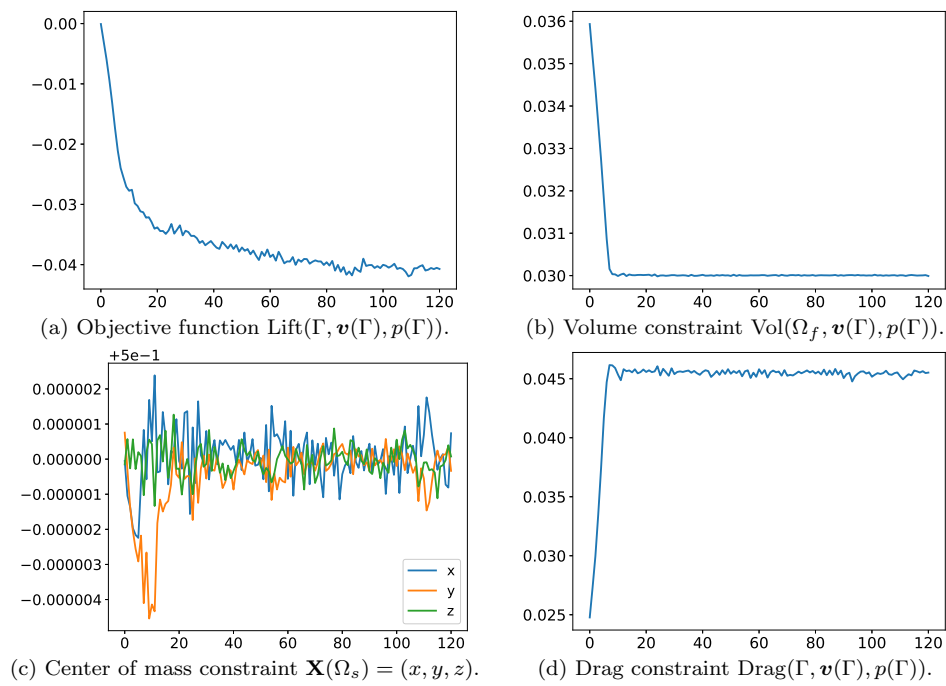


Figure 6.20: Convergence histories for the 3-d lift-drag optimization test case of section 6.2.3.

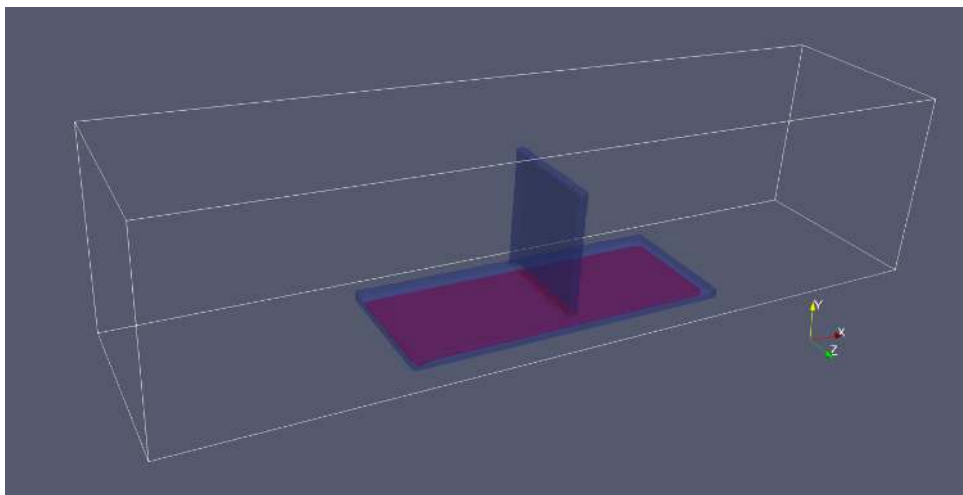


Figure 6.21: Setting of the fluid-structure interaction test case. A flow is entering from the left-hand side; a no-slip boundary condition is imposed at the bottom wall, and a slip boundary condition on the other side walls. The flow is pushing against a non optimizable vertical mechanical plate tightened to the bottom wall with a zero Dirichlet boundary condition (on the red surface).

Solving state equations	24 cpus	15'
Computing shape derivatives and gradients	24 cpus	2'
Advection of the level-set function	sequential	24"
Symmetrization and regularization of the level-set function	24 cpus	13"
Remeshing	sequential	2'
Computation of the signed distance function	sequential	3'

Table 6.2: Running times for the first iteration (220,283 mesh nodes) of the fluid-structure interaction test case of [section 6.2.4](#).

optimizable plate and  $\|\mathbf{v}_0\|_\infty = 1$ ). The Lamé coefficients of the mechanical structure are  $\lambda = 0.00529$  and  $\mu = 0.0476$ .

The optimized shape obtained with the *volume* expression of the shape derivative is plotted on [Figure 6.22](#). Note that for this example, we did not see a clear difference in performance between the results arising from the use of the volume and the surface expressions of the shape derivative. Not surprisingly, the final design has an aerodynamic profile in order to reduce the stress applied by the fluid flow.

This example is our most computationally involved test case. Finite element computations were run in parallel on 24 cpus (2.60GHz) . The number of mesh nodes varies from 220,283 for the first iteration (including 132,775 nodes in the fluid domain, which means approx. 1,7 millions degrees of freedom for the linearized fluid system) to 82,454 (including 66,021 nodes in the fluid domain) at the last iteration. Running cpu times for the first mesh (220,283 mesh nodes) are listed in [Table 6.2](#). The most intensive task is the resolution of the state equations (the running time mentioned includes all domain decomposition steps, finite element matrix assembly, the Newton loop for the Navier-Stokes system, etc. . .).

Although these results are promising, we are still far from realistic industrial system sizes (featuring the order of the billion of elements). We had difficulties in going to mesh sizes reaching the million of vertices because the cost of remeshing becomes prohibitive. So far, this step is performed sequentially; future works could benefit from making this step run in parallel as well.

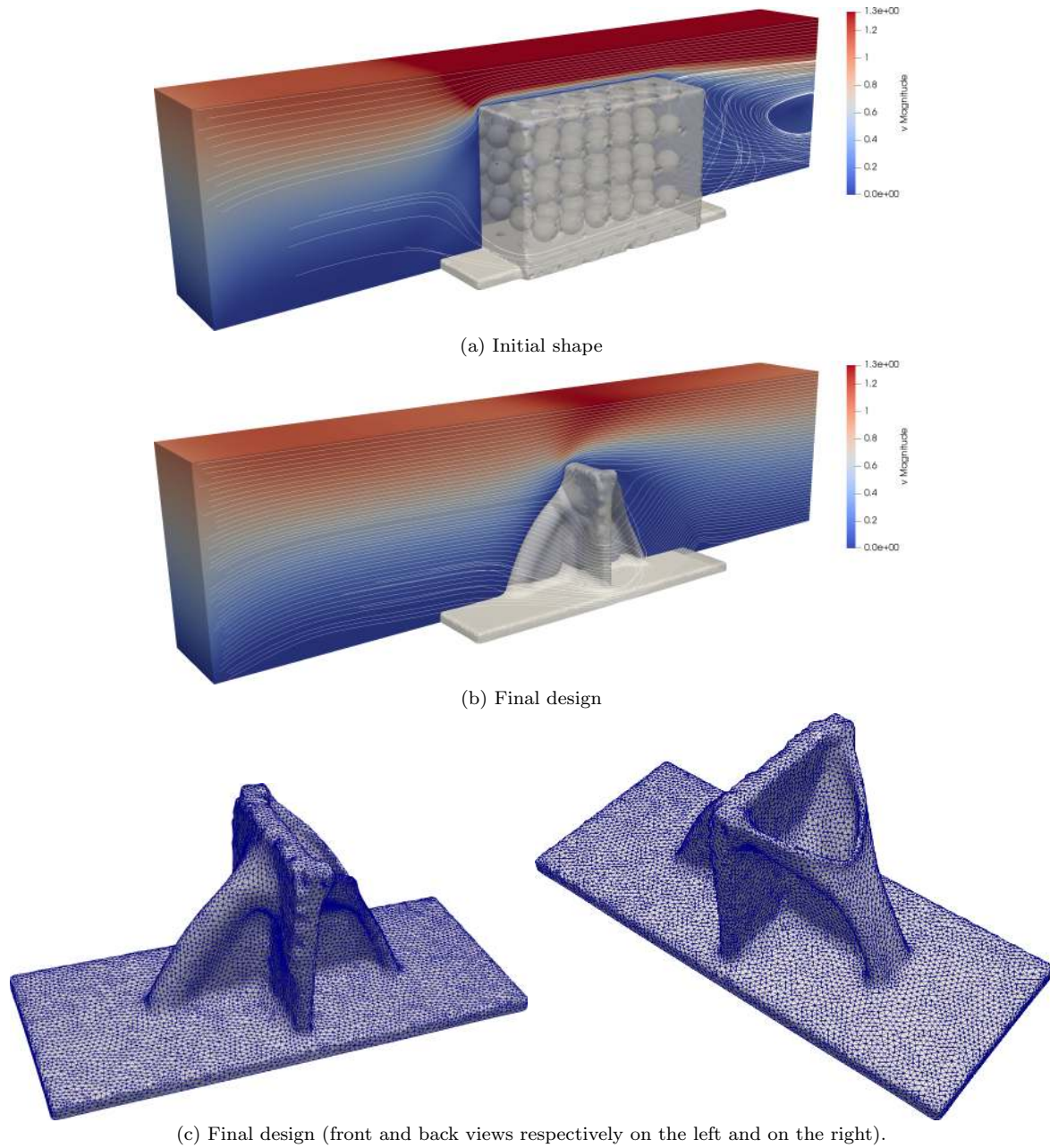


Figure 6.22: Optimized design for the 3-d fluid structure interaction test case of [section 6.2.4](#).

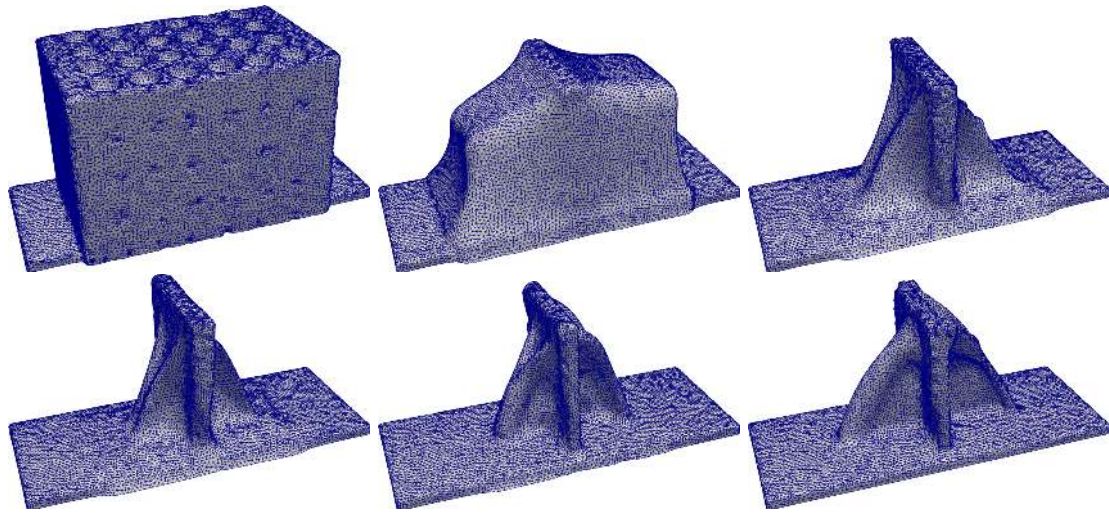


Figure 6.23: Intermediate iterations 0, 40, 100, 125, 175 and 300 for the fluid structure test case of (6.2.14).

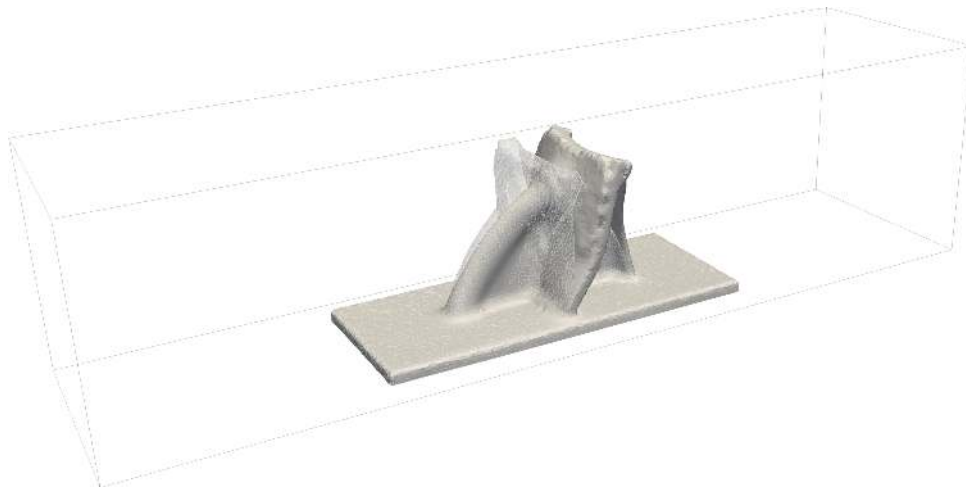


Figure 6.24: Linear elastic deformation of the solid structure under the load force applied by the fluid for the test case of (6.2.14).

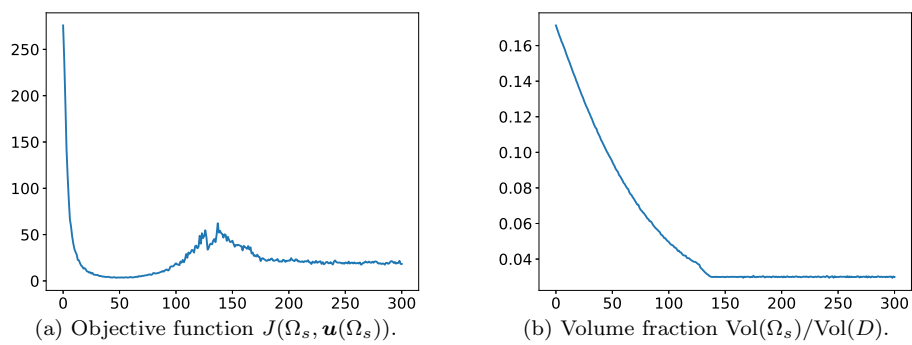


Figure 6.25: Convergence history for the fluid-structure interaction optimization test case of section 6.2.4.



## CHAPTER 7

# HIGH ORDER HOMOGENIZED EQUATIONS FOR PERFORATED PROBLEMS: TOWARDS FLUID TOPOLOGY OPTIMIZATION BY THE HOMOGENIZATION METHOD

### Contents

---

<b>7.1</b>	<b>Introduction</b> . . . . .	<b>239</b>
<b>7.2</b>	<b>Motivations from shape optimization and summary of results</b> . . . . .	<b>240</b>
7.2.1	Shape optimization by the homogenization method in linear elasticity . . . . .	241
7.2.2	Three homogenized regimes for a Stokes flow in a porous medium . . . . .	242
7.2.3	High order homogenized equations for the Stokes problem: summary of results . . . . .	243
7.2.4	Setting and notation conventions related to tensors . . . . .	246
<b>7.3</b>	<b>High order homogenization for the perforated Poisson problem</b> . . . . .	<b>249</b>
7.3.1	Formal infinite order two-scale expansions and tensors $\mathcal{X}^k$ . . . . .	251
7.3.2	Formal infinite order homogenized equation and criminal ansatz: tensors $M^k$ and $N^k$ . . . . .	255
7.3.3	Homogenized equations of order $2K + 2$ : tensors $\mathbb{B}_K$ and $\mathbb{D}_K$ . . . . .	258
7.3.4	Error estimates and justification of the higher order homogenization process . . . . .	263
7.3.5	Low volume fraction limits when the size of the obstacle tends to zero . . . . .	265
7.3.6	Simplifications for the tensors $\mathcal{X}^{k*}$ and $M^k$ in case of symmetries . . . . .	270
<b>7.4</b>	<b>High order homogenization for the perforated elasticity system</b> . . . . .	<b>271</b>
7.4.1	Formal infinite order two-scale asymptotics and matrix valued tensors $\mathcal{X}^k$ . . . . .	272
7.4.2	High order homogenized equations: tensors $M^k, N^k, \mathbb{B}_K$ and $\mathbb{D}_K$ . . . . .	275
7.4.3	Low volume fraction limits when the size of the obstacle tends to 0 . . . . .	279
7.4.4	Simplifications for the tensors $\mathcal{X}^{k*}$ and $M^k$ in case of symmetries . . . . .	283
7.4.5	Appendix: numerical evidences for the “very strange” tensors $\mathcal{X}^{1*}$ and $M^1$ being nonzero . . . . .	285
<b>7.5</b>	<b>High order homogenization for the Stokes system in a porous medium</b> . . . . .	<b>287</b>
7.5.1	Formal infinite order two scale expansions: tensors $(\mathcal{X}^k, \alpha^k)$ . . . . .	288
7.5.2	Higher order homogenized equations: tensors $M^k, N^k, \beta^k$ and $\mathbb{D}_K$ . . . . .	293
7.5.3	Error estimates and justifications of the higher order homogenization process . . . . .	298
7.5.4	Low volume fraction limits when the size of the obstacle tends to 0 . . . . .	304
7.5.5	Simplifications for the tensors $\mathcal{X}^{k*}$ and $M^k$ under symmetries . . . . .	308
7.5.6	Appendix: extension to multicomponent fluid domains . . . . .	310

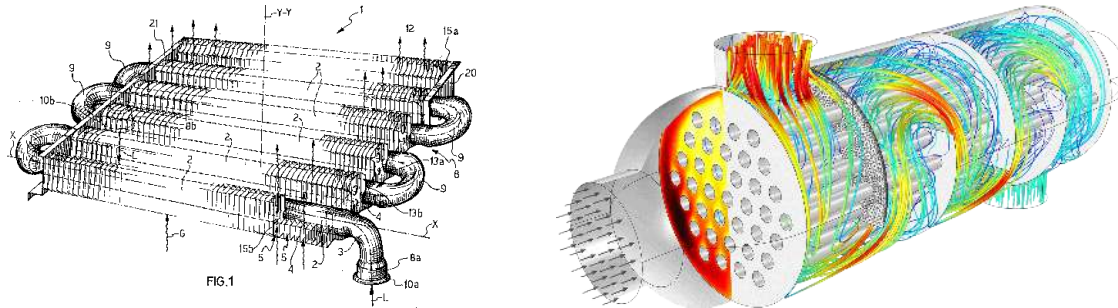
---

### 7.1 INTRODUCTION

This final chapter is an opening towards the use of topology optimization for the design of fluid systems by the homogenization method. Our motivation originates from the observation that many industrial applications in the aeronautic industry involve multi-scale designs. For instance, heat exchangers feature periodic patterns visible at a microscopic scale which are geometrically modulated over larger scales (Figure 7.1). They are integrated into a suitable macroscopic structure so as to maximize the exchange surface between hot and cold phases, while limiting the output pressure loss.

When it comes to automatically generate such multi-scale structures, the topology optimization techniques based on the method of Hadamard described in the previous chapters reach their limits. Indeed, the numerical optimization of the shapes of highly resolved composite structures would require the use of very fine meshes so as to capture the most microscopic details of the fluid flowing in the whole structure and to update the tiniest details of its boundaries; the numerical computation would be very costly and probably very slow to converge.

The goal of this chapter is to lay down theoretical material that would allow, ultimately, to design fluid systems by homogenization methods similar to those available in the context of topology optimization of mechanical structures [65, 58, 18, 254, 27]. In the latter contexts, they specifically allow to generate



(a) Industrial gas-liquid heat exchanger design featuring blade patterns (Figure from [169]). (b) Industrial air-water heat exchanger design featuring tube patterns (Figure from COMSOL multiphysics [239]).

Figure 7.1: Two examples of heat exchanger designs featuring periodically repeated patterns.

multi-scale designs by optimizing simultaneously both the microstructure of a mechanical system and its macroscopic shape. Unfortunately, these methods do not extend in a straightforward manner to fluid systems, because the homogenization theory which accounts for the effective physics of the fluid in a porous medium is different. Several homogenized models exist depending on various scaling regimes assumed by the microstructure pattern (namely, the Darcy, Brinkman, or Stokes regimes [272, 14]) which makes it unclear which effective model should be used to describe a context featuring all possible regimes simultaneously at different locations in the domain.

In this chapter, we derive high order homogenized models for the Stokes system (in periodic domains) which unify the three classical regimes. The mathematical methodology is inspired from the works of Bakhvalov and Panasenko [53], Smyshlyaev and Cherednichenko [287], and Allaire et. al. [33]; it is based on (non standard) two-scale asymptotic expansions and formal operations on related power series which give rise to several families of tensors and homogenized equations at any order. These formal expansions are then justified thanks to rigorous error estimates.

The obtained higher order homogenized models are higher order corrections of the Darcy model: in the low-volume fraction limit where the obstacle size assume one of the three classical regimes of the literature, they specialize to either the Stokes, the Brinkman, or the Darcy equation. Since these higher order formulations unify the three regimes, we expect these could open the way to the development of homogenization methods for the topology optimization of fluid systems. Furthermore, a by-product of the use of higher order models is to yield a more accurate description of the effective physics characterizing porous media in contexts where the size of the microstructure (a parameter  $\varepsilon$  which is assumed to be close to zero in the homogenized model) is not so small.

The remainder of the chapter is organized as follows. Section 7.2 outlines the motivations for the use of higher order homogenized model for the Stokes system and provides a summary of our main mathematical results. The notation and mathematical setting used in the remainder of the chapter are also introduced.

In order to highlight the main features of our derivation, and for pedagogical purposes, the next sections investigate higher order homogenized equations for several elliptic problems with an increasing order of complexity. We start with a study of the perforated Poisson problem in section 7.3, which is a simplified, scalar version of the Stokes system. In Section 7.4 we consider the elasticity system with Dirichlet boundary conditions on the holes (rather than Neumann ones as is more customarily considered). It is the direct vectorial analogue of the Poisson system which is characterized by an important difference: high order homogenized models of the scalar problem feature only differential operators of *even* order, while differential operators of *odd* orders are present in the vectorial context. However these vanish in case of symmetries. Finally, the Stokes system in a porous medium (7.2.13) is treated in section 7.5, which requires additional work due to the pressure variable and the associated divergence constraint.

## 7.2 MOTIVATIONS FROM SHAPE OPTIMIZATION AND SUMMARY OF RESULTS

This section motivates and present our main results developed in thorough details in the next parts. We start by providing in section 7.2.1 a brief account of the homogenization method for topology optimization of mechanical structures. The difficulties that arise when considering the application of these methods



for the design of fluid systems, related to the emergence of three different homogenized regimes, are stressed in section 7.2.2. We then provide a summary of our main results in section 7.2.3 regarding the high order homogenization of Stokes flow in porous media which could address some of these difficulties in future works. Finally, the general mathematical setting and the notation conventions used in the remainder of the work are introduced in section 7.2.4.

### 7.2.1 Shape optimization by the homogenization method in linear elasticity

The homogenization method is a powerful technique for the topology optimization of mechanical structures: it allows to capture multi-scale designs characterized by geometrically modulated micro-structures [65, 18, 254, 177, 27]. Let  $D = [-L, L]^d$  be a given computational rectangular domain and consider the problem of finding the shape assumed by a distribution of two materials characterized by constant elastic Hooke's tensors  $A_0$  and  $A_1$  in the respective domains  $\Omega \subset D$  and  $D \setminus \Omega$ :

$$\begin{aligned} \min_{\Omega \subset D} \quad & J(\Omega, \mathbf{u}(\Omega)) \\ \text{s.t.} \quad & \begin{cases} -\operatorname{div}(A(\mathbf{1}_\Omega)\nabla\mathbf{u}) = \mathbf{f} \text{ in } D \\ \mathbf{u} = 0 \text{ on } \partial D, \end{cases} \end{aligned} \quad (7.2.1)$$

The elastic tensor in (7.2.1) is given by  $A(\mathbf{1}_\Omega) := \mathbf{1}_\Omega A_0 + (1 - \mathbf{1}_\Omega)A_1$  where  $\mathbf{1}_\Omega$  denotes the characteristic function of  $\Omega$  and  $\mathbf{f}$  is a source term. The main point of the homogenization method is to replace the design problem (7.2.1) of finding *shapes*  $\Omega$  by a relaxed one expressed in terms of composite *microstructures* which are not shapes but rather limits of minimizing sequence of shapes. In theory, one can formulate a relaxed version of (7.2.1) posed over the whole set of such possible limits (the  $G$ -closure, see [18]). However, for practical applications, it may be sufficient to consider periodic microstructures with a geometrically modulated pattern: such microstructures can be parameterized by varying parameters  $(a_1(x), \dots, a_m(x))$ , as illustrated on Figure 7.2. This allows to replace (7.2.1) with a ‘‘homogenized’’

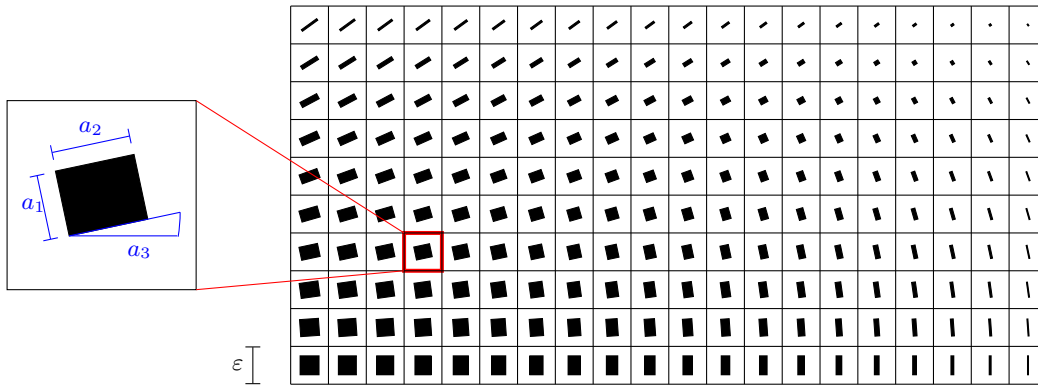


Figure 7.2: A periodic medium with varying microstructure parameterized by the size  $(a_1, a_2)$  and orientation  $a_3$  of a rectangular inclusion. Note that in [27], the parameter  $a_3$  rather parameterizes the orientation of the periodic cells themselves than only their inner rectangular pattern.

relaxed problem

$$\begin{aligned} \min_{(a_1, \dots, a_m) \in L^\infty(D, \mathbb{R}^d)} \quad & J^*(a_1, \dots, a_m, \mathbf{u}(a_1, \dots, a_m)) \\ \text{s.t.} \quad & \begin{cases} -\operatorname{div}(A^*(a_1, \dots, a_m)\nabla\mathbf{u}) = \mathbf{f} \text{ in } D \\ \mathbf{u} = 0 \text{ on } \partial D, \end{cases} \end{aligned} \quad (7.2.2)$$

where  $A^*(a_1, \dots, a_m)$  is an effective tensor characterizing the stiffness of the composite material obtained by mixing  $A_0$  and  $A_1$  according to the microstructure pattern determined by  $a_1, \dots, a_m$ .

The problem (7.2.2) is easier to solve than (7.2.1): it reduces to rather standard parametric optimization with respect to the  $(a_1, \dots, a_m)$ . Furthermore, the state variable  $\mathbf{u}$  can be evaluated in the ‘‘homogeneous’’ domain  $D$ , without the need for discretizing all the geometric details of the microstructure. The difficulty of the method lies the post-treatment needed to actually reconstruct a sequence of shapes approaching the optimal parameterized microstructure; see the contributions [254, 35, 166] for a description of some (delicate) algorithms that allow to perform this task in both 2-d and 3-d.

### 7.2.2 Three homogenized regimes for a Stokes flow in a porous medium

The ultimate motivation would be to apply similar methods for the topology optimization of fluid systems. The optimization problem at play would be formulated in terms of the Stokes (or the Navier-Stokes) equation (instead of (7.2.1)):

$$\begin{aligned} \min_{\Omega \subset D} \quad & J(\Omega, \mathbf{u}(\Omega), p(\Omega)), \\ \text{s.t.} \quad & \begin{cases} -\Delta \mathbf{u} + \nabla p = \mathbf{f} \text{ in } \Omega \\ \operatorname{div}(\mathbf{u}) = 0 \text{ in } \Omega, \\ \mathbf{u} = 0 \text{ on } \partial\Omega, \end{cases} \end{aligned} \tag{7.2.3}$$

where  $(\mathbf{u}, p)$  denotes the fluid velocity and pressure couple (extended by zero outside the optimization domain  $\Omega$ ). However, it is not straightforward to derive a relaxed formulation such as (7.2.3) because the available homogenization theory for the Stokes system (7.2.3) is substantially different to the one for the elasticity system (7.2.1). This difference comes from the fact that the shape  $\Omega$  characterizing the mixture (of material  $A_0$  and  $A_1$  for the elasticity system, or of fluid and solid in the solid system) arises in the elasticity problem (7.2.1) through the piecewise constant Hooke’s tensor  $A(\mathbf{1}_\Omega)$ , whereas it occurs in the Stokes system (7.2.3) in terms of a *zero Dirichlet* (no-slip) boundary condition  $\mathbf{u} = 0$  on the boundary of the holes  $\partial\Omega$ . In contrast with (7.2.1), there is no known theory that would characterize the  $G$ -closure of the Stokes system featured in (7.2.3), i.e. the set of all possible limit models obtained by considering minimizing sequences of domains. Instead, the literature [272, 12, 14] describes the occurrence of three possible homogenized regimes (described below) depending on the scaling ratio between the size of the periodic inclusions and their relative distance.

To be more precise, let us place ourselves in the classical context of periodic homogenization for (7.2.3), which as motivated above, may be sufficient for shape optimization purposes. We assume  $D := [-L, L]^d$  to be a  $d$ -dimensional box filled with periodic obstacles  $\omega_\varepsilon := \varepsilon(\mathbb{Z}^d + \eta T) \cap D$ . The fluid domain  $Y$  is assumed to be connected, see the introduction of section 7.5 for precise assumptions. we denote by  $P = (0, 1)^d$  the unit cell, and by  $Y := P \setminus \eta T$  the unit perforated cell. The boundary of the obstacle  $T$  is assumed to be smooth. The setting is illustrated on Figure 7.3. The parameter  $\varepsilon$  denotes the size of the periodic cell and is assumed to be small. The parameter  $\eta$  is another rescaling of the obstacle  $T$  within each cell of size  $\varepsilon$  (which will allow us later to consider the so-called low volume fraction limit case where  $\eta$  is small comparatively to  $\varepsilon$ ). The working, porous domain is denoted by  $D_\varepsilon := D \setminus \omega_\varepsilon$ , and we consider

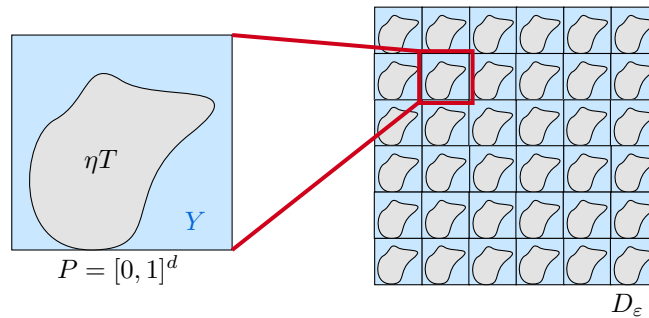


Figure 7.3: The perforated domain  $D_\varepsilon$  and the unit cell  $Y = P \setminus (\eta T)$ .

the solution  $(\mathbf{u}_\varepsilon, p_\varepsilon)$  of the following Stokes system:

$$\begin{cases} -\Delta \mathbf{u}_\varepsilon + \nabla p_\varepsilon = \mathbf{f} \text{ in } D_\varepsilon \\ \operatorname{div}(\mathbf{u}_\varepsilon) = 0 \\ \mathbf{u}_\varepsilon = 0 \text{ on } \partial\omega_\varepsilon \\ \mathbf{u}_\varepsilon \text{ is } D\text{-periodic,} \end{cases} \tag{7.2.4}$$

where  $\mathbf{f}$  is now required to be a smooth,  $D$ -periodic right hand-side. The periodicity assumption for  $\mathbf{u}_\varepsilon$  is classical in homogenization and is used to avoid difficulties related to the arising of boundary layers (see [216, 326, 19]). The literature accounts for several homogenized equations depending on how the size  $a_\varepsilon = \eta\varepsilon$  of the holes compares to the critical size  $\sigma_\varepsilon := \varepsilon^{d/(d-2)}$  in dimension  $d \geq 3$  or  $\sigma_\varepsilon := \exp(-1/\varepsilon^2)$  for  $d = 2$  [242, 105, 11, 15]:

- if  $a_\varepsilon = o(\sigma_\varepsilon)$ , then the holes are “too small” and  $(\mathbf{u}_\varepsilon, p_\varepsilon)$  converges as  $\varepsilon \rightarrow 0$  to the solution  $(\mathbf{u}, p)$  of the Stokes equation in the homogeneous domain  $D$ :

$$\begin{cases} -\Delta \mathbf{u} + \nabla p = \mathbf{f} & \text{in } D \\ \operatorname{div}(\mathbf{u}) = 0 \\ \mathbf{u} \text{ is } D\text{-periodic.} \end{cases} \quad (7.2.5)$$

- if  $a_\varepsilon = \sigma_\varepsilon$ , then  $(\mathbf{u}_\varepsilon, p_\varepsilon)$  converges as  $\varepsilon \rightarrow 0$  to the solution  $(\mathbf{u}, p)$  of the Brinkman equation

$$\begin{cases} -\Delta \mathbf{u} + \Psi^* \mathbf{u} + \nabla p = \mathbf{f} & \text{in } D \\ \operatorname{div}(\mathbf{u}) = 0 \\ \mathbf{u} \text{ is } D\text{-periodic,} \end{cases} \quad (7.2.6)$$

where the so-called *strange term*  $\Psi^* \mathbf{u}$  involves a symmetric positive definite  $d \times d$  matrix  $\Psi^*$  that can be computed by means of an exterior problem in  $\mathbb{R}^d \setminus T$  when  $d \geq 3$ , and which is equal to  $\pi I$  if  $d = 2$  (see [11]).

- if  $\sigma_\varepsilon = o(a_\varepsilon)$  and  $a_\varepsilon = \eta \varepsilon$  with  $\eta \rightarrow 0$  as  $\varepsilon \rightarrow 0$ , then the holes are “large” and  $(a_\varepsilon^{d-2} \varepsilon^{-d} \mathbf{u}_\varepsilon, p_\varepsilon)$  converges to the solution  $(\mathbf{u}, p)$  of the Darcy problem

$$\begin{cases} \Psi^* \mathbf{u} + \nabla p = \mathbf{f} & \text{in } D \\ \operatorname{div}(\mathbf{u}) = 0 & \text{in } D \\ \mathbf{u} \text{ is } D\text{-periodic,} \end{cases} \quad (7.2.7)$$

where  $\Psi^*$  is the same symmetric positive definite  $d \times d$  matrix as in (7.2.6).

- if  $a_\varepsilon = \eta \varepsilon$  with the ratio  $\eta$  fixed, then  $(\varepsilon^{-2} \mathbf{u}_\varepsilon, p_\varepsilon)$  converges to the solution  $(\mathbf{u}, p)$  of the Darcy problem

$$\begin{cases} M^0 \mathbf{u} + \nabla p = \mathbf{f} & \text{in } D \\ \operatorname{div}(\mathbf{u}) = 0 & \text{in } D \\ \mathbf{u} \text{ is } D\text{-periodic,} \end{cases} \quad (7.2.8)$$

where  $M^0$  is another positive symmetric  $d \times d$  matrix (which depends on  $\eta$ ). Furthermore it can be shown that  $M^0 / |\log(\eta)| \rightarrow \Psi^*$  if  $d = 2$ , and  $M^0 / \eta^{d-2} \rightarrow \Psi^*$  (if  $d \geq 3$ ) when  $\eta \rightarrow 0$ , so that there is a continuous transition from (7.2.8) to (7.2.7), see [13].

These three different regimes, namely Stokes, Brinkman, and Darcy, raise practical difficulties in view of applying the homogenization method for shape optimization. Indeed, it is not clear which regime to use in order to write down a relaxed problem (7.2.2) in a context where the shape of the holes  $\eta T \equiv \eta(x)T(x)$  would be allowed to be modulated along the domain  $D$ : in regions featuring very tiny obstacles, one should use the Stokes or Brinkman equation (7.2.5) and (7.2.6), however one should use the Darcy model (7.2.8) when the obstacles become large enough.

In fact, there is a continuous transition between all regimes which can be captured by higher order homogenized equations, and which is the purpose of the present chapter. Corrective terms scaled by higher powers of  $\varepsilon$  can be added to the Darcy equation (7.2.8), which reduce to one of (7.2.5) to (7.2.7) in the other scaling regimes considered.

### 7.2.3 High order homogenized equations for the Stokes problem: summary of results

The main result of this chapter is the derivation of higher order homogenized equations for the Stokes system (7.2.4). For a desired order  $K \in \mathbb{N}$ , there exists a homogenized model of order  $2K + 2$  which can be written

$$\begin{cases} \sum_{k=0}^{2K+2} \varepsilon^{k-2} \mathbb{D}_K^k \cdot \nabla^k \mathbf{v}_K^* + \nabla q_K^* = \mathbf{f}, \\ \operatorname{div}(\mathbf{v}_K^*) = 0, \\ \mathbf{v}_K^* \text{ is } D\text{-periodic,} \end{cases} \quad (7.2.9)$$

and which determines a so-called higher order homogenized solution  $(\mathbf{v}_K^*, q_K^*)$  (definition 7.8). Note that  $(\mathbf{v}_K^*, q_K^*)$  depends on  $\varepsilon$ , but this dependence is implicitly understood for simplicity of notation. The notation

$$\mathbb{D}_K^k := (\mathbb{D}_{K, i_1 \dots i_k, lm}^k)_{1 \leq i_1 \dots i_k \leq d, 1 \leq l, m \leq d}$$

designates a  $k$ -th order matrix valued tensor involving  $k$  indices  $i_1 \dots i_k$  associated with partial derivatives and two indices  $l, m$  associated with spatial coordinates; it can be computed thanks to a procedure involving the resolution of cell problems. The  $k$ -th order differential operator  $\mathbb{D}_K^k \cdot \nabla^k$  is obtained by contracting the  $k$  indices  $i_1 \dots i_k$  with  $k$  partial derivatives and performing a matrix product with the other two indices  $l, m$ :

$$(\mathbb{D}_K^k \cdot \nabla^k \mathbf{v}_K^*)_l := \mathbb{D}_{K, i_1 \dots i_k, lm}^k \partial_{i_1 \dots i_k}^k v_{K, m}^*,$$

where  $(v_{K, m}^*)_{1 \leq m \leq d}$  denote the components of  $\mathbf{v}_K^*$  and an implicit summation over the repeated indices  $i_1, \dots, i_k$  and  $m$  is assumed.

Equation (7.2.9) is said to be of “higher order” (namely of order  $2K + 2$  with  $K \in \mathbb{N}$ ) because  $\mathbf{v}_K^*$  (respectively  $q_K^*$ ) yields an approximation of  $\mathbf{u}_\varepsilon$  (respectively  $p_\varepsilon$ ) of order  $\varepsilon^{K+3}$  (respectively  $\varepsilon^{K+1}$ ) in the  $L^2(D)$  norm: in proposition 7.39 below, we prove the following error bounds with a constant  $C_K(\mathbf{f})$  that depends only on  $\mathbf{f}$ ,  $K$  (and *a priori* on the shape of the obstacle  $\eta T$ ):

$$\left\| \mathbf{u}_\varepsilon - \sum_{k=0}^K \varepsilon^k N^k(\cdot/\varepsilon) \cdot \nabla^k \mathbf{v}_K^* \right\|_{H^1(D, \mathbb{R}^d)} \leq C_K(\mathbf{f}) \varepsilon^{K+2}$$

$$\left\| p_\varepsilon - \left( q_K^* + \sum_{k=0}^{K-1} \varepsilon^{k-1} \boldsymbol{\beta}^k(\cdot/\varepsilon) \cdot \nabla^k \mathbf{v}_K^* \right) \right\|_{L^2(D)} \leq C_K(\mathbf{f}) \varepsilon^{K+1}.$$

The variables  $N^k(\cdot/\varepsilon)$  and  $\boldsymbol{\beta}^k(\cdot/\varepsilon)$  are respectively matrix and vector valued corrector tensors which are periodic in the cell  $P$  (Figure 7.3) and which do not depend on  $\mathbf{f}$ ,  $\mathbf{v}_K^*$ ,  $p_K^*$ . Furthermore,  $N^k$  vanishes on the boundary of the obstacle  $\partial(\eta T)$  so that a function multiplied by  $N^k(\cdot/\varepsilon)$  vanishes on the boundary of the holes  $\partial\omega_\varepsilon$ .

Our methodology is based on the existence of “criminal” ansatz for the velocity and pressure solution  $(\mathbf{u}_\varepsilon, p_\varepsilon)$ . The “classical” ansatz reads formally

$$\mathbf{u}_\varepsilon = \sum_{i=0}^{+\infty} \varepsilon^{i+2} \mathbf{u}_i(x, x/\varepsilon), \quad p_\varepsilon(x) = \sum_{i=0}^{+\infty} \varepsilon^i (p_i^*(x) + \varepsilon p_i(x, x/\varepsilon)), \quad x \in D_\varepsilon. \quad (7.2.10)$$

where the functions  $\mathbf{u}_i(x, y)$  and  $p_i(x, y)$  are periodic in the variable  $y \in P$  (see eqn. (7.5.2) below). To our knowledge, very few works have investigated higher order homogenized models for the Stokes equation. Most of the available works have considered situations with low regularity for  $\mathbf{f}$  and  $D$  (see [272, 12]), where the homogenization process can be justified only for the first terms of the ansatz (7.2.10). Error bounds for the higher order terms (namely the result of proposition 7.37) have been obtained in [228, 74]. A few additional works have sought corrector terms from physical modelling considerations [159, 47, 46].

By introducing a suitable family of  $k$ -th order tensors  $(\mathcal{X}^k, \boldsymbol{\alpha}^k)$  obtained as the solutions to cell problems, we obtain that (7.5.2) rewrites more explicitly (see proposition 7.29)

$$\begin{cases} \mathbf{u}_\varepsilon(x) = \sum_{i=0}^{+\infty} \varepsilon^i \mathcal{X}^i(x/\varepsilon) \cdot \nabla^i (\mathbf{f}(x) - \nabla p_\varepsilon^*(x)), \\ p_\varepsilon(x) = p_\varepsilon^*(x) + \sum_{i=0}^{+\infty} \varepsilon^{i+1} \boldsymbol{\alpha}^i(x/\varepsilon) \cdot \nabla^i (\mathbf{f}(x) - \nabla p_\varepsilon^*(x)), \end{cases} \quad x \in D_\varepsilon. \quad (7.2.11)$$

The ansatz for  $\mathbf{u}_\varepsilon$  is a first instance of what [33] called an “asymptotic crime”, because the function  $p_\varepsilon^*$  it features is a homogenized pressure which is itself a formal power series in  $\varepsilon$  (see (7.5.3)).

Despite being explicit, the ansatz (7.2.11) is not fully satisfactory because it requires the knowledge of the partial derivatives of  $f$  at any order, which may be difficult to obtain numerically. One of our

main results is [proposition 7.33](#), where we show that  $(\mathbf{u}_\varepsilon, p_\varepsilon)$  can be formally decomposed as

$$\begin{cases} \mathbf{u}_\varepsilon(x) = \sum_{i=0}^{+\infty} \varepsilon^i N^i(x/\varepsilon) \cdot \nabla^i \mathbf{u}_\varepsilon^*(x), \\ p_\varepsilon(x) = p_\varepsilon^*(x) + \sum_{i=0}^{+\infty} \varepsilon^{i-1} \beta^i(x/\varepsilon) \cdot \nabla^i \mathbf{u}_\varepsilon^*(x) \end{cases} \quad \forall x \in D_\varepsilon. \quad (7.2.12)$$

The ansatz (7.2.12) is not standard and is different from (7.2.11); it is “even more” criminal because it expresses *both* oscillating solution  $\mathbf{u}_\varepsilon, p_\varepsilon$  in terms of their formal, non-oscillating averaged  $\mathbf{u}_\varepsilon^*, p_\varepsilon^*$  (which both depend themselves on  $\varepsilon$ , see (7.5.3)). Since the average of the tensor  $N^0$  over the unit cell  $P$  is the  $2 \times 2$  identity tensor while the average of the other tensors  $(N^k)_{k \geq 1}$  and  $(\beta^k)_{k \geq 0}$  is zero (see [proposition 7.34](#)), formally averaging (7.2.12) with respect to the fast variable  $x/\varepsilon$  shows that the variable  $(\mathbf{u}_\varepsilon^*, p_\varepsilon^*)$  can be indeed interpreted as a formal homogenized average of  $(\mathbf{u}_\varepsilon, p_\varepsilon)$  of “infinite” order.

Criminal ansatz have been first derived in Bakhvalov and Panasenko (1989) [53] concerned with the conductivity equation with heterogeneous diffusivity tensor (and no holes). They have then been used in Smyshlyaev and Cherednichenko (2000) [287] to obtain higher order homogenized equation for a scalar elasticity model with discontinuous Hooke’s tensor, and in [274, 33] in the context of the wave equation.

In [proposition 7.32](#), we obtain that  $(\mathbf{u}_\varepsilon^*, p_\varepsilon^*)$  solves the following formal “infinite-order” homogenized equation,

$$\begin{cases} \sum_{k=0}^{+\infty} \varepsilon^{k-2} M^k \cdot \nabla^k \mathbf{u}_\varepsilon^* + \nabla p_\varepsilon^* = \mathbf{f}, \\ \operatorname{div}(\mathbf{u}_\varepsilon^*) = 0, \\ \mathbf{u}_\varepsilon^* \text{ is } D\text{-periodic.} \end{cases} \quad (7.2.13)$$

which involves a different family of constant tensors  $(M^k)_{k \in \mathbb{N}}$ . Our higher order homogenized models (7.2.9) of order  $2K + 2$  turn out to be truncations of (7.2.13) in the sense that the first  $K + 1$  coefficients coincide:  $\mathbb{D}_K^k = M^k$  for  $0 \leq k \leq K$ . Truncating directly (7.2.13) yields, in general, an ill-posed model [20]. This fact is classical and several methods have been proposed to address this issue and obtain nonetheless well-posed higher order models see e.g. [35, 22, 3, 4, 33]. In our case, we adapt an idea from [287], whereby the coefficients  $\mathbb{D}_K^k$  for  $K + 1 \leq k \leq 2K + 2$  are obtained thanks to a minimization principle which makes indeed (7.2.9) well-posed. Let us stress, however, that this choice is not unique.

A rather surprising feature lies in the fact that (7.2.9) and (7.2.13) feature operators of *odd* order terms. This fact is to relate the vectorial context: the tensors  $\mathbb{D}_K^k$  and  $M^k$  are symmetric and anti-symmetric valued matrices for respectively even and odd  $k$  ([corollary 7.9](#)). This property ensures that eventually, the operators  $\mathbb{D}_K^k \cdot \nabla^k$  and  $M^k \cdot \nabla^k$  are symmetric operators (see [remark 7.20](#)). Tensors  $\mathbb{D}_K^k$  (with  $k = 2p + 1, p \in \mathbb{N}$ ) are *a priori* not zero (this assertion may be difficult to prove, however numerical results supporting this assertion are provided in [section 7.4.5](#) below); we shall see nevertheless that these tensors vanish in the case where the obstacle  $\eta T$  is symmetric with respect to the cell axes.

A legitimate question is to ask whether (7.2.9) or (7.2.13) reduces indeed to either of the equations (7.2.5) to (7.2.8) under the various possible scaling regimes for the sizes of the periodic obstacles. The answer is affirmative let us first mention that keeping the terms of lowest powers in (7.2.13) yields the Darcy model (7.2.8) (the tensor  $M^0$  is the same). Furthermore, we prove the following asymptotics in the low volume fraction limit  $\eta \rightarrow 0$  (in [corollary 7.12](#)) and assuming  $d \geq 3$ :

$$\varepsilon^{-2} M^0 \sim \eta^{d-2} / \varepsilon^2 \Psi^*, \quad (7.2.14)$$

$$\varepsilon^{-1} M^1 = o(\varepsilon(\eta^{d-2} / \varepsilon^2)), \quad (7.2.15)$$

$$\varepsilon^0 M^2 \rightarrow -I, \quad (7.2.16)$$

$$\forall k \geq 1, \varepsilon^{2k-2} M^{2k} = o\left(\left(\frac{\varepsilon^2}{\eta^{(d-2)}}\right)^{k-1}\right), \quad (7.2.17)$$

$$\forall k \geq 1, \varepsilon^{2k-1} M^{2k+1} = o\left(\varepsilon \left(\frac{\varepsilon^2}{\eta^{(d-2)}}\right)^{k-1}\right). \quad (7.2.18)$$

where  $I$  is the  $d$ -by- $d$  identity matrix. The first convergence result (7.2.14) has been obtained in [13], it expresses the continuous transition of the Darcy tensor  $M^0$  in (7.2.7) towards the Brinkman tensor  $\Psi^*$  in

(7.2.6). These asymptotics involve the ratio  $\eta^{d-2}/\varepsilon^2$ , which naturally brings into play the classical critical scaling  $\eta \sim \varepsilon^{2/(d-2)}$  (i.e.  $\sigma_\varepsilon = \varepsilon\varepsilon^{2/(d-2)} = \varepsilon^{d/(d-2)}$  for the holes). Therefore, the Stokes, Brinkman, or Darcy regimes (7.2.5) to (7.2.7) are retrieved in sense that the coefficients  $M^k$  of (7.2.13) converge to those of these three equations as  $\eta, \varepsilon \rightarrow 0$ .

To conclude, let us outline how such higher order homogenized models could be used for topology optimization of fluid flows. For  $K = 0$ , (7.2.9) rewrites in the form of a second order elliptic problem:

$$\begin{cases} \mathbb{D}_0^2 \cdot \nabla^2 \mathbf{v}_0^* + \varepsilon^{-1} \mathbb{D}_0^1 \cdot \nabla \mathbf{v}_0^* + \varepsilon^{-2} \mathbb{D}_0^0 \mathbf{v}_0^* + \nabla q_K^* = \mathbf{f}, \\ \operatorname{div}(\mathbf{v}_0^*) = 0, \\ \mathbf{v}_0^* \text{ is } D\text{-periodic.} \end{cases} \quad (7.2.19)$$

We shall see that (7.2.19) yields the same order of approximation  $O(\varepsilon^3)$  of  $\mathbf{u}_\varepsilon$  as the more classical Darcy equation (7.2.8). However, (7.2.19) makes sense for all three regimes simultaneously: it provides missing corrector terms that are not visible, at first order, in the Darcy regime (7.2.7). Further, it can be verified that the tensors  $\mathbb{D}_0^k$  satisfy the same asymptotics as the tensors  $M^k$  for  $0 \leq k \leq 2$ . The formulation (7.2.19) is expected to be amenable for the development of a numerical topology optimization method: upon a suitable choice of parameterization  $(a_1, \dots, a_m)$  of the microstructure, a relaxed formulation of (7.2.3) could be

$$\begin{aligned} & \min_{(a_1, \dots, a_m) \in L^\infty(D, \mathbb{R}^d)} J^*(a_1, \dots, a_m, \mathbf{v}_K^*(a_1, \dots, a_m), q_K^*(a_1, \dots, a_m)) \\ \text{s.t. } & \begin{cases} \mathbb{D}_0^2(a_1, \dots, a_m) \cdot \nabla^2 \mathbf{v}_0^* + \varepsilon^{-1} \mathbb{D}_0^1(a_1, \dots, a_m) \cdot \nabla \mathbf{v}_0^* + \varepsilon^{-2} \mathbb{D}_0^0(a_1, \dots, a_m) \mathbf{v}_0^* + \nabla q_K^* = \mathbf{f}, \\ \mathbf{v}_K^* \text{ is } D\text{-periodic,} \end{cases} \end{aligned} \quad (7.2.20)$$

where  $J^*$  would be a consistent relaxation of  $J$  and the functions  $(a_1, \dots, a_m) \mapsto \mathbb{D}_0^k(a_1, \dots, a_m)$  could be computed in an off-line step (following [27]). This model is second order for any type of microstructure and contains all three regimes simultaneously, which solves our initial issue. Higher order models could be used similarly to capture contexts where the size  $\varepsilon$  is only moderately small. Our work opens interesting perspectives for the topology optimization of liquid-liquid heat exchangers because our model can be extended in a straightforward manner to multicomponent fluid domains, see the appendix of section 7.5.6) for a discussion.

Finally let us notice that (7.2.20) is a generalization of the model of Borrvall and Petersson [71, 72] which is commonly used in density based methods for fluid topology optimization: it is a version of (7.2.6) where the Brinkman tensor  $\Psi^*$  is assumed to be a scalar and is used as a penalization to enforce  $\mathbf{v} = 0$  in “solid” regions modeled by large values of  $\Psi^*$ . Some variants, however, have been proposed such as the Darcy model [324].

### 7.2.4 Setting and notation conventions related to tensors

In all what follows, the setting considered is that of Figure 7.3.  $D := [-L, L]^d$  is a  $d$ -dimensional box filled with periodic obstacles  $\omega_\varepsilon := \varepsilon(\mathbb{Z}^d + \eta T) \cap D$ .  $P = (0, 1)^d$  is the unit cell, and  $Y := P \setminus \eta T$  the unit perforated cell. The boundary of the obstacle  $T$  is assumed to be smooth. The perforated domain is  $D_\varepsilon := D \setminus \omega_\varepsilon$ .

Below and further on, we consider scalar and vectorial functions such as

$$u : D \times P \rightarrow \mathbb{R}, \quad \mathbf{u} : D \times P \rightarrow \mathbb{R}^d$$

which are both  $D$  and  $P$ -periodic with respect to respectively the first and the second variable. Their respective values are denoted by  $u(x, y)$  and  $\mathbf{u}(x, y)$  where the arguments  $x$  and  $y$  are called respectively the “slow” and the “fast” variable. The partial derivative with respect to the variable  $y_j$  (respectively  $x_j$ ) is written simply  $\partial_j$  instead of  $\partial_{y_i}$  (respectively  $\partial_{x_j}$ ) where the context is clear, i.e. when the function to which it is applied depends only on  $y$  (respectively only on  $x$ ).

The star-“\*”-symbol is used to denote the average of such functions with respect to the  $y$  variable:

$$\mathbf{u}^*(x) := \int_P \mathbf{u}(x, y) dy, \quad u^*(x) := \int_P u(x, y) dy.$$

For functions depending only on the  $y$  variable, it will be sometimes more convenient to write this average with the usual angle bracket symbols. For instance, if  $\mathcal{X} : P \rightarrow \mathbb{R}^d$  is a vector field in the cell  $P$ , then its average is defined by

$$\langle \mathcal{X} \rangle := \int_P \mathcal{X}(y) dy.$$

In all what follows, unless otherwise specified, the Einstein summation convention over repeated *subscript* indices is assumed (but never on *superscript* indices). Vectors of  $\mathbb{R}^d$  are written in bold face notation.

The notation conventions used for tensor related operations are summarized in the nomenclature below.

### Scalar, vector, and matrix valued tensors and their coordinates

$\mathbf{b}$	Vector of $\mathbb{R}^d$
$(b_j)_{1 \leq j \leq d}$	Coordinates of the vector $\mathbf{b}$ .
$b^k$	Scalar valued tensor of order $k$ ( $b_{i_1 \dots i_k}^k \in \mathbb{R}$ for $1 \leq i_1, \dots, i_k \leq d$ )
$\mathbf{b}^k$	Vector valued tensor of order $k$ ( $\mathbf{b}_{i_1 \dots i_k}^k \in \mathbb{R}^d$ for $1 \leq i_1, \dots, i_k \leq d$ )
$B^k$	Matrix valued tensor of order $k$ ( $B_{i_1 \dots i_k}^k \in \mathbb{R}^{d \times d}$ for $1 \leq i_1, \dots, i_k \leq d$ )
$(b_j^k)_{1 \leq j \leq d}$	Coordinates of the vector valued tensor $\mathbf{b}^k$ seen as <i>scalar</i> tensors of order $k$ .
$(B_{lm}^k)_{1 \leq l, m \leq d}$	Coefficients of the matrix valued tensor $B^k$ seen as <i>scalar</i> tensors of order $k$ .
$b_{i_1 \dots i_k, j}^k$	Coefficient of the vector valued tensor $\mathbf{b}^k$ ( $1 \leq i_1, \dots, i_k \leq d$ and $1 \leq j \leq d$ ).
$B_{i_1 \dots i_k, lm}^k$	Coefficient of the matrix valued tensor $B^k$ ( $1 \leq i_1, \dots, i_k \leq d$ and $1 \leq l, m \leq d$ ).

### Tensor products

$b^p \otimes c^{k-p}$  Tensor product of scalar tensors of order  $p$  and  $k$ :

$$(b^p \otimes c^{k-p})_{i_1 \dots i_p} := b_{i_1 \dots i_p}^p c_{i_{p+1} \dots i_k}^{k-p}. \quad (7.2.21)$$

$a^p \otimes \mathbf{b}^{k-p}$  Tensor product of a  $p$ -th order scalar tensor  $a^p$  and a  $k-p$ -th order vector valued tensor  $\mathbf{b}^{k-p}$ :

$$(a^p \otimes \mathbf{b}^{k-p})_{i_1 \dots i_k} := a_{i_1 \dots i_p}^p \mathbf{b}_{i_{p+1} \dots i_k}^{k-p}. \quad (7.2.22)$$

$B^p \otimes C^{k-p}$  Tensor product of the  $p$ -th order and a  $k-p$ -th order  $d \times d$  matrix valued tensors  $B^p = (B_{lm}^p)_{1 \leq l, m \leq d}$  and  $C^{k-p} = (C_{lm}^{k-p})_{1 \leq l, m \leq d}$ :

$$(B^p \otimes C^{k-p})_{i_1 \dots i_k, lm} := B_{i_1 \dots i_p, lj}^p C_{i_{p+1} \dots i_k, jm}^{k-p}, \quad (7.2.23)$$

Hence a matrix product is implicitly assumed in the notation  $B^p \otimes C^{k-p}$ .

$B^p : C^{k-p}$  Tensor product of two  $p$ -th order and  $k-p$ -th order matrix valued tensors:

$$(B^p : C^{k-p})_{i_1 \dots i_k} := B_{i_1 \dots i_p, lm}^p C_{i_{p+1} \dots i_k, lm}^{k-p}. \quad (7.2.24)$$

$\mathbf{b}^p \cdot \mathbf{c}^{k-p}$  Tensor product of two vector valued tensors  $\mathbf{b}^p$  and  $\mathbf{c}^{k-p}$ :

$$(\mathbf{b}^p \cdot \mathbf{c}^{k-p})_{i_1 \dots i_k} := b_{i_1 \dots i_p, m}^p c_{i_{p+1} \dots i_k, m}^{k-p}. \quad (7.2.25)$$

$B^p \cdot \mathbf{c}^{k-p}$  Tensor product of a matrix valued tensor  $B^p$  and a vector valued tensors  $\mathbf{c}^{k-p}$ :

$$(B^p \cdot \mathbf{c}^{k-p})_{i_1 \dots i_k, l} := B_{i_1 \dots i_p, lm}^p c_{i_{p+1} \dots i_k, m}^{k-p}. \quad (7.2.26)$$

### Contraction with partial derivatives

$b^k \cdot \nabla^k$  Differential operator of order  $k$  associated with a scalar tensor  $b^k$ : for any smooth scalar field  $v$ ,

$$b^k \cdot \nabla^k := b_{i_1 \dots i_k}^k \partial_{i_1 \dots i_k}^k. \quad (7.2.27)$$

$\mathbf{b}^k \cdot \nabla^k$  Differential operator of order  $k$  associated with a vector tensor  $\mathbf{b}^k$ : for any smooth vector field  $\mathbf{v}$ ,

$$\mathbf{b}^k \cdot \nabla^k \mathbf{v} = b_{i_1 \dots i_k, l}^k \partial_{i_1 \dots i_k}^k v_l. \quad (7.2.28)$$

$B^k \cdot \nabla^k$  Differential operator of order  $k$  associated with a matrix valued tensor  $B^k$ : for any smooth vector field  $\mathbf{v}$ ,

$$(B^k \cdot \nabla^k \mathbf{v})_l = B_{i_1 \dots i_k, lm}^k \partial_{i_1 \dots i_k}^k v_m. \quad (7.2.29)$$

### Special tensors

$\delta_{ij}$  Kronecker symbol:  $\delta_{ij} = 1$  if  $i = j$  and  $\delta_{ij} = 0$  if  $i \neq j$ .

$I$  Identity tensor of order 2:

$$I_{i_1 i_2} = \delta_{i_1 i_2} = e_j \otimes e_j.$$

The identity tensor is another notation for the Kronecker tensor.

$I^{2k}$  Identity tensor of order  $2k$ :

$$I^{2k} := \overbrace{I \otimes I \otimes \dots \otimes I}^{k \text{ times}}.$$

$(e_j)_{1 \leq j \leq d}$  Vectors of the canonical basis of  $\mathbb{R}^d$ .

$e_j$  Scalar valued tensor of order 1 given by  $e_{j, i_1} := \delta_{i_1 j}$  (with  $1 \leq j \leq d$ ).

$\mathbb{B}_K^{k, l}$  Bilinear tensor of order  $k + l$  which may be either scalar or matrix valued:

- if  $\mathbb{B}_K^{k, l}$  is a *scalar* tensor, then for any smooth scalar fields  $v$  and  $w$ :

$$\mathbb{B}_K^{k, l} \nabla^k v \nabla^l w := \mathbb{B}_{K, i_1 \dots i_k, j_1 \dots j_l}^{k, l} \partial_{i_1 \dots i_k}^k v \partial_{j_1 \dots j_l}^l w.$$

- if  $\mathbb{B}_K^{k, l}$  is a *matrix* valued tensor, then for any smooth vector fields  $\mathbf{v}, \mathbf{w} \in \mathcal{C}^\infty(\mathbb{R}^d)$ :

$$\mathbb{B}_K^{k, l} \nabla^k \mathbf{v} \nabla^l \mathbf{w} := \mathbb{B}_{K, i_1 \dots i_k, j_1 \dots j_l, m}^{p, k} \partial_{i_1 \dots i_k}^p v_l \partial_{j_1 \dots j_l}^k w_m.$$

With a small abuse of notation, we consider zeroth order tensors  $b^0$  to be constants (i.e.  $b^0 \in \mathbb{R}$ ) and we still denote by  $b^0 \otimes c^k := b^0 c^k$  the tensor product with a  $k$ -th order tensor  $c^k$ . The same convention also applies to vector valued and matrix valued tensors.

Since a  $k$ -th order tensor (scalar, vector, or matrix valued)  $B^k$  makes sense when contracted with  $k$  partial derivatives as in (7.2.27) to (7.2.29), the order in which the indices  $i_1, \dots, i_k$  are written in  $B_{i_1, \dots, i_k}^k$  does not matter in general. However the ordering of spatial indices  $l, m$  of a matrix valued tensor  $B^k = (B_{i_1 \dots i_k, lm}^k)$  *does* matter.

Finally, in the whole work, we write  $C, C_K$  or  $C_K(\mathbf{f})$  to denote universal constants that do not depend on  $\varepsilon$  but whose values may change from lines to lines (and which may depend on  $\eta$ ).

**Remark 7.1.** In a limited number of places, the superscript or subscript indices  $p$  and  $q$  are used. Naturally, these are not to be confused with the pressure variables  $p_\varepsilon$  or  $q_\varepsilon$  introduced in section 7.5.



## 7.3 HIGH ORDER HOMOGENIZATION FOR THE PERFORATED POISSON PROBLEM

Our first study is concerned with the high order homogenization of the perforated Poisson problem:

$$\begin{cases} -\Delta u_\varepsilon = f \text{ in } D_\varepsilon \\ u_\varepsilon = 0 \text{ on } \partial\omega_\varepsilon \\ u_\varepsilon \text{ is } D\text{-periodic.} \end{cases} \quad (7.3.1)$$

where  $f \in C^\infty(D)$  is a smooth  $D$ -periodic right-hand side. The system (7.3.1) can be considered as a simplified scalar version of the Stokes problem (7.2.4). Our objective is to derive well-posed high order homogenized equations for (7.3.1) of arbitrary order, which reduce to the three classical regimes (see [103]) analogous to (7.2.5) to (7.2.7) in the low volume fraction limit  $\eta \rightarrow 0$ .

**Remark 7.2.** The  $D$ -periodicity assumption for  $f$  and  $u_\varepsilon$  is classical in homogenization; it is used in order to avoid difficulties related to boundary layers. There exist other settings in which we expect our analysis would also work, for instance if the domain  $D$  is smooth, if  $f \in C_c^\infty(D)$  is compactly supported in  $D$ , and if (7.3.1) is supplemented with a Dirichlet boundary condition on  $\partial D$  (see [217]).

Let us summarize the main steps of our analysis. We start by writing the traditional (see e.g. [217]) two-scale ansatz for the solution of (7.3.1):

$$u_\varepsilon(x) = \sum_{i=0}^{+\infty} \varepsilon^{i+2} u_i(x, x/\varepsilon), \quad (7.3.2)$$

where  $u_i(x, y)$  is a  $P$ -periodic function in  $y$  satisfying  $u_i(x, y) = 0$  for  $y \in \partial(\eta T)$ . Here and in all what follows, equalities involving infinite power series such as (7.3.2) are formal and without a precise meaning of convergence (which shall rather be justified in the section 7.3.4 dedicated to error estimates). We search for higher order ‘‘homogenized’’ equations for the ‘‘homogenized’’ macroscopic approximation

$$u_\varepsilon^*(x) := \sum_{i=0}^{+\infty} \varepsilon^{i+2} \int_Y u_i(x, y) dy. \quad (7.3.3)$$

Our method can then be decomposed into the following steps:

1. introducing a family of  $k$ -th order scalar tensors  $(\mathcal{X}^k(y))_{k \in \mathbb{N}}$  (vanishing on  $\eta T$  and obtained as the solutions of cell problems (see definition 7.1)), we show in proposition 7.1 that (7.3.2) reads explicitly:

$$u_\varepsilon(x) = \sum_{i=0}^{+\infty} \varepsilon^{i+2} \mathcal{X}^i(x/\varepsilon) \cdot \nabla^i f(x). \quad (7.3.4)$$

Introducing the averaged tensors of order  $i$ ,  $\mathcal{X}^{i*} := \int_Y \mathcal{X}^i(y) dy$ , (7.3.3) reads similarly

$$u_\varepsilon^*(x) = \sum_{i=0}^{+\infty} \varepsilon^{i+2} \mathcal{X}^{i*} \cdot \nabla^i f(x). \quad (7.3.5)$$

2. We construct constant tensors  $M^i$  by inversion of the formal equality

$$\left( \sum_{i=0}^{+\infty} \varepsilon^{i-2} M^i \cdot \nabla^i \right) \left( \sum_{i=0}^{+\infty} \varepsilon^{i+2} \mathcal{X}^{i*} \cdot \nabla^i \right) = I, \quad (7.3.6)$$

which yields (by left multiplication in (7.3.5)) the following formal, infinite order homogenized equation for  $u_\varepsilon^*(x)$ :

$$\sum_{i=0}^{+\infty} \varepsilon^{i-2} M^i \cdot \nabla^i u_\varepsilon^*(x) = f(x). \quad (7.3.7)$$

3. We substitute the expression for  $f(x)$  given by (7.3.7) into the ansatz (7.3.4) so as to recognize a formal double series product

$$u_\varepsilon(x) = \left( \sum_{i=0}^{+\infty} \varepsilon^i \mathcal{X}^i(x/\varepsilon) \cdot \nabla^i \right) \left( \sum_{i=0}^{+\infty} \varepsilon^i M^i \cdot \nabla^i \right) u_\varepsilon^*(x). \quad (7.3.8)$$

Introducing a new family of tensors  $N^k(y)$  defined by the Cauchy product

$$N^k(y) := \sum_{p=0}^k \mathcal{X}^p(y) \otimes M^{k-p}, \quad y \in Y,$$

we are able to express the oscillating solution  $u_\varepsilon$  in terms of its formal average  $u_\varepsilon^*$  as follows:

$$u_\varepsilon(x) = \sum_{i=0}^{+\infty} \varepsilon^i N^i(x/\varepsilon) \cdot \nabla^i u_\varepsilon^*(x). \tag{7.3.9}$$

4. In order to obtain well-posed homogenized equations of finite order, we consider truncated versions of functions of the form of (7.3.9):

$$w_\varepsilon(v)(x) := \sum_{i=0}^K \varepsilon^i N^i(x/\varepsilon) \cdot \nabla^i v(x).$$

We then formulate a minimization problem for  $v$  by restriction of the one associated with  $u_\varepsilon$ :

$$\begin{aligned} \min_{v \in H^{K+1}(D)} \int_D \left( \frac{1}{2} |\nabla w_\varepsilon(v)|^2 - f(x)w_\varepsilon(v)(x) \right) dx \\ \text{s.t. } v \text{ is } D\text{-periodic.} \end{aligned} \tag{7.3.10}$$

Averaging over the fast variable  $x/\varepsilon$  (by using lemma 7.3), we obtain a limit minimization problem involving an approximate energy  $J_K^*$  (see section 7.3.3) which does not depend on  $x/\varepsilon$ ,

$$\begin{aligned} \min_{v \in H^{K+1}(D)} J_K^*(v, f, \varepsilon) \\ \text{s.t. } v \text{ is } D\text{-periodic.} \end{aligned} \tag{7.3.11}$$

Its Euler-Lagrange equation (see definition 7.5) yields a well-posed homogenized equation of order  $2K + 2$ :

$$\sum_{k=0}^{K+1} \varepsilon^{2k-2} \mathbb{D}_K^{2k} \cdot \nabla^{2k} v_K^* = f, \tag{7.3.12}$$

where the tensors  $\mathbb{D}_K^{2k}$  are inferred from  $J_K^*$ . We then prove in proposition 7.13 that its solution  $v_K^* \in H^{K+1}(D)$  yields a high order approximation of  $u_\varepsilon$  thanks to the following error estimate:

$$\left\| u_\varepsilon - \sum_{k=0}^K \varepsilon^k N^k(\cdot/\varepsilon) \cdot \nabla^k v_K^* \right\|_{L^2(D)} \leq C_K(f) \varepsilon^{K+3}, \tag{7.3.13}$$

for a constant  $C_K(f)$  which depends only on  $K$  and  $f$ .

The most essential step of the methodology is the derivation of the non-classical ansatz (7.3.9). We shall see (proposition 7.9) that the tensor  $N^0$  is of average  $\int_Y N^0(y)dy = 1$  and  $N^k$  is of average  $\int_Y N^k(y)dy = 0$  for  $k \geq 1$ : averaging (7.3.9) with respect to the fast variable  $x/\varepsilon$  yields consistently the formal ‘‘homogenized average’’  $u_\varepsilon^*$ .

Let us stress that in the available works of the literature concerned with high order homogenization of scalar conductivity equations and its variants [53, 287, 33] (where the inhomogeneity arises through the discontinuity of the conductivity coefficient), the criminal ansatz (analogous to (7.3.9)) is readily obtained from the classical one (analogous to (7.3.2)) because the tensors  $N^k$  and  $\mathcal{X}^k$  coincide in these contexts (check for instance [53, 20]). This does not occur in our case because of the Dirichlet boundary condition on  $\partial(\eta T)$ .

The remainder of this part is organized as follows: in section 7.3.1, we examine the definition of the family of tensors  $\mathcal{X}^k$  yielding the classical ansatz (7.3.3). We state several properties of the tensor  $\mathcal{X}^k$  and review classical error estimate results for the traditional ansatz. In section 7.3.2, we detail the procedure which allows to obtain the family of constant tensors  $M^k$  and the criminal ansatz (7.3.9) involving the new tensors  $N^k(y)$ . This allows us to derive in section 7.3.3 the higher order ‘‘homogenized energy’’  $J_K^*$  and its associated high order homogenized equation (7.3.12). We prove the well-posedness of this

equation and that its first  $K + 1$  coefficients coincide with those of (7.3.7). Section 7.3.4 is dedicated to the proof of the error estimate (7.3.13) justifying the higher order homogenization process. We show in section 7.3.5 that the coefficients of the infinite order homogenized equation (7.3.7) converge to those of the original Poisson equation (7.3.1) in the low-volume fraction limit when the obstacle's size  $\eta \rightarrow 0$ , which allows us to retrieve the three classical regimes and the arising of the celebrated ‘‘strange term’’ (see [103]) for the critical size  $\eta \sim \varepsilon^{2/(d-2)}$ . Lastly, we establish in section 7.3.6 symmetry properties for the homogenized tensors  $\mathcal{X}^{k*}$  and  $M^k$  when the obstacle  $\eta T$  is invariant with respect to cell symmetries.

### 7.3.1 Formal infinite order two-scale expansions and tensors $\mathcal{X}^k$

The first step of our methodology is to insert formally the ansatz (7.3.2) into the Poisson system (7.3.1). Because it will help highlight the occurrence of double series structures, we also assume (although it is not fully necessary) that the right-hand side  $f$  depends on  $\varepsilon$  and admits the following formal expansion:

$$f(x) = \sum_{i=0}^{+\infty} \varepsilon^i f_i(x).$$

Evaluating the Laplace operator against (7.3.2) yields then formally

$$-\Delta u_\varepsilon = \sum_{i=-2}^{+\infty} \varepsilon^{i+2} (-\Delta_{yy} u_{i+2} - \Delta_{xy} u_{i+1} - \Delta_{xx} u_i)$$

where we use the convention  $u_{-2}(x, y) = u_{-1}(x, y) = 0$ , and where  $-\Delta_{yy}$ ,  $-\Delta_{xy}$ ,  $-\Delta_{xx}$  are the operators

$$-\Delta_{xx} := -\operatorname{div}_x(\nabla_x \cdot), \quad -\Delta_{xy} := -\operatorname{div}_x(\nabla_y \cdot) - \operatorname{div}_y(\nabla_x \cdot), \quad -\Delta_{yy} := -\operatorname{div}_y(\nabla_y \cdot).$$

Identifying all powers in  $\varepsilon$  yields then the traditional cascade of equations (obtained e.g. in [217]):

$$\begin{cases} -\Delta_{yy} u_{i+2} = f_{i+2} + \Delta_{xy} u_{i+1} + \Delta_{xx} u_i & \text{for all } i \geq -2 \\ u_{-2}(x, y) = u_{-1}(x, y) = 0. \end{cases} \quad (7.3.14)$$

**Definition 7.1.** We define the family of tensors  $(\mathcal{X}_{i_1 \dots i_k}^k(y))_{k \in \mathbb{N}}$  of order  $k$  by recurrence as follows:

$$\begin{cases} -\Delta_{yy} \mathcal{X}^0 = 1 \text{ in } Y \\ -\Delta_{yy} \mathcal{X}^1 = 2\partial_j \mathcal{X}^0 \otimes e_j \text{ in } Y \\ -\Delta_{yy} \mathcal{X}^{k+2} = 2\partial_j \mathcal{X}^{k+1} \otimes e_j + \mathcal{X}^k \otimes I \text{ in } Y, & \text{for all } k \geq 0 \\ \mathcal{X}^k = 0 \text{ on } \partial(\eta T) \\ \mathcal{X}^k \text{ is } P\text{-periodic.} \end{cases} \quad (7.3.15)$$

**Lemma 7.1.** *Assume that the boundary of the perforated cell  $Y$  is smooth. Then, for any  $k \in \mathbb{N}$ , the tensor  $\mathcal{X}^k$  is well-defined and is smooth, namely it holds  $\mathcal{X}^k \in \mathcal{C}^\infty(\bar{Y})$ . In particular,  $\mathcal{X}^k \in L^\infty(Y) \cap H^1(Y)$ .*

*Proof.* Since the constant function 1 is smooth, standard regularity theory for the Laplace operator  $-\Delta_{yy}$  (see [176, 76, 146]) implies  $\mathcal{X}^0 \in \mathcal{C}^\infty(\bar{Y})$ . The result follows by induction by repeating this argument to  $\mathcal{X}^1$  and  $\mathcal{X}^{k+2}$  for any  $k \geq 0$ .  $\square$

**Proposition 7.1.** *The solutions  $(u_i(x, y))_{i \geq 0}$  to the cascade of equations (7.3.14) are given by*

$$\forall i \geq 0, u_i(x, y) = \sum_{k=0}^i \mathcal{X}^k(y) \cdot \nabla^k f_{i-k}(x), \quad x \in D, y \in Y. \quad (7.3.16)$$

*Recognizing a Cauchy product, the ansatz (7.3.2) can be formally written as the following infinite power series product:*

$$u_\varepsilon(x) = \sum_{i=0}^{+\infty} \varepsilon^{i+2} \mathcal{X}^i(x/\varepsilon) \cdot \nabla^i f(x) = \varepsilon^2 \left( \sum_{i=0}^{+\infty} \varepsilon^i \mathcal{X}^i(x/\varepsilon) \cdot \nabla^i \right) \left( \sum_{i=0}^{+\infty} \varepsilon^i f_i(x) \right). \quad (7.3.17)$$

*Proof.* We proceed by induction. The cases  $i = 0$  and  $i = 1$  are simple consequences of (7.3.14). If the result holds until rank  $i + 1$ , then

$$\begin{aligned} -\Delta_{yy}u_{i+2}(x, y) &= f_{i+2}(x) + \sum_{k=0}^{i+1} 2\partial_j \mathcal{X}^k(y) \cdot \nabla^k (\partial_j f_{i+1-k}(x)) + \sum_{k=0}^i \mathcal{X}^k \cdot \nabla^k (\Delta f_{i-k}(x)) \\ &= f_{i+2}(x) + \sum_{k=-1}^i 2(\partial_j \mathcal{X}^{k+1}(y) \otimes e_j) \cdot \nabla^{k+2} f_{i-k}(x) + \sum_{k=0}^i (\mathcal{X}^k(y) \otimes I) \cdot \nabla^{k+2} f_{i-k}(x) \\ &= f_{i+2}(x) + 2\partial_j \mathcal{X}^0(y) \partial_j f_{i+1}(x) + \sum_{k=0}^i (-\Delta_{yy} \mathcal{X}^{k+2}(y)) \cdot \nabla^{k+2} f_{i-k}(x). \end{aligned}$$

Whence

$$u_{i+2}(x, y) = \mathcal{X}^0(y) f_{i+2}(x) + \mathcal{X}_j^1(y) \partial_j f_{i+1}(x) + \sum_{k=0}^i \mathcal{X}^{k+2}(y) \cdot \nabla^{k+2} f_{i-k}(x) = \sum_{k=0}^{i+2} \mathcal{X}^k(y) \cdot \nabla^k f_{i+2-k}(x).$$

□

In what follows, we are going to derive high order PDEs for the “infinite order” homogenized solution (7.3.3). Following the notation conventions introduced in section 7.2.4, we find convenient to introduce averaged functions  $u_i^*$  and tensors  $\mathcal{X}^{i*}$  labelled with the star “\*” notation:

**Definition 7.2.** For any  $i \in \mathbb{N}$ , the averaged functions and tensors with respect to the fast variable are denoted by:

$$u_i^*(x) := \int_Y u_i(x, y) dy. \quad (7.3.18)$$

$$\mathcal{X}^{i*} := \int_Y \mathcal{X}^i(y) dy. \quad (7.3.19)$$

**Remark 7.3.** These are average on the unit cell  $P = [0, 1]^d$ , although the integral is written on the domain  $Y$ , because  $u_i(x, \cdot)$  and  $\mathcal{X}^i$  vanish on  $P \setminus Y = \overline{\eta T}$  (see Figure 7.3).

In the next proposition we show that  $\mathcal{X}^{2p*}$  depends only on the lower order tensors  $\mathcal{X}^p$  and  $\mathcal{X}^{p-1}$  and that  $\mathcal{X}^{2p+1*}$  is null for any  $p \in \mathbb{N}$ . Similar formulas have been obtained for the wave equation in heterogeneous media, see e.g. Theorem 3.5 in [4], and also [3, 264].

**Proposition 7.2.** For any  $0 \leq p \leq k$ , the following identity holds for the tensor  $\mathcal{X}^{k*}$ :

$$\mathcal{X}^{k*} = \int_Y \mathcal{X}^k dy = (-1)^p \int_Y (\mathcal{X}^{k-p} \otimes (-\Delta_{yy} \mathcal{X}^p) - \mathcal{X}^{k-p-1} \otimes \mathcal{X}^{p-1} \otimes I) dy, \quad (7.3.20)$$

with the convention that  $\mathcal{X}^{-1} = 0$ . In particular, for any  $p \in \mathbb{N}$ :

- $\mathcal{X}^{2p+1*} = 0$
- $\mathcal{X}^{2p*}$  depends only on the tensors  $\mathcal{X}^p$  and  $\mathcal{X}^{p-1}$ :

$$\mathcal{X}^{2p*} = (-1)^p \int_Y (\partial_j \mathcal{X}^p \otimes \partial_j \mathcal{X}^p - \mathcal{X}^{p-1} \otimes \mathcal{X}^{p-1} \otimes I) dy. \quad (7.3.21)$$

*Proof.* We proceed again by induction. Formula (7.3.20) holds true for  $p = 0$  by using the convention  $\mathcal{X}^{-1} = 0$  and  $-\Delta_{yy} \mathcal{X}^0 = 1$ . Assuming now the result true for  $p < k$ , we may perform the following integration by parts thanks to the boundary conditions satisfied by the tensors  $\mathcal{X}^k$ :

$$\begin{aligned} \mathcal{X}^{k*} &= (-1)^p \int_Y (\mathcal{X}^{k-p} \otimes (-\Delta_{yy} \mathcal{X}^p) - \mathcal{X}^{k-p-1} \otimes \mathcal{X}^{p-1} \otimes I) dy \\ &= (-1)^p \int_Y (-\Delta_{yy} \mathcal{X}^{k-p} \otimes \mathcal{X}^p - \mathcal{X}^{k-p-1} \otimes \mathcal{X}^{p-1} \otimes I) dy \\ &= (-1)^p \int_Y ((2\partial_j \mathcal{X}^{k-p-1} \otimes e_j + \mathcal{X}^{k-p-2} \otimes I) \otimes \mathcal{X}^p - \mathcal{X}^{k-p-1} \otimes \mathcal{X}^{p-1} \otimes I) dy \\ &= (-1)^p \int_Y (-2\partial_j \mathcal{X}^p \otimes e_j - \mathcal{X}^{p-1} \otimes I) \otimes \mathcal{X}^{k-p-1} + \mathcal{X}^{k-p-2} \otimes \mathcal{X}^p \otimes I) dy \\ &= (-1)^{p+1} \int_Y ((-\Delta_{yy} \mathcal{X}^{p+1}) \otimes \mathcal{X}^{k-p-1} - \mathcal{X}^{k-p-2} \otimes \mathcal{X}^p \otimes I) dy. \end{aligned}$$

Hence the formula is proved at order  $p + 1$ .  
Now, the formula at order  $p = k$  reads

$$\mathcal{X}^{k*} = - \int_Y \mathcal{X}^0(-\Delta_{yy}\mathcal{X}^k)dy = - \int_Y \mathcal{X}^k(-\Delta_{yy}\mathcal{X}^0)dy = -\mathcal{X}^{k*},$$

which implies  $\mathcal{X}^{k*} = 0$  if  $k$  is odd. Formula (7.3.21) follows easily from (7.3.20) with  $k = 2p$ .  $\square$

The following result demonstrates that  $\mathcal{X}^k(y)$  is not identically equal to zero (although some components  $\mathcal{X}_{i_1 \dots i_k}^k(y)$  could vanish for some set of indices  $i_1, \dots, i_k$ , e.g. in case of symmetries of the obstacle  $\eta T$ ).

**Proposition 7.3.** *The following identity holds:*

$$-\Delta_{yy}(\partial_{i_1 \dots i_k}^k \mathcal{X}_{i_1 \dots i_k}^k) = (-1)^k(k+1). \quad (7.3.22)$$

*Proof.* The results naturally holds true for  $k = 0$ . For  $k = 1$ , it holds

$$-\Delta_{yy}\partial_i \mathcal{X}_i^1 = \partial_i(2\partial_i \mathcal{X}^0) = 2\Delta \mathcal{X}^0 = -2.$$

Assuming the result till rank  $k - 1$ , the formula still holds at rank  $k$  because

$$\begin{aligned} -\Delta_{yy}\partial_{i_1 \dots i_k}^k \mathcal{X}_{i_1 \dots i_k}^k &= \partial_{i_1 \dots i_k}^k (2\partial_{i_k} \mathcal{X}_{i_1 \dots i_{k-1}}^{k-1} + \mathcal{X}_{i_1 \dots i_{k-2}}^{k-2} \delta_{i_{k-1} i_k}) \\ &= 2\Delta_{yy}(\partial_{i_1 \dots i_{k-1}}^{k-1} \mathcal{X}_{i_1 \dots i_{k-1}}^{k-1}) + \Delta_{yy}(\partial_{i_1 \dots i_{k-2}}^{k-2} \mathcal{X}_{i_1 \dots i_{k-2}}^{k-2}) \\ &= -2(-1)^{k-1}k - (-1)^{k-2}(k-1) \\ &= (-1)^k(k+1). \end{aligned}$$

$\square$

The next result can be found in classical textbooks, see e.g. [217]. It states error estimates for the truncated classical ansatz (7.5.2):

**Proposition 7.4.** *Denote  $u_{\varepsilon, K}$  the truncated ansatz of (7.3.17) at order  $K$ :*

$$\forall x \in D, u_{\varepsilon, K}(x) := \sum_{i=0}^K \varepsilon^{i+2} \mathcal{X}^i(x/\varepsilon) \cdot \nabla^i f(x). \quad (7.3.23)$$

*Then assuming  $f$  is  $D$ -periodic, the following bound holds:*

$$\|u_\varepsilon - u_{\varepsilon, K}\|_{H^1(D_\varepsilon)} \leq C_K \varepsilon^{K+2} \|f\|_{H^{K+2}(D)} \quad (7.3.24)$$

*for a constant  $C_K$  independent of  $f$  and  $\varepsilon$  (but depending on  $K$ ).*

*Proof.* It is sufficient to observe that the remainder  $r_\varepsilon(x) = u_\varepsilon(x) - u_{\varepsilon, K}(x)$  satisfies

$$-\Delta r_\varepsilon = (\varepsilon^{K+1}(\Delta_{xy}u_K + \Delta_{xx}u_{K-1}) + \varepsilon^{K+2}\Delta_{xx}u_K)(x, x/\varepsilon).$$

and the terms  $(\Delta_{xy}u_K)(x, x/\varepsilon)$ ,  $(\Delta_{xx}u_{K-1})(x, x/\varepsilon)$  and  $(\Delta_{xx}u_K)(x, x/\varepsilon)$  can be bounded in the  $L^2(D)$  norm by  $C_K \|f\|_{H^{K+2}}$ . The result follows from

$$\|\nabla r_\varepsilon\|_{L^2(D_\varepsilon, \mathbb{R}^d)}^2 \leq C_K \varepsilon^{K+1} \|f\|_{H^{K+2}(D)} \|r_\varepsilon\|_{L^2(D_\varepsilon)} \leq C_K \varepsilon^{K+2} \|f\|_{H^{K+2}(D)} \|\nabla r_\varepsilon\|_{L^2(D_\varepsilon, \mathbb{R}^d)},$$

where we have used the classical Poincaré inequality of lemma 7.2 below.  $\square$

**Lemma 7.2** (see e.g. Lions 1981 [217]). *There exists a constant  $C$  independent of  $\varepsilon$  such that for any  $\phi \in H^1(D_\varepsilon)$  satisfying  $\phi = 0$  on the boundary  $\partial\omega_\varepsilon$  of the holes, the following Poincaré inequality holds:*

$$\|\phi\|_{L^2(D_\varepsilon)} \leq C\varepsilon \|\nabla \phi\|_{L^2(D_\varepsilon, \mathbb{R}^d)}.$$

In the final part of this section, we provide a few insights highlighting in which sense the homogenized average  $u_\varepsilon^*$  of (7.3.3) is an approximation of  $u_\varepsilon$ . Let us recall the following important lemma (see e.g. Appendix C of [287] or [16, 138]):

**Lemma 7.3.** *Let  $\phi$  a  $P = [0, 1]^d$ -periodic function and  $f \in C^\infty(D^d)$  a smooth  $D$ -periodic function. Then for any  $k \in \mathbb{N}$  arbitrarily high, the following inequality holds:*

$$\left| \int_D f(x)\phi(x/\varepsilon)dx - \int_D \int_P f(x)\phi(y)dydx \right| \leq \frac{(2L)^{d/2}}{|2\pi|^k} \left\| \phi - \int_P \phi dy \right\|_{L^2(P)} \|f\|_{H^k(D)} \varepsilon^k. \quad (7.3.25)$$

*Proof.* The key idea is to perform integration by parts by differentiating  $f$  and integrating  $\phi$ . Without loss of generality, we may assume  $\phi$  to be of average zero, i.e.  $\int_P \phi dy = 0$ . We show by induction that for any  $k \geq 0$ , there exists a function  $\phi_{i_1 \dots i_k}^k$  with  $1 \leq i_1 \dots i_k \leq d$  satisfying

$$\sum_{1 \leq i_1, \dots, i_k \leq d} \|\phi_{i_1 \dots i_k}^k\|_{L^2(P)}^2 \leq \frac{1}{|4\pi^2|^k} \|\phi\|_{L^2(P)}^2, \quad (7.3.26)$$

$$\int_P \phi_{i_1 \dots i_k}^k dy = 0, \quad (7.3.27)$$

and such that

$$\int_D f(x)\phi(x/\varepsilon)dy = (-1)^k \int_D \varepsilon^k \partial_{i_1 \dots i_k}^k f \phi_{i_1 \dots i_k}^k(x/\varepsilon)dx. \quad (7.3.28)$$

Obviously, the above identities hold true for  $k = 0$ . Assuming now that the results holds till rank  $k \geq 0$ , we introduce  $\psi_{i_1 \dots i_k}^k \in H^1(P)$  the unique solution to the following Laplace problem:

$$\begin{cases} -\Delta \psi_{i_1 \dots i_k}^k = \phi_{i_1 \dots i_k}^k & \text{in } P \\ \psi_{i_1 \dots i_k}^k & \text{is } P\text{-periodic} \\ \int_P \psi_{i_1 \dots i_k}^k dy = 0. \end{cases} \quad (7.3.29)$$

Such function actually exists because of the compatibility condition (7.3.27); this allows us to define

$$\phi_{i_1 \dots i_{k+1}}^{k+1} := \partial_{i_{k+1}} \psi_{i_1 \dots i_k}^k.$$

It is clear from the periodicity of  $\psi_{i_1 \dots i_k}^k$  that (7.3.27) holds at rank  $k+1$ , and it holds from the variational formulation of (7.3.29)

$$\sum_{i_{k+1}=1}^d \|\phi_{i_1 \dots i_{k+1}}^{k+1}\|_{L^2(P)}^2 = \|\nabla \psi_{i_1 \dots i_k}^k\|_{L^2(P, \mathbb{R}^d)}^2 \leq \frac{1}{4\pi^2} \|\phi_{i_1 \dots i_k}^k\|_{L^2(P)}^2,$$

because  $1/(4\pi^2)$  is the Poincaré constant of  $[0, 1]^d$ . Summing over the indices  $i_1 \dots i_k$  and applying hypothesis (7.3.26) at rank  $k$  yields the same estimate (7.3.26) at rank  $k+1$ . We may subsequently write

$$\begin{aligned} \int_D \partial_{i_1 \dots i_k}^k f \phi_{i_1 \dots i_k}^k(x/\varepsilon)dx &= \int_D \varepsilon \partial_{i_1 \dots i_k}^k f \operatorname{div}(\phi_{i_1 \dots i_{k+1}}^{k+1}(x/\varepsilon) \mathbf{e}_{i_{k+1}})dx \\ &= - \int_D \varepsilon \partial_{i_1 \dots i_{k+1}}^{k+1} f(x) \phi_{i_1 \dots i_{k+1}}^{k+1}(x/\varepsilon)dx, \end{aligned}$$

from where (7.3.28) follows at rank  $k+1$ .

Let us finally prove (7.3.25). From (7.3.28) and the Cauchy-Schwartz inequality, we obtain

$$\left| \int_D f(x)\phi(x/\varepsilon)dx \right|^2 \leq \varepsilon^{2k} \|f\|_{H^k(D)}^2 \sum_{1 \leq i_1 \dots i_k} \int_D |\phi_{i_1 \dots i_k}^k(x/\varepsilon)|^2 dx \leq \varepsilon^{2k} \|f\|_{H^k(D)}^2 \frac{(2L)^d}{|4\pi^2|^k} \|\phi\|_{L^2(P)}^2,$$

whence the result. □

**Remark 7.4.** The constant  $2\pi$  in (7.3.25) is related to the fact that the unit cell is of size 1 (it would be one for a cell of length  $2\pi$ ), it is therefore no surprise that it scales with the same power  $k$  than  $\varepsilon$ .

This lemma implies that  $u_\varepsilon^*$  is an approximation of  $u_\varepsilon$  in a distributional sense:

**Proposition 7.5.** Denote by  $u_{\varepsilon,K}^*$  the average of the truncated ansatz  $u_{\varepsilon,K}$  (eqn. (7.3.23)) with respect to the fast variable, namely

$$u_{\varepsilon,K}^*(x) := \sum_{i=0}^K \varepsilon^{i+2} u_i^*(x) = \sum_{i=0}^K \varepsilon^{i+2} \mathcal{X}^{i*} \cdot \nabla^i f(x). \quad (7.3.30)$$

Then  $u_{\varepsilon,K}^*$  is a  $(K+3)$ -th order approximation of  $u_\varepsilon$  in the following weak sense:

$$\forall \psi \in C_c^\infty(D), \left| \int_D u_\varepsilon(x) \psi(x) dx - \int_D u_{\varepsilon,K}^*(x) \psi(x) dx \right| \leq C_K(\psi) \|f\|_{H^{K+2}(D)} \varepsilon^{K+3},$$

for a constant  $C_K(\psi)$  depending only on  $K$  and  $\psi$ .

*Proof.* We write

$$\begin{aligned} & \left| \int_D u_\varepsilon(x) \psi(x) dx - \int_D u_{\varepsilon,K}^*(x) \psi(x) dx \right| \\ & \leq \int_D |(u_\varepsilon(x) - u_{\varepsilon,K}(x)) \psi(x)| dx + \left| \int_D (u_{\varepsilon,K}(x) - u_{\varepsilon,K}^*(x)) \psi(x) dx \right| \end{aligned} \quad (7.3.31)$$

Let  $M_K(\psi) > 0$  a constant such that  $\|\nabla^k \psi\|_{L^\infty(D, \mathbb{R}^d)} \leq M_K(\psi)$  for any  $0 \leq k \leq K+1$ . The first term is bounded by  $\|u_\varepsilon - u_{\varepsilon,K}\|_{L^2(D)} \|\psi\|_{L^2(D)} \leq C_K M_K(\psi) \varepsilon^{K+3} \|f\|_{H^{K+2}(D)}$  according to proposition 7.4. The second term is bounded by  $C_K M_K(\psi) \|f\|_{H^{K+1}(D)} \varepsilon^{K+3}$  (up to the use of a larger constant  $C_K$ ) as a consequence of the important Lemma 7.3 stated above.  $\square$

**Remark 7.5.** More elaborate arguments can be proposed for providing  $u_\varepsilon^*(x)$  with a physical interpretation, e.g. by considering shifted cell averages (see [287, 102]).

### 7.3.2 Formal infinite order homogenized equation and criminal ansatz: tensors $M^k$ and $N^k$

We now detail the steps (2) and (3) of the procedure outlined in the introduction of this section. Let us recall that the first tensor  $\mathcal{X}^{0*}$  is a strictly positive number, since (7.3.21) implies  $\mathcal{X}^{0*} = \int_Y |\nabla \mathcal{X}^0|^2 dy$ .

**Proposition 7.6.** Let  $(M^k)_{i \in \mathbb{N}}$  be the family of  $k$ -th order tensors defined by induction as follows:

$$\begin{cases} M^0 = (\mathcal{X}^{0*})^{-1}, \\ M^k = -(\mathcal{X}^{0*})^{-1} \sum_{p=0}^{k-1} \mathcal{X}^{k-p*} \otimes M^p. \end{cases} \quad (7.3.32)$$

Then it holds, given the definitions (7.3.14) and (7.3.18) for  $u_i^*$ :

$$\forall i \geq 0, f_i(x) = \sum_{k=0}^i M^k \cdot \nabla^k u_{i-k}^*(x). \quad (7.3.33)$$

Recognizing a Cauchy product, this can be rewritten formally in terms of the following ‘infinite order’ homogenized equation for  $u_\varepsilon^*$  (defined in (7.3.5)):

$$\sum_{i=0}^{+\infty} \varepsilon^{i-2} M^i \cdot \nabla^i u_\varepsilon^* = f. \quad (7.3.34)$$

*Proof.* We proceed by induction. The case  $i=0$  results from the identity  $u_0^*(x) = \mathcal{X}^{0*} f_0(x)$  which yields  $f_0(x) = (\mathcal{X}^{0*})^{-1} u_0^*(x)$ . Then, assuming the result (7.3.33) holds till rank  $i-1$ , we average equation (7.3.16) with respect to the  $y$  variable to obtain

$$u_i^* = \sum_{p=0}^i \mathcal{X}^{p*} \cdot \nabla^p f_{i-p} = \mathcal{X}^{0*} f_i + \sum_{p=1}^i \mathcal{X}^{p*} \cdot \nabla^p f_{i-p}.$$

By inversion of the nonzero coefficient  $\mathcal{X}^{0*}$ , we obtain the following expression for  $f_i$ :

$$\begin{aligned}
 f_i &= (\mathcal{X}^{0*})^{-1} \left( u_i^* - \sum_{p=1}^i \sum_{q=0}^{i-p} (\mathcal{X}^{p*} \otimes M^q) \cdot \nabla^{p+q} u_{i-p-q}^* \right) \\
 &= (\mathcal{X}^{0*})^{-1} \left( u_i^* - \sum_{p=1}^i \sum_{k=p}^i (\mathcal{X}^{p*} \otimes M^{k-p}) \cdot \nabla^k u_{i-k}^* \right) \quad (\text{change of indices } k = p + q) \\
 &= (\mathcal{X}^{0*})^{-1} \left( u_i^* - \sum_{k=1}^i \sum_{p=1}^k (\mathcal{X}^{p*} \otimes M^{k-p}) \cdot \nabla^k u_{i-k}^* \right) \quad (\text{inversion of summation}) \\
 &= (\mathcal{X}^{0*})^{-1} \left( u_i^* - \sum_{k=1}^i \left( \sum_{p=0}^{k-1} \mathcal{X}^{k-p*} \otimes M^p \right) \cdot \nabla^k u_{i-k}^* \right) \quad (\text{change of index } p \leftrightarrow k - p) \\
 &= M^0 u_i^* + \sum_{k=1}^i M^k \cdot \nabla^k u_{i-k}^*,
 \end{aligned}$$

which concludes the proof.  $\square$

**Corollary 7.1.**  $M^k = 0$  for any odd value of  $k$ .

*Proof.* This follows from [proposition 7.2](#) and the recurrence formula [\(7.3.32\)](#).  $\square$

It is possible to write a more explicit formula for the tensors  $M^k$ :

**Proposition 7.7.** The tensors  $M^k$  are explicitly given by  $M^0 = (\mathcal{X}^{0*})^{-1}$  and:

$$\forall k \geq 1, M^k = \sum_{p=1}^k \frac{(-1)^p}{(\mathcal{X}^{0*})^{p+1}} \sum_{\substack{i_1 + \dots + i_p = k \\ 1 \leq i_1 \dots i_p \leq k}} \mathcal{X}^{i_1*} \otimes \dots \otimes \mathcal{X}^{i_p*}. \quad (7.3.35)$$

*Proof.* We prove it by induction. For  $k = 1$ , the result is true because

$$M^1 = -(\mathcal{X}^{0*})^{-1} M^0 \mathcal{X}^{1*} = -(\mathcal{X}^{0*})^{-2} \mathcal{X}^{1*}$$

which is exactly [\(7.3.35\)](#). Assuming [\(7.3.35\)](#) holds till rank  $k \geq 1$ , we now compute

$$\begin{aligned}
 M^{k+1} &= -(\mathcal{X}^{0*})^{-1} \sum_{p=0}^k \mathcal{X}^{k+1-p*} \otimes M^p \\
 &= -(\mathcal{X}^{0*})^{-1} M^0 \mathcal{X}^{k+1*} - (\mathcal{X}^{0*})^{-1} \sum_{p=1}^k \sum_{q=1}^p \frac{(-1)^q}{(\mathcal{X}^{0*})^{q+1}} \mathcal{X}^{k+1-p*} \otimes \sum_{\substack{i_1 + \dots + i_q = p \\ 1 \leq i_1 \dots i_q \leq p}} \mathcal{X}^{i_1*} \otimes \dots \otimes \mathcal{X}^{i_q*} \\
 &= -(\mathcal{X}^{0*})^{-2} \mathcal{X}^{k+1*} - (\mathcal{X}^{0*})^{-1} \sum_{q=1}^k \frac{(-1)^q}{(\mathcal{X}^{0*})^{q+1}} \sum_{p=q}^k \sum_{\substack{i_1 + \dots + i_q = p \\ 1 \leq i_1 \dots i_q \leq p}} \mathcal{X}^{k+1-p*} \otimes \mathcal{X}^{i_1*} \otimes \dots \otimes \mathcal{X}^{i_q*} \\
 &= -(\mathcal{X}^{0*})^{-2} \mathcal{X}^{k+1*} - (\mathcal{X}^{0*})^{-1} \sum_{q=1}^k \frac{(-1)^q}{(\mathcal{X}^{0*})^{q+1}} \sum_{\substack{i_1 + \dots + i_{q+1} = k+1 \\ 1 \leq i_1 \dots i_{q+1} \leq k+1}} \mathcal{X}^{i_{q+1}*} \otimes \mathcal{X}^{i_1*} \otimes \dots \otimes \mathcal{X}^{i_q*} \\
 &= -(\mathcal{X}^{0*})^{-2} \mathcal{X}^{k+1*} + \sum_{q=2}^{k+1} \frac{(-1)^q}{(\mathcal{X}^{0*})^{q+1}} \sum_{\substack{i_1 + \dots + i_q = k+1 \\ 1 \leq i_1 \dots i_q \leq k+1}} \mathcal{X}^{i_1*} \otimes \dots \otimes \mathcal{X}^{i_q*},
 \end{aligned}$$

from where the result follows.  $\square$

**Remark 7.6.** This result essentially states that  $\sum_{k=0}^{+\infty} \varepsilon^k M^k \cdot \nabla^k$  is the formal series expansion of

$$\left( \sum_{k=0}^{+\infty} \varepsilon^k \mathcal{X}^k \cdot \nabla^k \right)^{-1}.$$



Indeed, it is elementary to show the following identity for the inverse of a power series  $\sum_{k=0}^{+\infty} a_k z^k$  with  $(a_k) \in \mathbb{C}^N$ ,  $z \in \mathbb{C}$  and radius of convergence  $R > 0$ :

$$\left( \sum_{k=0}^{+\infty} a_k z^k \right)^{-1} = a_0^{-1} + \sum_{k=1}^{+\infty} \left( \sum_{p=1}^k \frac{(-1)^p}{a_0^{p+1}} \sum_{\substack{i_1+\dots+i_p=k \\ 1 \leq i_1 \dots i_p \leq k}} a_{i_1} a_{i_2} \dots a_{i_p} \right) z^k. \quad (7.3.36)$$

The infinite order homogenized equation (7.3.34) is not fully satisfactory because it is not immediately clear how to truncate the operator  $\sum_{i=0}^{+\infty} \varepsilon^i M^i \cdot \nabla^i$  so as to obtain a well-posed problem. Several techniques have been proposed in the literature to obtain well posed equations in the context of the homogenization of the conductivity or wave equation; we can mention among them a Boussinesq trick [35, 22, 3, 4] or a filtering method [33]. Here, we propose to follow a “variational” method inspired by [287], which is based on the existence of a “criminal” ansatz of the form (7.3.9). This ansatz is obtained by writing the oscillatory part  $u_i(x, y)$  in terms of the non oscillatory part  $u_i^*(x)$ , which is obtained from the formal equality (7.3.8):

**Proposition 7.8.** *Given the previous definitions of  $\mathcal{X}^i, M^i, u_i^*$ , the following equality holds:*

$$\forall i \geq 0, u_i(x, y) = \sum_{k=0}^i \left( \sum_{p=0}^k M^p \otimes \mathcal{X}^{k-p}(y) \right) \cdot \nabla^k u_{i-k}^*(x). \quad (7.3.37)$$

*Proof.* It suffices to substitute (7.3.33) into (7.3.16), which yields

$$\begin{aligned} u_i(x, y) &= \sum_{p=0}^i \sum_{q=0}^{i-p} (\mathcal{X}^p(y) \otimes M^q) \cdot \nabla^{p+q} u_{i-p-q}^*(x) \\ &= \sum_{p=0}^i \sum_{k=p}^i (\mathcal{X}^p(y) \otimes M^{p-k}) \cdot \nabla^k u_{i-k}^*(x) \quad (\text{change of indices } k = p + q) \\ &= \sum_{k=0}^i \sum_{p=0}^k (\mathcal{X}^p(y) \otimes M^{p-k}) \cdot \nabla^k u_{i-k}^*(x) \quad (\text{intersion of summation}) \end{aligned} \quad (7.3.38)$$

The result follows by performing a last change of indices  $p \rightarrow k - p$ .  $\square$

We are now able to introduce the tensors  $N^i$  of (7.3.9) and their related differential operators:

**Definition 7.3.** For any  $k \geq 0$ , we denote by  $N^k(y)$  the  $k$ -th order tensor

$$N^k(y) := \sum_{p=0}^k M^p \otimes \mathcal{X}^{k-p}(y). \quad (7.3.39)$$

Recognizing a Cauchy product, the identity (7.3.37) can be formally written in terms of the following “criminal” ansatz which expresses the oscillating solution  $u_\varepsilon$  in terms of its formal homogenized averaged  $u_\varepsilon^*$  (defined in (7.3.5)):

$$u_\varepsilon(x) = \sum_{i=0}^{\infty} \varepsilon^i N^i(x/\varepsilon) \cdot \nabla^i u_\varepsilon^*(x), \quad x \in D_\varepsilon. \quad (7.3.40)$$

The last proposition of this section gathers several important properties for the tensors  $N^k$  that are dual to those of the tensors  $\mathcal{X}^k$ .

**Proposition 7.9.** *The tensor  $N^k(y)$  satisfies:*

1.  $\int_Y N^0(y) dy = 1$  and  $\int_Y N^k(y) dy = 0$  for any  $k \geq 1$ ,
2. For any  $k \geq 0$ , using the convention  $N^{-1} = N^{-2} = 0$ :

$$\begin{cases} -\Delta_{yy} N^0 = M^0 \\ -\Delta_{yy} N^1 = 2\partial_j N^0 \otimes e_j + M^1 \\ -\Delta_{yy} N^{k+2} = 2\partial_j N^{k+1} \otimes e_j + N^k \otimes I + M^{k+2}. \end{cases} \quad (7.3.41)$$

3. For any  $k \geq 0$ ,

$$-\Delta_{yy}(\partial_{i_1 \dots i_k}^k N_{i_1 \dots i_k}^k) = (-1)^k (k+1) M^0. \quad (7.3.42)$$

4. For any  $1 \leq p \leq k-1$ ,

$$M^k = (-1)^{p+1} \int_Y (N^{k-p} \otimes (-\Delta_{yy} N^p) - N^{k-p-1} \otimes N^{p-1} \otimes I) dy, \quad (7.3.43)$$

In particular,  $M^{2p}$  depends only on the tensors  $N^p$  and  $N^{p-1}$ , which depend themselves only on the first  $p+1$  tensors  $\mathcal{X}^0 \dots \mathcal{X}^p$ .

*Proof.* 1. is a consequence of the definition (7.3.32) for the tensor  $M^k$  which can be rewritten as

$$\forall k \geq 1, \int_Y N^k(y) dy = \sum_{p=0}^k \mathcal{X}^{k-p*} \otimes M^p = 0,$$

and for  $k=0$ , it holds  $\int_Y N^0(y) dy = M^0 \mathcal{X}^{0*} = 1$ .

2. The first two equalities of (7.3.41) are easily verified. The third line is obtained by writing

$$\begin{aligned} -\Delta_{yy} N^k &= \sum_{p=0}^k M^{k-p} \otimes (-\Delta_{yy} \mathcal{X}^p) \\ &= \sum_{p=2}^k M^{k-p} \otimes (2\partial_j \mathcal{X}^{p-1} \otimes e_j + \mathcal{X}^{p-2} \otimes I) + M^{k-1} \otimes (2\partial_j \mathcal{X}^0 \otimes e_j) + M^k \\ &= 2\partial_j \left( \sum_{p=1}^k M^{k-p} \otimes \mathcal{X}^{p-1} \right) \otimes e_j + \left( \sum_{p=0}^{k-2} M^p \otimes \mathcal{X}^{k-p-2} \right) \otimes I + M^k, \end{aligned}$$

from where the result follows.

3. The proof of (7.3.42) is identical to that of proposition 7.3.

4. We start by proving the result for  $p=1$  with  $k > 1$ : using the point 1., we may write

$$\begin{aligned} M^k &= \int_Y N^0 \otimes M^k dy = \int_Y N^0 \otimes (-\Delta N^k - 2\partial_j N^{k-1} \otimes e_j - N^{k-2} \otimes I) dy \\ &= \int_Y M^0 \otimes N^k dy + \int_Y (2\partial_j N^0 \otimes e_j \otimes N^{k-1} - N^0 \otimes N^{k-2} \otimes I) dy \\ &= \int_Y ((-\Delta_{yy} N^1) \otimes N^{k-1} - N^0 \otimes N^{k-2} \otimes I) dy. \end{aligned}$$

If the result now holds until rank  $p$  with  $1 \leq p \leq k-2$ , we obtain the result at rank  $p+1$  with analogous computations:

$$\begin{aligned} M^k &= (-1)^{p+1} \int_Y (M^{k-p} + 2\partial_j N^{k-p-1} \otimes e_j + N^{k-p-2} \otimes I) \otimes N^p - N^{k-p-1} \otimes N^{p-1} \otimes I) dy \\ &= (-1)^{p+1} \int_Y ((-2\partial_j N^p \otimes e_j - N^{p-1} \otimes I) \otimes N^{k-p-1} + N^{k-p-2} \otimes N^p \otimes I) dy \\ &= (-1)^{p+1} \int_Y ((\Delta_{yy} N^{p+1} + M^{p+1}) \otimes N^{k-p-1} + N^{k-p-2} \otimes N^p \otimes I) dy. \end{aligned}$$

□

### 7.3.3 Homogenized equations of order $2K+2$ : tensors $\mathbb{B}_K$ and $\mathbb{D}_K$

We now detail the step (4) of the introduction of this section which allows to obtain a well-posed variational problem of order  $2K+2$  by truncating the infinite homogenized equation (7.3.7). Error estimates are then proved in section 7.3.4.

For any  $u \in H^1(D)$  and  $f \in L^2(D)$ , we denote by  $J(u, f)$  the energy

$$J(u, f) := \int_D \left( \frac{1}{2} |\nabla u|^2 - fu \right) dx.$$

It is well known that  $u_\varepsilon$  (identified with its extension by 0 in  $D \setminus D_\varepsilon$ ) is the solution to the following minimization problem:

$$u_\varepsilon = \arg \min_{w \in H^1(D_\varepsilon)} J(w, f) \quad \text{s.t.} \quad \begin{cases} w = 0 \text{ on } \partial\omega_\varepsilon \\ w \text{ is } D\text{-periodic.} \end{cases}$$

In this part, we follow the method of [53, 287] in order to obtain a minimization problem for an approximation of the formal homogenized average  $u_\varepsilon^*$ . The main idea to obtain well-posed homogenized equations of finite order is to consider truncations  $w_{\varepsilon, K}$  of the ‘‘criminal’’ ansatz (7.3.9) of the form:

$$w_{\varepsilon, K}(v)(x) := \sum_{k=0}^K \varepsilon^k N^k(x/\varepsilon) \cdot \nabla^k v(x), \quad x \in D_\varepsilon, \quad (7.3.44)$$

where the function  $v$  does not depend on the fast variable  $x/\varepsilon$  and is sought to approximate the formal homogenized average  $u_\varepsilon^*$ . Note that  $N^k(y)$  is extended by 0 in  $D \setminus D_\varepsilon$  for (7.3.44) to make sense. We consider the following approximate minimization problem for the function  $v$ :

$$\min_{v \in H^{K+1}(D)} J(w_\varepsilon(v), f) \quad \text{s.t.} \quad v \text{ is } D\text{-periodic.} \quad (7.3.45)$$

The next step is to eliminate the fast variable  $x/\varepsilon$  in (7.3.45) (by using lemma 7.3) and to read a higher order homogenized equation from the first order optimality condition. In order to do so, we introduce several additional tensors that arise in the averaging process.

**Definition 7.4** (Tensors  $\tilde{N}_j^k$  and  $\mathbb{B}_K^{l, m}$ ). For any  $K \in \mathbb{N}$ ,  $1 \leq j \leq d$  and  $0 \leq k \leq K + 1$ , let  $\tilde{N}_j^k(y)$  (with implicit dependence with respect to  $K$ ) be the  $k$ -th order tensor defined by

$$\tilde{N}_j^k(y) = \begin{cases} \partial_j N^0(y) & \text{if } k = 0 \\ \partial_j N^k(y) + N^{k-1}(y) \otimes e_j & \text{if } 1 \leq k \leq K \\ N^K(y) \otimes e_j & \text{if } k = K + 1. \end{cases} \quad (7.3.46)$$

We define a family of constant bilinear tensors  $\mathbb{B}_K^{l, m}$  of order  $l + m$  by the formula

$$\mathbb{B}_K^{l, m} := \int_Y \tilde{N}_j^l(y) \otimes \tilde{N}_j^m(y) dy, \quad \text{for any } 0 \leq l, m \leq K + 1, \quad (7.3.47)$$

where the Einstein summation convention still assumed over the repeated subscript index  $1 \leq j \leq d$ .

Remember that, following the convention of section 7.2.4, we denote by  $\mathbb{B}_K^{l, m} \nabla^l v \nabla^m v$  the bilinear contraction

$$\mathbb{B}_K^{l, m} \nabla^l v \nabla^m v := (\mathbb{B}_K^{l, m})_{i_1 \dots i_l j_1 \dots j_m} (\partial_{i_1 \dots i_l}^l v) (\partial_{j_1 \dots j_m}^m v).$$

The bilinear tensors allow us to formulate an approximate energy  $J_K^*(v, f, \varepsilon)$  defined for any periodic function  $v \in H^{K+1}(D)$  by:

$$J_K^*(v, f, \varepsilon) := \int_D \left( \frac{1}{2} \sum_{l, m=0}^{K+1} \varepsilon^{l+m-2} \mathbb{B}_K^{l, m} \nabla^l v \nabla^m v - fv \right) dx. \quad (7.3.48)$$

The definition of the energy  $J_K^*(v, f, \varepsilon)$  is motivated by the following asymptotic provided by lemma 7.3, which holds with any  $m \geq 0$  arbitrarily high:

$$J(w_{\varepsilon, K}(v), f) = J_K^*(v, f, \varepsilon) + o(\varepsilon^m).$$

More precisely, the following result holds:

**Proposition 7.10.** *Assume  $f \in C^\infty(D)$  and  $D$ -periodic. Let  $v \in C^\infty(D)$  be a smooth  $D$ -periodic function and  $w_{\varepsilon,K}(v) \in C^\infty(D_\varepsilon)$  be the truncated ansatz of the form of (7.3.44). The following asymptotic estimate holds true with  $m \geq 0$  arbitrarily high:*

$$|J(w_{\varepsilon,K}(v), f) - J_K^*(v, f, \varepsilon)| \leq C_{K,m} (\|\nabla^{m+2}v\|_{L^\infty(D, \mathbb{R}^d)}^2 + \|f\|_{H^m(D)}^2) \varepsilon^m.$$

for a constant  $C_{K,m}$  depending only on  $m$  and  $K$ .

*Proof.* For any  $1 \leq j \leq d$ , the partial derivative  $\partial_{x_j} w_{\varepsilon,K}(v)$  reads

$$\begin{aligned} \partial_{x_j} w_{\varepsilon,K}(v) &= \sum_{k=0}^K (\varepsilon^{k-1} \partial_{y_j} N^k(\cdot/\varepsilon) \cdot \nabla^k v + \varepsilon^k N^k(\cdot/\varepsilon) \otimes e_j \cdot \nabla^{k+1} v) \\ &= \sum_{i=0}^{K+1} \varepsilon^{k-1} \tilde{N}_j^i(\cdot/\varepsilon) \cdot \nabla^k v, \end{aligned}$$

by definition (7.3.46) of the tensors  $\tilde{N}_j^k$ . The computation of the energy  $J(w_{\varepsilon,K}(v), f)$  yields then

$$\begin{aligned} J(w_{\varepsilon,K}(v), f) &= \int_D \left( \frac{1}{2} \sum_{l,m=0}^{K+1} \varepsilon^{l+m-2} (\tilde{N}_j^l(x/\varepsilon) \cdot \nabla^l v) (\tilde{N}_j^m(x/\varepsilon) \cdot \nabla^m v) - \sum_{l=0}^K \varepsilon^l (N^l(x/\varepsilon) \cdot \nabla^l v(x)) f(x) \right) dx. \end{aligned}$$

The result follows from the application of lemma 7.3 and by using that  $N^k(x/\varepsilon)$  is of average 1 if  $k = 0$  and 0 otherwise (proposition 7.9).  $\square$

The approximate energy (7.3.48) is used (instead of  $J(w_\varepsilon(v), f)$  in (7.3.45)) in order to build a (well-posed) higher order homogenized equation:

**Definition 7.5.** For any  $K \in \mathbb{N}$ , we call homogenized equation of order  $2K + 2$  the Euler-Lagrange equation associated with the minimization problem

$$\min_{\substack{v \in H^{K+1}(D), \\ D\text{-periodic}}} J_K^*(v, f, \varepsilon). \quad (7.3.49)$$

This equation reads explicitly in terms of a higher order homogenized solution  $v_K^* \in H^{K+1}(D)$  as

$$\begin{cases} \sum_{k=0}^{K+1} \varepsilon^{2k-2} \mathbb{D}_K^{2k} \cdot \nabla^{2k} v_K^* = f \\ v_K^* \text{ is } D\text{-periodic,} \end{cases} \quad (7.3.50)$$

where the constant tensors  $\mathbb{D}_K^{2k}$  are defined by the formula for any  $0 \leq k \leq K + 1$ :

$$\mathbb{D}_K^{2k} := \sum_{l=0}^{2k} (-1)^l \mathbb{B}_K^{l, 2k-l}, \quad (7.3.51)$$

assuming the convention  $\mathbb{B}_K^{l,m} = 0$  for any  $l > K + 1$  or  $m > K + 1$ .

*Proof.* Let us detail slightly the derivation of (7.3.50). The Euler-Lagrange Equation of (7.3.49) reads, after a suitable integration by parts:

$$\sum_{l,m=0}^{l+1} \varepsilon^{l+m-2} \frac{1}{2} ((-1)^l + (-1)^m) \mathbb{B}_K^{l,m} \nabla^{l+m} v_K^* = f.$$

Since  $(-1)^k + (-1)^l$  vanishes when  $k$  and  $l$  are not of the same parity, only terms such that  $k + l$  is even are not zero in the above equation. Hence, it rewrites as (7.3.50) with

$$\mathbb{D}_K^{2k} = \sum_{l+m=2k} \frac{1}{2} ((-1)^l + (-1)^m) \mathbb{B}_K^{l,m} = \sum_{l=0}^{2k} \frac{1}{2} ((-1)^l + (-1)^{2k-l}) \mathbb{B}_K^{l, 2k-l},$$

which eventually yields the desired expression (7.3.51).  $\square$

**Remark 7.7.** Naturally, the higher order homogenized solution  $v_K^*$  of (7.3.50) depends on  $\varepsilon$ , but we omit this dependence for the sake of notational simplicity.

**Remark 7.8.** Let us examine equation (7.3.50) at order  $K = 0$ . It reads

$$\mathbb{D}_0^2 \cdot \nabla^2 v_0^* + \varepsilon^{-2} \mathbb{D}_0^0 v_0^* = f. \quad (7.3.52)$$

The constant coefficient  $\mathbb{D}_0^0$  is given by

$$\mathbb{D}_0^0 = \mathbb{B}_0^{0,0} = \int_Y \tilde{N}_{0,j}^0 \tilde{N}_{0,j}^0 dy = \int_Y |\nabla N^0|^2 dy = (M^0)^2 \int_Y |\nabla \mathcal{X}^0|^2 dy = (\mathcal{X}^{0*})^{-1}.$$

The tensor  $\mathbb{D}_0^{2*}$  reads

$$\begin{aligned} \mathbb{D}_0^2 &= -\mathbb{B}_0^{1,1} = - \int_Y \tilde{N}_{0,j}^1 \otimes \tilde{N}_{0,j}^1 dy = - \int_Y |N^0(y)|^2 e_j \otimes e_j dy \\ &= - \left( \int_Y |N^0(y)|^2 dy \right) I = -(M^0)^2 \left( \int_Y |\mathcal{X}^0(y)|^2 dy \right) \\ &= - \left( \frac{1}{(\mathcal{X}^{0*})^2} \int_Y |\mathcal{X}^0(y)|^2 dy \right) I. \end{aligned}$$

Importantly, (7.3.52) does not coincide with the two term truncation of the infinite order homogenized equation (7.3.34), which would read instead:

$$M^2 \cdot \nabla^2 v_0^* + \varepsilon^{-2} M^0 v_0^* = f.$$

Indeed, formulas (7.3.21) and (7.3.35) imply

$$M^2 = \frac{1}{(\mathcal{X}^{0*})^2} \int_Y (\partial_j \mathcal{X}^1 \otimes \partial_j \mathcal{X}^1 - |\mathcal{X}^0|^2 I) dy \neq \mathbb{D}_0^2.$$

However, the coefficients  $M^0$  and  $\mathbb{D}_0^0$  do coincide,

Before stating the next proposition establishing the existence and uniqueness of the high order homogenized solution  $v_K^*$ , let us underline the following result which is an obvious, but somewhat important consequence of the definition (7.3.47):

**Lemma 7.4.** *The dominant tensor  $\mathbb{B}_K^{K+1, K+1}$  is symmetric and non-negative.*

**Proposition 7.11.** *Assume further that the dominant tensor  $\mathbb{B}_K^{K+1, K+1}$  is non-degenerate, that is there exists a constant  $\nu > 0$  such that*

$$\forall \xi = \xi_{i_1 \dots i_{K+1}} \in \mathbb{R}^{d^{K+1}}, \mathbb{B}_K^{K+1, K+1} \cdot \xi \xi \geq \nu |\xi|^2. \quad (7.3.53)$$

*Then there exists a unique periodic solution  $v_K^* \in H^{K+1}(D)$  to the homogenized equation (7.3.50) of order  $2K + 2$ .*

*Proof.* Let us consider the space  $V_K := \{v \in H^{K+1}(D) \mid v \text{ is } D\text{-periodic}\}$  and introduce  $a : V_K \times V_K \rightarrow \mathbb{R}$  and  $b : V_K \rightarrow \mathbb{R}$  the respective bilinear and linear forms defined for any  $v \in V_K$  by

$$a(v, v) = \int_D \sum_{k,l=0}^{K+1} \varepsilon^{k+l-2} \mathbb{B}_K^{k,l} \nabla^k v \nabla^l v dx, \quad (7.3.54)$$

$$b(v) = \int_D f v dx. \quad (7.3.55)$$

The homogenized equation (7.3.50) reduces then to the following variational problem:

$$\text{Find } v_K^* \in V_K \text{ such that } \forall v \in V_K, a(v_K^*, v) = b(v). \quad (7.3.56)$$

From there, one could directly rely on the theory of Fredholm operators [231] to conclude. However we shall show with elementary arguments that  $a$  is coercive, which will allow us to apply Lax-Milgram theorem.

Under the non-degeneracy assumption (7.3.53), it is readily obtained that there exists a constant  $C_\varepsilon$  (depending on  $\varepsilon$ ) such that

$$\forall v \in V_K(D), a(v, v) \geq (\varepsilon^{2K}\nu) \|\nabla^{K+1}v\|_{L^2(D, \mathbb{R}^d)}^2 + \varepsilon^{-2}M^0 \|v\|_{L^2(D)}^2 - C_\varepsilon \|v\|_{H^{K+1}(D)} \|v\|_{H^K(D)}.$$

Remembering  $M^0 > 0$  and applying the following Young's inequality

$$\forall x, y \in \mathbb{R}, |xy| \leq \frac{x^2}{2\gamma} + \frac{\gamma y^2}{2}$$

for a sufficiently small  $\gamma > 0$ , we obtain the existence of two constants  $\alpha_{\varepsilon, K} > 0$  and  $\beta_{\varepsilon, K} > 0$  (that depend on  $\varepsilon$  and  $K$ ) such that

$$\forall v \in V_K(D), a(v, v) \geq \alpha_{\varepsilon, K} \|v\|_{H^{K+1}(D)}^2 - \beta_{\varepsilon, K} \|v\|_{H^K(D)}^2. \tag{7.3.57}$$

Furthermore, (7.3.47) together with the proof of proposition 7.10 allow to rewrite  $a(v, v)$  as

$$a(v, v) = \int_D \int_Y \left\| (\varepsilon^{-1} \nabla_y + \nabla_x) \left( \sum_{k=0}^K \varepsilon^k N^k(y) \cdot \nabla^k v(x) \right) \right\|^2 dx.$$

Then,  $\int_D \int_Y u(x, y)^2 dy dx \geq \int_D |\int_Y u(x, y) dy|^2 dx$  implies the following inequality:

$$\forall v \in V_K, a(v, v) \geq \|\nabla v\|_{L^2(D, \mathbb{R}^d)}^2. \tag{7.3.58}$$

We shall now prove that (7.3.57) and (7.3.58) together imply the coercivity of  $a$  on the space  $V_K$ , that is we claim there exists a constant  $c_{\varepsilon, K} > 0$  such that

$$\forall v \in V_K, a(v, v) \geq c_{\varepsilon, K} \|v\|_{H^{K+1}(D)}^2. \tag{7.3.59}$$

Assume the contrary is true, then one can find a sequence  $(v_n)$  of functions satisfying  $\|v_n\|_{H^{K+1}(D)} = 1$  and such that  $a(v_n, v_n) \rightarrow 0$ . Up to extracting a relevant subsequence, we may assume that  $v_n \rightharpoonup v$  weakly in  $H^{K+1}(D)$  and  $v_n \rightarrow v$  strongly in  $H^K(D)$ . Then the polarization identity together with (7.3.57) and the positivity of  $a$  allow to show that  $(v_n)$  is a Cauchy sequence in  $V_K$ :

$$\begin{aligned} \forall l, m \in \mathbb{N}, \alpha_{\varepsilon, K} \|v_l - v_m\|_{H^{K+1}(D)}^2 &\leq a(v_l - v_m, v_l - v_m) + \beta_{\varepsilon, K} \|v_l - v_m\|_{H^K(D)}^2 \\ &= 2a(v_l, v_l) + 2a(v_m, v_m) - a(v_l + v_m, v_l + v_m) + \beta_{\varepsilon, K} \|v_l - v_m\|_{H^K(D)}^2 \\ &\leq 2a(v_l, v_l) + 2a(v_m, v_m) + \beta_{\varepsilon, K} \|v_l - v_m\|_{H^K(D)}^2 \xrightarrow{l, m \rightarrow \infty} 0. \end{aligned}$$

Therefore  $v_n \rightarrow v$  strongly in  $V_K$ . The continuity of  $a$  implies then  $a(v, v) = \lim_{n \rightarrow +\infty} a(v_n, v_n) = 0$ . The property (7.3.58) yields then that  $v$  is a constant. Therefore,  $0 = a(v, v) = \varepsilon^{-2}M^0 \|v\|_{L^2(D)}^2$ , which implies  $v = 0$ . This is in contradiction with the fact that  $\|v_n\|_{H^{K+1}(D)} = 1$  for any  $n \geq 0$  and the strong convergence of  $(v_n)$ . Finally, the coercivity (7.3.59) and the continuity of  $a$  and  $b$  over  $V_K$  ensure that all the assumptions of the Lax-Milgram theorem are fulfilled, which yields existence and uniqueness to the problem (7.3.56).  $\square$

**Remark 7.9.** The non degeneracy assumption (7.3.53) is automatically fulfilled for any shape of obstacle  $\eta T$  when  $K = 0$ . It could fail to be satisfied for particular obstacle shapes for  $K \geq 1$  (e.g. under strong symmetries with respect to the cell axes).

Before going to the proof of error estimates for the higher homogenized solution  $v_K^*$ , we provide a last result which shows that the high order homogenized equation (7.3.50) is in some sense a truncation of the formal infinite order homogenized equation (7.3.34). This fact was also observed by [287] for a (scalar) antiplane elasticity model.

**Proposition 7.12.** *The first  $K + 1$  homogenized coefficients of the homogenized equation (7.3.50) of order  $2K + 2$  coincide with those of the formal equation (7.3.34):*

$$\forall k \leq \lfloor K/2 \rfloor, \mathbb{D}_K^{2k} = M^{2k}.$$

*Proof.* We show the following, slightly more general, result:

$$\forall k \leq K, M^k = \sum_{l=0}^k (-1)^l \mathbb{B}_K^{l,k-l}, \quad (7.3.60)$$

which is sufficient for our purpose because of (7.3.51). For  $k, l \leq K$ , the coefficient  $\mathbb{B}_K^{l,k-l}$  is given by (from (7.3.46))

$$\mathbb{B}_K^{l,k-l} = \int_Y (\partial_j N^l + N^{l-1} \otimes e_j) \otimes (\partial_j N^{k-l} + N^{k-l-1} \otimes e_j) dy,$$

where we use the convention  $N^{-1} = N^{-2} = 0$ . After an integration by parts, we rewrite  $\mathbb{B}_K^{l,k-l}$  as follows:

$$\begin{aligned} \mathbb{B}_K^{l,k-l} &= \int_Y (-\Delta N^l - 2\partial_j N^{l-1} \otimes e_j - N^{l-2} \otimes I) \otimes N^{k-l} dy \\ &+ \int_Y (\partial_j N^l \otimes N^{k-l-1} \otimes e_j + N^{l-1} \otimes N^{k-l-1} \otimes I + \partial_j N^{l-1} \otimes N^{k-l} \otimes e_j + N^{l-2} \otimes N^{k-l} \otimes I) dy \\ &= \int_Y (M^l \otimes N^{k-l}) dy + B^{k,l} + B^{k,l-1} \end{aligned}$$

where  $B^{k,l}$  is the  $k$ -th order tensor defined by

$$B^{k,l} := \int_Y (\partial_j N^l \otimes N^{k-l-1} \otimes e_j + N^{l-1} \otimes N^{k-l-1} \otimes I) dy.$$

Using now the point 1. of proposition 7.9 and recognizing a telescopic series, we finally obtain

$$\begin{aligned} \sum_{l=0}^k (-1)^l \mathbb{B}_K^{l,k-l} &= (-1)^k M^k + \sum_{l=0}^k ((-1)^l B^{k,l} - (-1)^{l-1} B^{k,l-1}) \\ &= (-1)^k M^k + (-1)^k B^{k,k} - (-1)^{-1} B^{k,-1}. \end{aligned}$$

The result follows from the fact that  $M^k = 0$  when  $k$  is odd and  $B^{k,k} = B^{k,-1} = 0$  with our convention  $N^{-1} = 0$ .  $\square$

### 7.3.4 Error estimates and justification of the higher order homogenization process

The main result of this subsection is proposition 7.13 below: we show that the truncated ansatz  $w_{\varepsilon,K}(v_K^*)$  (equation (7.3.44)) built from the solution  $v_K^*$  of the homogenized equation (7.3.50) yields an approximation of the original solution  $u_\varepsilon$  of order  $K+2$  in the  $H^1(D)$  norm, and of order  $K+3$  in the  $L^2(D)$  norm.

The first step of the proof consists in showing that under periodicity assumptions,  $u_\varepsilon$  can be effectively approximated by truncated ansatz of the form of (7.3.44):

**Lemma 7.5.** *Let  $u_{\varepsilon,K}^*$  be the average of the truncated expansion  $u_{\varepsilon,K}$  (7.3.23):*

$$\forall x \in D, u_{\varepsilon,K}^*(x) := \sum_{k=0}^K \varepsilon^{k+2} u_k^*(x).$$

There exists a constant  $C_K(f)$  independent of  $\varepsilon$  such that

$$\left\| u_\varepsilon - \sum_{k=0}^K \varepsilon^k N^k(\cdot/\varepsilon) \cdot \nabla^k u_{\varepsilon,K}^* \right\|_{H^1(D)} \leq C_K(f) \varepsilon^{K+2}.$$

*Proof.* We use proposition 7.8 to rewrite

$$\begin{aligned} u_{\varepsilon,K}(x) &= \sum_{i=0}^K \sum_{k=0}^i \varepsilon^{2+i} N^k(x/\varepsilon) \cdot \nabla^k u_{i-k}^*(x) = \sum_{k=0}^K \sum_{i=k}^K \varepsilon^k \varepsilon^{i-k+2} N^k(x/\varepsilon) \cdot \nabla^k u_{i-k}^*(x) \\ &= \sum_{k=0}^K \varepsilon^k N^k(x/\varepsilon) \cdot \nabla^k \left( \sum_{i=0}^{K-k} \varepsilon^{i+2} u_i^*(x) \right) \\ &= \sum_{k=0}^K \varepsilon^k N^k(x/\varepsilon) \cdot \nabla^k u_{\varepsilon,K}^*(x) - \sum_{k=0}^K \sum_{i=K-k+1}^K \varepsilon^{k+i+2} N^k(x/\varepsilon) \cdot \nabla^k u_i^*(x). \end{aligned}$$

Hence using the result of [proposition 7.4](#),

$$\begin{aligned} \left\| u_\varepsilon - \sum_{k=0}^K \varepsilon^k N^k(\cdot/\varepsilon) \cdot \nabla^k u_{\varepsilon,K}^* \right\|_{H^1(D)} &\leq \|u_\varepsilon - u_{\varepsilon,K}\|_{H^1(D)} + \left\| u_{\varepsilon,K} - \sum_{k=0}^K \varepsilon^k N^k(x/\varepsilon) \cdot \nabla^k u_{\varepsilon,K}^*(x) \right\|_{H^1(D)} \\ &\leq C_K \|f\|_{H^{K+2}(D)} \varepsilon^{K+2} + \varepsilon^{K+3} \sum_{k=0}^K \sum_{i=K-k+1}^K \|N^k(\cdot/\varepsilon) \cdot \nabla^k u_i^*\|_{H^1(D)} \\ &\leq C_K \|f\|_{H^{K+2}(D)} \varepsilon^{K+2} + C_K \varepsilon^{K+2} \|f\|_{H^{2K+1}(D)} \end{aligned}$$

whence the result.  $\square$

In order to state our final result, we need the following lemma providing uniform estimates of the partial derivatives of  $v_K^*$ :

**Lemma 7.6.** *The solution  $v_K^*$  of (7.3.50) belongs to  $C^\infty(D)$  and for any  $m \in \mathbb{N}$ , there exists a constant  $C_m$  that does not depend on  $\varepsilon$  such that*

$$\|v_K^*\|_{H^{m+2}(D)} \leq C_m \|f\|_{H^m(D)} \varepsilon^2. \quad (7.3.61)$$

*Proof.* The result is obvious by solving (7.3.50) explicitly with Fourier expansions.  $\square$

Using the previous results, ideas that resemble those of Cea's Lemma allow finally to derive the following error estimate for the higher order ansatz  $w_{\varepsilon,K}(v_K^*)$  (eqn. (7.3.44)):

**Proposition 7.13.** *Let  $v_K^*$  be the solution to the homogenized equation (7.3.50) of order  $2K+2$ . There exists a constant  $C_K(f)$  independent of  $\varepsilon$  (but depending on  $K$  and  $f$ ) such that the following error estimates hold:*

$$\left\| u_\varepsilon - \sum_{k=0}^K \varepsilon^k N^k(\cdot/\varepsilon) \cdot \nabla^k v_K^* \right\|_{H^1(D)} \leq C_K(f) \varepsilon^{K+2}. \quad (7.3.62)$$

$$\left\| u_\varepsilon - \sum_{k=0}^K \varepsilon^k N^k(\cdot/\varepsilon) \cdot \nabla^k v_K^* \right\|_{L^2(D)} \leq C_K(f) \varepsilon^{K+3}. \quad (7.3.63)$$

*Proof.* From, [lemma 7.3](#) and by the definition (7.3.49) of  $w_{\varepsilon,K}(v_K^*)$ , itnfer holds for any smooth periodic function  $\phi \in C^\infty(D)$  and with any  $m \geq 0$  arbitrarily high:

$$\begin{aligned} \left| \int_D \nabla v_{\varepsilon,K} \cdot \nabla \left( \sum_{k=0}^K \varepsilon^k N^k(x/\varepsilon) \cdot \nabla^k \phi(x) \right) dx - \int_D f \phi dx \right| &\leq C_m \varepsilon^m \|\phi v_K^*\|_{H^{m+2}(D)} \\ &\leq C'_m \varepsilon^m \|\phi\|_{H^{m+2}(D)} \|f\|_{H^m(D)}. \end{aligned}$$

Furthermore by definition of  $u_\varepsilon$  it also holds:

$$\begin{aligned} \left| \int_D \nabla u_\varepsilon \cdot \nabla \left( \sum_{k=0}^K \varepsilon^k N^k(x/\varepsilon) \cdot \nabla^k \phi(x) \right) dx - \int_D f \phi dx \right| \\ \leq C_m \varepsilon^m \|f \phi\|_{H^m(D)} \leq C'_m \varepsilon^m \|f\|_{H^m(D)} \|\phi\|_{H^{m+2}(D)}. \end{aligned}$$

Therefore  $u_\varepsilon - v_{\varepsilon,K}$  is nearly orthogonal in the  $H^1$  scalar product (up to some terms of order  $O(\varepsilon^m) \leq C'_m \varepsilon^m (\|\phi\|_{H^{m+2}(D)} \|f\|_{H^m(D)})$ ) to the test fields of the form  $w_{\varepsilon,K}(\phi) := \sum_{k=0}^K \varepsilon^k N^k(\cdot/\varepsilon) \cdot \nabla^k \phi$ . Hence with successively  $\phi = v_K^*$  and  $\phi = u_{\varepsilon,K}^*$ :

$$\begin{aligned} \|\nabla(u_\varepsilon - w_{\varepsilon,K}(v_K^*))\|_{L^2(D, \mathbb{R}^d)}^2 &= \int_D |\nabla(u_\varepsilon - w_{\varepsilon,K}(v_K^*))|^2 dx \\ &\leq \int_D \nabla(u_\varepsilon - w_{\varepsilon,K}(v_K^*)) \cdot \nabla u_\varepsilon dx + C'_m \varepsilon^m \|f\|_{H^m(D)} \|v_K^*\|_{H^{m+2}(D)} \\ &\leq \int_D \nabla(u_\varepsilon - w_{\varepsilon,K}(v_K^*)) \cdot \nabla(u_\varepsilon - w_{\varepsilon,K}(u_{\varepsilon,K}^*)) dx \\ &\quad + C'_m \varepsilon^m \|f\|_{H^m(D)} \|v_K^*\|_{H^{m+2}(D)} + C''_m \varepsilon^m \|f\|_{H^m(D)} \|u_{\varepsilon,K}^*\|_{H^{m+2}(D)}. \end{aligned}$$



From the definition (7.3.30) and the regularity estimate lemma 7.6, we infer that

$$\|u_{\varepsilon,K}^*\|_{H^{m+1}(D)} \leq C_m \|f\|_{H^{K+m+1}} \text{ and } \|v_K^*\|_{H^{m+1}(D)} \leq C_m \|f\|_{H^{m+3}(D)}.$$

Hence we finally obtain

$$\begin{aligned} \|\nabla(u_\varepsilon - w_{\varepsilon,K}(v_K^*))\|_{L^2(D,\mathbb{R}^d)}^2 &\leq \frac{1}{2} \|\nabla(u_\varepsilon - w_{\varepsilon,K}(v_K^*))\|_{L^2(D,\mathbb{R}^d)}^2 + \frac{1}{2} \|\nabla(u_\varepsilon - w_{\varepsilon,K}(u_{\varepsilon,K}^*))\|_{L^2(D,\mathbb{R}^d)}^2 \\ &\quad + C_m \varepsilon^m \|f\|_{H^{K+m+1}(D)} \\ &\leq \frac{1}{2} \|\nabla(u_\varepsilon - w_{\varepsilon,K}(v_K^*))\|_{L^2(D,\mathbb{R}^d)}^2 + C_K \varepsilon^{2K+4} \|f\|_{H^{K+2}(D)} + C_m \varepsilon^m \|f\|_{H^{K+m+1}(D)} \end{aligned}$$

where we have applied Young’s inequality  $|ab| \leq \frac{1}{2}(a^2 + b^2)$  and the error bound (7.3.24) in the last two lines. Setting  $m = 2K + 4$  implies the first error estimate (7.3.62). Then (7.3.63) follows from the Poincaré inequality of lemma 7.2. □

**Remark 7.10.** We knew already from proposition 7.4 with  $K = 0$  that  $\|u_\varepsilon - N^0(x/\varepsilon)\varepsilon^2 \mathcal{X}^{0*} f\|_{H^1(D)} = O(\varepsilon^2)$ . Therefore, the solution  $w_{\varepsilon,K}(v_K^*)$  obtained from the homogenized equation (7.3.50) of order 2 with  $K = 0$  does not provide a better order of magnitude. However, we may argue that the approximation  $N^0(x/\varepsilon)v_0^*$  is better in the energy norm, since from the minimization principle (7.3.49), it must hold

$$\|u_\varepsilon - N^0(\cdot/\varepsilon)v_0^*\|_{H^1(D)} \leq \|u_\varepsilon - N^0(x/\varepsilon)\varepsilon^2 \mathcal{X}^{0*} f\|_{H^1(D)} + O(\varepsilon^m)$$

for any arbitrarily large value of  $m$ . The next section will provide additional supporting arguments that (7.3.50) provides more robust homogenized approximations since this equation does not degenerate when the size  $\eta$  of the obstacle  $\eta T$  vanishes (it converges to the Poisson problem without holes).

**Remark 7.11.** We could have hoped, at first glance, to obtain similar bounds for (7.3.1) in case where the periodicity condition on  $\partial D$  is replaced by a Dirichlet boundary conditions  $v_K^* = 0$  on  $\partial D$ . Indeed, the truncated ansatz  $w_{\varepsilon,K}(v_K^*)$  would then satisfy simultaneously the boundary conditions  $v_{\varepsilon,K} = 0$  on  $\partial D$  and  $\partial\omega_\varepsilon$ . However our proposed error analysis does not extend in a straightforward manner: one of the difficulty lies in that the result of lemma 7.6 does not hold because of the arising of boundary layers of  $v_{\varepsilon,K}$  and  $u_\varepsilon$  near  $\partial D$ . See [216, 19] for some examples of treatments of related problems.

**Remark 7.12.** There are other ways to build well-posed higher order homogenized equations, for instance by filtering the right-hand side  $f$  [33], or by resorting to a Boussinesq trick [20, 3, 264]. These approaches yield order  $K$  estimates by solving order  $K$  equations only (instead of  $2K + 2$  as in (7.3.50)), however they resort to the evaluation of some partial derivatives of  $f$  to the right-hand side (which may be a numerically delicate task).

**Remark 7.13.** The proof of proposition 7.13 is inspired from [287], however it does not extend straightforwardly to the context of the Stokes system (7.2.3) because of the divergence constraint. A different proof is proposed in the dedicated section 7.5.3.

**7.3.5 Low volume fraction limits when the size of the obstacle tends to zero**

The purpose of this section is to establish asymptotics for the tensors  $\mathcal{X}^{k*}$  and  $M^k$  in the low volume fraction limit, namely when the size  $\eta$  of the obstacle vanishes to zero. Our main result is stated in corollary 7.2, which implies that the infinite order homogenized equation (7.3.34) converges formally to either of the three classical regimes of the literature (namely, to the original Laplace equation (7.3.1), or to the analogue of the Brinkman or Darcy equation (7.2.6) and (7.2.7)).

In this whole subsection, it is assumed, for simplicity, that the space dimension is greater than 3:

$$d \geq 3. \tag{7.3.64}$$

The case  $d = 2$  requires a specific treatment (see e.g. [13]), although very similar results could be stated.

The hole  $\eta T$  is assumed to be strictly included in the unit cell for any  $\eta \leq 1$  (it does not touch the boundary):  $\eta T \subset\subset P$ . For a given function  $\tilde{v} \in L^2(\eta^{-1}P)$ , we denote by  $\langle \tilde{v} \rangle$  the average  $\langle \tilde{v} \rangle := \eta^d \int_{\eta^{-1}P} \tilde{v}(y) dy$ . Following the methodology of [13] (and also [196]), we shall use several times the following useful lemma:

**Lemma 7.7.** *Assume  $d \geq 3$ . There exists a constant  $C > 0$  independent of  $\eta > 0$  such that for any  $\tilde{v} \in H^1(\eta^{-1}P \setminus T)$  which vanishes on the hole  $\partial T$  and which is  $\eta^{-1}P$  periodic, the following inequalities hold:*

$$\|\tilde{v}\|_{L^2(\eta^{-1}P \setminus T)} \leq C\eta^{-d/2}\|\nabla\tilde{v}\|_{L^2(\eta^{-1}P \setminus T, \mathbb{R}^d)}, \quad (7.3.65)$$

$$|\langle \tilde{v} \rangle| \leq C\|\nabla\tilde{v}\|_{L^2(\eta^{-1}P \setminus T, \mathbb{R}^d)}, \quad (7.3.66)$$

$$\|\tilde{v} - \langle \tilde{v} \rangle\|_{L^2(\eta^{-1}P \setminus T)} \leq C\eta^{-1}\|\nabla\tilde{v}\|_{L^2(\eta^{-1}P \setminus T, \mathbb{R}^d)}, \quad (7.3.67)$$

$$\|\tilde{v} - \langle \tilde{v} \rangle\|_{L^{2d/(d-2)}(\eta^{-1}P \setminus T)} \leq C\|\nabla\tilde{v}\|_{L^2(\eta^{-1}P \setminus T, \mathbb{R}^d)}. \quad (7.3.68)$$

*Proof.* See [13, 210]. □

Let us also recall that for any  $v \in H^1(P \setminus (\eta T))$ , if  $\tilde{v}$  is taken to be the rescaled function  $\tilde{v}(y) \equiv v(\eta y)$  in the rescaled cell  $\eta^{-1}P \setminus T$ , then the  $L^2$  norms of the functions and their gradients are related by the following identities:

$$\|v\|_{L^2(P \setminus (\eta T))} = \eta^{d/2}\|\tilde{v}\|_{L^2(\eta^{-1}P \setminus T)} \text{ and } \|\nabla v\|_{L^2(P \setminus (\eta T), \mathbb{R}^d)} = \eta^{d/2-1}\|\nabla\tilde{v}\|_{L^2(\eta^{-1}P \setminus T, \mathbb{R}^d)}.$$

We also need to consider the so-called Deny-Lions space  $\mathcal{D}^{1,2}(\mathbb{R}^d \setminus T)$  for which we recall the definition (the reader is referred to [13, 11, 15] and also [243], p.59. for more details).

**Definition 7.6** (Deny-Lions space). The Deny-Lions space  $\mathcal{D}^{1,2}(\mathbb{R}^d \setminus T)$  is the completion of the space of smooth functions by the  $L^2$  norm of their gradients:

$$\mathcal{D}^{1,2}(\mathbb{R}^d \setminus T) := \overline{\mathcal{D}(\mathbb{R}^d \setminus T)}^{\|\nabla \cdot\|_{L^2(\mathbb{R}^d \setminus T, \mathbb{R}^d)}}.$$

When  $d \geq 3$ , it admits the following characterization:

$$\mathcal{D}^{1,2}(\mathbb{R}^d \setminus T) = \{\phi \text{ measurable} \mid \|\phi\|_{L^{2d/(d-2)}(\mathbb{R}^d \setminus T)} < +\infty \text{ and } \|\nabla\phi\|_{L^2(\mathbb{R}^d \setminus T, \mathbb{R}^d)} < +\infty\}.$$

We introduce  $\Psi$  the unique solution to the exterior problem

$$\begin{cases} -\Delta\Psi = 0 \text{ in } \mathbb{R}^d - T \\ \Psi = 0 \text{ on } \partial T \\ \Psi \rightarrow 1 \text{ at } \infty, \end{cases} \quad (7.3.69)$$

and we denote by  $\Psi^*$  the normal flux

$$\Psi^* := \int_{\mathbb{R}^d \setminus T} |\nabla\Psi|^2 dx = - \int_{\partial T} \nabla\Psi \cdot \mathbf{n} ds,$$

where  $\mathbf{n}$  is the normal pointing *inward*  $T$ . The condition  $\Psi \rightarrow 1$  at  $\infty$  is to be understood in the sense that  $\Psi - 1 \in \mathcal{D}^{1,2}(\mathbb{R}^d \setminus T)$ .

The following result provides asymptotics for the tensors  $\mathcal{X}^k$  and their averages  $\mathcal{X}^{k*}$ . The special case  $k = 0$  is a reformulation of Theorem 3.1 of [13] (see also [196]) to the case of the Poisson system.

**Proposition 7.14.** *Assume  $d \geq 3$ . For any  $k \geq 0$ , denote by  $\tilde{\mathcal{X}}^{2k}$  and  $\tilde{\mathcal{X}}^{2k+1}$  the rescaled tensors in  $\eta^{-1}P \setminus T$  defined by:*

$$\forall x \in \eta^{-1}P \setminus T, \tilde{\mathcal{X}}^{2k}(x) := \eta^{(d-2)(k+1)}\mathcal{X}^{2k}(\eta x) \text{ and } \tilde{\mathcal{X}}^{2k+1}(x) := \eta^{(d-2)(k+1)}\mathcal{X}^{2k+1}(\eta x).$$

*Then:*

1. *there exists a constant  $C > 0$  independent of  $\eta > 0$  such that:*

$$\forall \eta > 0, \|\nabla\tilde{\mathcal{X}}^{2k}\|_{L^2(\eta^{-1}P \setminus T, \mathbb{R}^d)} \leq C \text{ and } \|\nabla\tilde{\mathcal{X}}^{2k+1}\|_{L^2(\eta^{-1}P \setminus T, \mathbb{R}^d)} \leq C; \quad (7.3.70)$$

2. *the following asymptotic convergences hold:*

$$\tilde{\mathcal{X}}^{2k} \rightharpoonup \frac{\Psi}{|\Psi^*|^{k+1}} I^{2k}, \text{ weakly in } H_{loc}^1(\mathbb{R}^d \setminus T) \quad (7.3.71)$$

$$\tilde{\mathcal{X}}^{2k+1} \rightharpoonup 0 \text{ weakly in } H_{loc}^1(\mathbb{R}^d \setminus T) \quad (7.3.72)$$

$$\mathcal{X}^{2k*} \sim \frac{1}{\eta^{(d-2)(k+1)} |\Psi^*|^{k+1}} I^{2k}, \quad (7.3.73)$$

where we recall the definition  $I^{2k} := \overbrace{I \otimes I \otimes \dots \otimes I}^{k \text{ times}}$  of [section 7.2.4](#).

**Remark 7.14.** Let us recall that we already know  $\mathcal{X}^{2k+1*} = 0$  for any  $k \in \mathbb{N}$  ([proposition 7.2](#)). However this is not the case in the vectorial context; in that case we shall adapt the arguments in order to obtain asymptotics also for  $\mathcal{X}^{2k+1*}$  ([proposition 7.27](#) below).

*Proof.* The result is proved by induction.

1. *Case  $2k$  for  $k = 0$ .* The tensors  $\tilde{\mathcal{X}}^0$  satisfies  $-\Delta \tilde{\mathcal{X}}^0 = \eta^d$  in  $\eta^{-1}P \setminus T$ , hence for any  $\Phi \in H^1(\eta^{-1}P \setminus T)$  which is  $\eta^{-1}P$ -periodic, it holds

$$\int_{\eta^{-1}P \setminus T} \nabla \tilde{\mathcal{X}}^0 \cdot \nabla \Phi \, dy = \eta^d \int_{\eta^{-1}P \setminus T} \Phi \, dy + \int_{\partial T} (\nabla \tilde{\mathcal{X}}^0 \cdot \mathbf{n}) \Phi \, ds. \quad (7.3.74)$$

Setting  $\Phi = \tilde{\mathcal{X}}^0$  in (7.3.74) and using successively the Cauchy-Schwartz inequality and [lemma 7.7](#) yields

$$\|\nabla \tilde{\mathcal{X}}^0\|_{L^2(\eta^{-1}P \setminus T, \mathbb{R}^d)}^2 = \eta^d \int_{\eta^{-1}P \setminus T} \tilde{\mathcal{X}}^0 \, dy \leq \eta^d \eta^{-d/2} \|\tilde{\mathcal{X}}^0\|_{L^2(\eta^{-1}P \setminus T)} \leq C \|\nabla \tilde{\mathcal{X}}^0\|_{L^2(\eta^{-1}P \setminus T, \mathbb{R}^d)}.$$

This implies the first equality of (7.3.70) for  $k = 0$ . Furthermore, (7.3.66) also shows that  $\langle \tilde{\mathcal{X}}^0 \rangle$  is bounded by a constant independent of  $\eta$ . Hence, up to extracting a subsequence, there exists a constant  $c^0 \in \mathbb{R}$  and a function  $\widehat{\Psi}^0$  such that

$$\langle \tilde{\mathcal{X}}^0 \rangle \rightharpoonup c^0 \text{ and } \tilde{\mathcal{X}}^0 \rightharpoonup \widehat{\Psi}^0 \text{ weakly in } H_{loc}^1(\mathbb{R}^d \setminus T) \text{ when } \eta \rightarrow 0.$$

Furthermore, the lower semi-continuity of the  $H_{loc}^1(\mathbb{R}^d \setminus T)$  norm and (7.3.68) imply that  $\widehat{\Psi}^0 - c^0$  belongs to  $\mathcal{D}^{1,2}(\mathbb{R}^d \setminus T)$  (see [13] for a detailed justification). Setting now  $\Phi$  with compact support in  $\mathbb{R}^d \setminus \bar{T}$  in (7.3.74) and passing to the limit when  $\eta \rightarrow 0$  yields that  $\widehat{\Psi}^0$  is solution to the exterior problem

$$\begin{cases} -\Delta \widehat{\Psi}^0 = 0 \text{ in } \mathbb{R}^d \setminus T \\ \widehat{\Psi}^0 = 0 \text{ on } \partial T \\ \widehat{\Psi}^0 \rightarrow c^0 \text{ at } \infty. \end{cases} \quad (7.3.75)$$

Obviously, the unique solution to this problem is given by  $\widehat{\Psi} = c^0 \Psi$ . Finally, the constant  $c^0$  can be identified by setting  $\Phi = 1$  in (7.3.74), which yields

$$-\int_{\partial T} \nabla \widehat{\Psi}^0 \cdot \mathbf{n} \, ds = \lim_{\eta \rightarrow 0} -\int_{\partial T} \nabla \tilde{\mathcal{X}}^0 \cdot \mathbf{n} \, ds = 1.$$

Therefore,  $c_0 = 1/\Psi^*$  from where (7.3.71) follows for  $k = 0$ . Since the obtained limit is unique, the convergence holds for the whole sequence. Then (7.3.73) follows from a simple change of variable.

2. *Case  $2k + 1$  for  $k = 0$ .* It holds  $-\Delta \tilde{\mathcal{X}}^1 = 2\eta \partial_j \tilde{\mathcal{X}}^0 \otimes e_j$ , hence for any  $\Phi \in H^1(\eta^{-1}P \setminus T)$  which is  $\eta^{-1}P$ -periodic:

$$\int_{\eta^{-1}P \setminus T} \nabla \tilde{\mathcal{X}}^1 \cdot \nabla \Phi \, dy = \int_{\eta^{-1}P \setminus T} 2\eta (\partial_j \tilde{\mathcal{X}}^0 \otimes e_j) \Phi \, dy + \int_{\partial T} (\nabla \tilde{\mathcal{X}}^1 \cdot \mathbf{n}) \Phi \, ds. \quad (7.3.76)$$

Taking  $\Phi = \tilde{\mathcal{X}}_j^1$  for a given fixed  $j$  yields

$$\begin{aligned} \|\nabla \tilde{\mathcal{X}}_j^1\|_{L^2(\eta^{-1}P \setminus T, \mathbb{R}^d)}^2 &= \int_{\eta^{-1}P \setminus T} 2\eta \partial_j \tilde{\mathcal{X}}^0 (\tilde{\mathcal{X}}_j^1 - \langle \tilde{\mathcal{X}}_j^1 \rangle) \, dy \\ &\leq 2\eta \|\partial_j \tilde{\mathcal{X}}^0\|_{L^2(\eta^{-1}P \setminus T)} \|\tilde{\mathcal{X}}_j^1 - \langle \tilde{\mathcal{X}}_j^1 \rangle\|_{L^2(\eta^{-1}P \setminus T)} \\ &\leq C\eta \eta^{-1} \|\nabla \tilde{\mathcal{X}}_j^1\|_{L^2(\eta^{-1}P \setminus T, \mathbb{R}^d)}. \end{aligned} \quad (7.3.77)$$

This implies the second part of (7.3.70) for  $k = 0$ . As before, this means that up to extracting a subsequence, we may assume the existence of a constant vector  $c^1 \in \mathbb{R}^d$  and a first order tensor  $\widehat{\Psi}^1$  such that  $\langle \widetilde{\mathcal{X}}^1 \rangle \rightarrow c^1$  and  $\widetilde{\mathcal{X}}^1 \rightharpoonup \Psi^1$  in  $H_{loc}^1(\mathbb{R}^d \setminus T)$  (note that we already know that  $\langle \widetilde{\mathcal{X}}^1 \rangle = 0$  and  $c^1 = 0$  but the coming ideas will extend to the vectorial case) with  $\widetilde{\mathcal{X}}^1 - c^1$  belonging to  $\mathcal{D}^{1,2}(\mathbb{R}^d \setminus T)$ . Setting  $\Phi$  with compact support in  $\mathbb{R}^d \setminus \overline{T}$  in (7.3.76) yields that  $\widehat{\Psi}^1$  is solution to the exterior problem

$$\begin{cases} -\Delta \widehat{\Psi}^1 = 0 & \text{in } \mathbb{R}^d \setminus T \\ \widehat{\Psi}^1 = 0 & \text{on } \partial T \\ \widehat{\Psi}^1 \rightarrow c^1 & \text{at } \infty. \end{cases} \quad (7.3.78)$$

Therefore,  $\widehat{\Psi}^1 = c^1 \Psi$  and the constant  $c^1$  is identified by taking  $\Phi = 1$  in (7.3.76), which states that:

$$-\int_{\partial T} \nabla \widehat{\Psi}^1 \cdot \mathbf{n} ds = \lim_{\eta \rightarrow 0} -\int_{\partial T} \nabla \widetilde{\mathcal{X}}^1 \cdot \mathbf{n} ds = 0.$$

This implies  $c^1 = 0$  and (7.3.72) is obtained for  $k = 0$ .

3. *General case.* We now complete the proof by induction on  $k$ . Assuming the result holds till rank  $k$ , we compute

$$\begin{aligned} -\Delta \widetilde{\mathcal{X}}^{2k+2} &= 2\eta^{d-1} \partial_j \widetilde{\mathcal{X}}^{2k+1} \otimes e_j + \eta^d \widetilde{\mathcal{X}}^{2k} \otimes I, \\ -\Delta \widetilde{\mathcal{X}}^{2k+3} &= 2\eta \partial_j \widetilde{\mathcal{X}}^{2k+2} \otimes e_j + \eta^d \widetilde{\mathcal{X}}^{2k+1} \otimes I. \end{aligned}$$

Hence for any  $\Phi \in H^1(\eta^{-1}P \setminus T)$  which is  $\eta^{-1}P$ -periodic:

$$\int_{\eta^{-1}P \setminus T} \nabla \widetilde{\mathcal{X}}^{2k+2} \cdot \nabla \Phi dy = \int_{\eta^{-1}P \setminus T} (2\eta^{d-1} \partial_j \widetilde{\mathcal{X}}^{2k+1} \otimes e_j + \eta^d \widetilde{\mathcal{X}}^{2k} \otimes I) \Phi dy + \int_{\partial T} (\nabla \widetilde{\mathcal{X}}^{2k+2} \cdot \mathbf{n}) \Phi ds, \quad (7.3.79)$$

$$\int_{\eta^{-1}P \setminus T} \nabla \widetilde{\mathcal{X}}^{2k+3} \cdot \nabla \Phi dy = \int_{\eta^{-1}P \setminus T} (2\eta \partial_j \widetilde{\mathcal{X}}^{2k+2} \otimes e_j + \eta^d \widetilde{\mathcal{X}}^{2k+1} \otimes I) \Phi dy + \int_{\partial T} (\nabla \widetilde{\mathcal{X}}^{2k+3} \cdot \mathbf{n}) \Phi ds. \quad (7.3.80)$$

Setting  $\Phi = \widetilde{\mathcal{X}}^{2k+2}$  and  $\Phi = \widetilde{\mathcal{X}}^{2k+3}$  respectively in the above variational formulations (for a fixed set of indices), applying the Cauchy-Schwartz inequality and using lemma 7.7 yield

$$\begin{aligned} \|\nabla \widetilde{\mathcal{X}}^{2k+2}\|_{L^2(\eta^{-1}P \setminus T, \mathbb{R}^d)}^2 &\leq C\eta^{d-1} \|\nabla \widetilde{\mathcal{X}}^{2k+1}\|_{L^2(\eta^{-1}P \setminus T, \mathbb{R}^d)} \|\widetilde{\mathcal{X}}^{2k+2} - \langle \widetilde{\mathcal{X}}^{2k+2} \rangle\|_{L^2(\eta^{-1}P \setminus T)} \\ &\quad + \eta^d \|\widetilde{\mathcal{X}}^{2k}\|_{L^2(\eta^{-1}P \setminus T)} \|\widetilde{\mathcal{X}}^{2k+2}\|_{L^2(\eta^{-1}P \setminus T)} \\ &\leq (C'\eta^{d-2} + C'') \|\nabla \widetilde{\mathcal{X}}^{2k+2}\|_{L^2(\eta^{-1}P \setminus T, \mathbb{R}^d)}, \end{aligned} \quad (7.3.81)$$

$$\begin{aligned} \|\nabla \widetilde{\mathcal{X}}^{2k+3}\|_{L^2(\eta^{-1}P \setminus T, \mathbb{R}^d)}^2 &\leq 2\eta \|\nabla \widetilde{\mathcal{X}}^{2k+2}\|_{L^2(\eta^{-1}P \setminus T, \mathbb{R}^d)} \|\widetilde{\mathcal{X}}^{2k+3} - \langle \widetilde{\mathcal{X}}^{2k+3} \rangle\|_{L^2(\eta^{-1}P \setminus T)} \\ &\quad + \eta^d \|\widetilde{\mathcal{X}}^{2k+1}\|_{L^2(\eta^{-1}P \setminus T)} \|\widetilde{\mathcal{X}}^{2k+3}\|_{L^2(\eta^{-1}P \setminus T)} \\ &\leq C \|\nabla \widetilde{\mathcal{X}}^{2k+3}\|_{L^2(\eta^{-1}P \setminus T, \mathbb{R}^d)}. \end{aligned} \quad (7.3.82)$$

This implies (7.3.70) at rank  $k+1$ . Still up to extracting a subsequence, there exists tensor functions  $\widehat{\Psi}^{2k+2}$ ,  $\widehat{\Psi}^{2k+3}$  and constant tensors  $c^{2k+2}$  and  $c^{2k+3}$  such that

$$\begin{aligned} \langle \widetilde{\mathcal{X}}^{2k+2} \rangle &\rightarrow c^{2k+2} \quad \text{and} \quad \langle \widetilde{\mathcal{X}}^{2k+3} \rangle \rightarrow c^{2k+3} \quad \text{when } \eta \rightarrow 0 \\ \widetilde{\mathcal{X}}^{2k+2} &\rightharpoonup \widehat{\Psi}^{2k+2} \quad \text{and} \quad \widetilde{\mathcal{X}}^{2k+3} \rightharpoonup \widehat{\Psi}^{2k+3} \quad \text{weakly in } H_{loc}^1(\mathbb{R}^d \setminus T) \quad \text{when } \eta \rightarrow 0. \end{aligned}$$

The very same previous arguments yield eventually  $\widehat{\Psi}^{2k+2} = c^{2k+2} \Psi$  and  $\widehat{\Psi}^{2k+3} = c^{2k+3} \Psi$ . Finally, setting  $\Phi = 1$  in (7.3.79) and (7.3.80) entails

$$\begin{aligned} -\int_{\partial T} \nabla \widehat{\Psi}^{2k+2} \cdot \mathbf{n} ds &= \lim_{\eta \rightarrow 0} -\int_{\partial T} \nabla \widetilde{\mathcal{X}}^{2k+2} \cdot \mathbf{n} ds = \lim_{\eta \rightarrow 0} \langle \widetilde{\mathcal{X}}^{2k} \rangle \otimes I = \frac{1}{|\Psi^*|^{k+1}} I^{2k+2} \\ -\int_{\partial T} \nabla \widehat{\Psi}^{2k+3} \cdot \mathbf{n} ds &= \lim_{\eta \rightarrow 0} -\int_{\partial T} \nabla \widetilde{\mathcal{X}}^{2k+3} \cdot \mathbf{n} ds = \lim_{\eta \rightarrow 0} \langle \widetilde{\mathcal{X}}^{2k+1} \rangle \otimes I = 0. \end{aligned}$$

This implies  $c^{2k+2} = I^{2k+2}/|\Psi^*|^{k+2}$  and  $c^{2k+3} = 0$ , which concludes the proof.

□

**Remark 7.15.** This result seems to indicate that we may have not found the best scaling for the odd order tensors  $\mathcal{X}^{2k+1}$  since we were unable to identify a non zero weak limit. However the established bounds are sufficient for our purpose and extend to the vectorial case.

We are now able to identify the asymptotic behavior of the constant tensors  $M^k$ . Recall we already know that  $M^{2k+1} = 0$  in this scalar context (corollary 7.1).

**Corollary 7.2.** Assume  $d \geq 3$ . The following convergences hold for the tensors  $M^k$  as  $\eta \rightarrow 0$ :

$$M^0 \sim \eta^{d-2} |\Psi^*|, \tag{7.3.83}$$

$$M^2 \rightarrow -I, \tag{7.3.84}$$

$$\forall k > 1, M^{2k} = o\left(\frac{1}{\eta^{(d-2)(k-1)}}\right). \tag{7.3.85}$$

*Proof.* We replace the asymptotics of proposition 7.14 in the explicit formula (7.3.35) for the tensor  $M^k$ . (7.3.83) is a consequence of the definition  $M^0 = (\mathcal{X}^{0*})^{-1}$ . (7.3.84) is obtained by writing

$$M^2 = -((\mathcal{X}^{0*})^{-1})^2 \mathcal{X}^{2*} \sim -\frac{\eta^{2(d-2)} |\Psi^*|^2}{\eta^{2(d-2)} |\Psi^*|^2} I = -I.$$

Finally, by eliminating terms of odd orders in (7.3.35), we may write for any  $k \geq 1$ ,

$$\begin{aligned} M^{2k} &= \sum_{p=1}^{2k} \frac{(-1)^p}{(\mathcal{X}^{0*})^{p+1}} \sum_{\substack{i_1+\dots+i_p=k \\ 1 \leq i_1 \dots i_p \leq k}} \mathcal{X}^{2i_1*} \otimes \dots \otimes \mathcal{X}^{2i_p*} \\ &= \sum_{p=1}^{2k} (-1)^p \eta^{(p+1)(d-2)} |\Psi^*|^{p+1} \sum_{\substack{i_1+\dots+i_p=k \\ 1 \leq i_1 \dots i_p \leq k}} \frac{I^{2i_1}}{\eta^{(d-2)(i_1+1)} |\Psi^*|^{i_1+1}} \otimes \dots \otimes \frac{I^{2i_p}}{\eta^{(d-2)(i_p+1)} |\Psi^*|^{i_p+1}} \\ &\quad + o\left(\frac{1}{\eta^{(k-1)(d-2)}}\right) \\ &= \frac{I^{2k}}{\eta^{(k-1)(d-2)} |\Psi^*|^{k-1}} \left( \sum_{p=1}^{2k} (-1)^p \sum_{\substack{i_1+\dots+i_p=k \\ 1 \leq i_1 \dots i_p \leq k}} 1 \right) + o\left(\frac{1}{\eta^{(k-1)(d-2)}}\right). \end{aligned}$$

Then (7.3.85) results from the last summation being zero:

$$\forall k > 1, \sum_{p=1}^k (-1)^p \sum_{\substack{i_1+\dots+i_p=k \\ 1 \leq i_1 \dots i_p \leq k}} 1 = 0. \tag{7.3.86}$$

There are several ways to obtain the latter formula. A rather direct argument in the spirit of the proof of proposition 7.7 is to apply the identity (7.3.36) to the power series  $1/(1-z) = \sum_{k \in \mathbb{N}} z^k$  which yields

$$1 - z = 1 + \sum_{k=1}^{+\infty} \left( \sum_{p=1}^k (-1)^p \sum_{\substack{i_1+\dots+i_p=k \\ 1 \leq i_1 \dots i_p \leq k}} 1 \right) z^k,$$

from where (7.3.86) follows by identifying the powers in  $z^k$ . □

**Remark 7.16.** We can retrieve formally in corollary 7.2 the different classical asymptotic regimes for the perforated problem (7.3.1) depending on the size of the hole  $a_\varepsilon = \eta\varepsilon$  [13, 11, 103, 273]:

1. If  $1 \geq a_\varepsilon \gg \varepsilon^{d/(d-2)}$ , that is  $\eta$  goes to zero at a smaller rate than  $\varepsilon^{2/(d-2)}$ , then the zeroth order term  $\varepsilon^{-2} M^0 \sim \varepsilon^{-2} \eta^{d-2} |\Psi^*|$  is dominant in the infinite order homogenized equation (7.3.34), which is analogous to the ‘‘Darcy’’ regime

$$\varepsilon^{-2} M^0 u_\varepsilon^* \simeq f.$$

2. If  $a_\varepsilon = \varepsilon^{d/(d-2)}$ , the coefficients of the infinite order homogenized equation (7.3.34) converge to those of

$$-\Delta u^* + \Psi^* u^* = f.$$

This corresponds to the “Brinkman” regime with the well-known “strange” reaction term [103].

3. If  $a_\varepsilon = o(\varepsilon^{d/(d-2)})$ , that is the size of the hole  $\eta$  goes to zero at a rate faster than  $\varepsilon^{2/(d-2)}$ , then there is no strange term and the coefficients of (7.3.34) converge to those of the initial Poisson problem (7.3.1):

$$-\Delta u^* = f.$$

### 7.3.6 Simplifications for the tensors $\mathcal{X}^{k*}$ and $M^k$ in case of symmetries

In this last subsection, we analyze how the homogenized tensors  $\mathcal{X}^{k*}$  and  $M^k$  simplify when the obstacle  $\eta T$  is symmetric with respect to the cell axes. This is a classical step in the theory of homogenization; our methodology follows e.g. section 6 in [35].

In all what follows, we denote by  $S := (S_{ij})_{1 \leq i, j \leq d}$  an arbitrary orthogonal symmetry (satisfying  $S = S^T$  and  $SS = I$ ). We shall in corollary 7.4 below specialize  $S$  to either of the two kinds of following cell symmetries:

- for  $1 \leq l \leq d$ , we consider the symmetry  $S^l$  with respect to the hyperplane orthogonal to  $e_l$ :

$$S^l := I - 2e_l e_l^T; \quad (7.3.87)$$

- for  $1 \leq m \neq l \leq d$ , we consider the symmetry  $S^{lm}$  with respect to the diagonal hyperplane orthogonal to  $e_l - e_m$ :

$$S^{lm} := I - e_l e_l^T - e_m e_m^T + e_l e_m^T + e_m e_l^T. \quad (7.3.88)$$

Recall the Laplace operator is invariant under orthogonal symmetries  $S$  (satisfying  $SS = I$  and  $S^T = S$ ): for any smooth scalar field  $\mathcal{X}$ ,

$$-\Delta(\mathcal{X} \circ S) = -(\Delta \mathcal{X}) \circ S. \quad (7.3.89)$$

**Proposition 7.15.** *If the cell  $Y = P \setminus \eta T$  is invariant with respect to a symmetry  $S$ , i.e.  $S(Y) = Y$ , then the following identity holds for the cell solutions  $\mathcal{X}^k$  of eqn. (7.3.15):*

$$\mathcal{X}_{i_1 \dots i_k}^k \circ S = S_{i_1 j_1} \dots S_{i_k j_k} \mathcal{X}_{j_1 \dots j_k}^k. \quad (7.3.90)$$

*Proof.* The result is proved by induction on  $k$ . For  $k = 0$ , it holds

$$-\Delta_{yy}(\mathcal{X}^0 \circ S) = 1 \circ S = 1$$

and the symmetry of  $Y$  implies that  $\mathcal{X}^0 \circ S$  also satisfies the boundary conditions of (7.4.7). This implies  $\mathcal{X}^0 \circ S = \mathcal{X}^0$ . For  $k = 1$ , we write

$$-\Delta_{yy}(\mathcal{X}_{i_1}^1 \circ S) = 2(\partial_{i_1} \mathcal{X}^0) \circ S = 2\partial_{j_1}(\mathcal{X}^0 \circ S)S_{i_1 j_1} = 2\partial_{j_1} \mathcal{X}^0 S_{i_1 j_1}.$$

This implies similarly  $\mathcal{X}_{i_1}^1 \circ S = S_{i_1 j_1} \mathcal{X}_{j_1}^1$ . Finally, if the result holds till rank  $k + 1$  with  $k \geq 0$ , then

$$\begin{aligned} -\Delta_{yy}(\mathcal{X}_{i_1 \dots i_{k+2}}^{k+2} \circ S) &= 2(\partial_{i_{k+2}} \mathcal{X}_{i_1 \dots i_{k+1}}^{k+1}) \circ S + \delta_{i_{k+1} i_{k+2}} \mathcal{X}_{i_1 \dots i_k}^k \circ S \\ &= 2S_{i_{k+2} j_{k+2}} \partial_{j_{k+2}}(\mathcal{X}_{i_1 \dots i_{k+1}}^{k+1} \circ S) + S_{i_{k+1} j_{k+1}} S_{i_{k+2} j_{k+2}} \delta_{j_{k+1} j_{k+2}} \mathcal{X}_{i_1 \dots i_k}^k \circ S \\ &= -S_{i_1 j_1} \dots S_{i_k j_k} \Delta_{yy} \mathcal{X}_{j_1 \dots j_{k+2}}^{k+2}, \end{aligned}$$

whence the result at rank  $k + 2$ . □

**Corollary 7.3.** *If the cell  $Y = P \setminus \eta T$  is invariant with respect to a symmetry  $S$ , then the constant tensors  $\mathcal{X}^{k*}$  and  $M^k$  satisfy, for any set of indices  $1 \leq i_1, \dots, i_k \leq d$ :*

$$\mathcal{X}_{i_1 \dots i_k}^{k*} = S_{i_1 j_1} \dots S_{i_k j_k} \mathcal{X}_{j_1 \dots j_k}^{k*} \quad (7.3.91)$$

$$M_{i_1 \dots i_k}^k = S_{i_1 j_1} \dots S_{i_k j_k} M_{j_1 \dots j_k}^k \quad (7.3.92)$$

with implicit summation over the repeated indices  $j_1 \dots j_k$ .

*Proof.* Equality (7.3.91) results from the previous proposition and the following change of variables:

$$\mathcal{X}_{i_1 \dots i_k}^{k*} = \int_Y \mathcal{X}_{i_1 \dots i_k}^k dy = \int_Y \mathcal{X}_{i_1 \dots i_k}^k \circ S dy.$$

Equality (7.3.92) can be obtained by applying (7.3.91) in the formula (7.3.35).  $\square$

**Corollary 7.4.** 1. If the cell  $Y$  is symmetric with respect to all cell axes  $\mathbf{e}_l$ , i.e.  $S^l(Y) = Y$  for any  $1 \leq l \leq d$ , then

$$\mathcal{X}_{i_1 \dots i_k}^{k*} = 0 \text{ and } M_{i_1 \dots i_k}^k = 0$$

whenever there exists a number  $r$  occurring with an odd multiplicity in the indices  $i_1 \dots i_k$ , i.e. whenever

$$\exists r \in \{1, \dots, d\}, \text{ Card}\{j \in \{1, \dots, k\} \mid i_j = r\} \text{ is odd.}$$

2. If the cell  $Y$  is symmetric with respect to all diagonal axes orthogonal to  $(\mathbf{e}_l - \mathbf{e}_m)$ , i.e.  $S^{l,m}(Y) = Y$  for any  $1 \leq l < m \leq d$ , then for any permutation  $\sigma \in \mathfrak{S}_d$ ,

$$\mathcal{X}_{\sigma(i_1) \dots \sigma(i_k)}^{k*} = \mathcal{X}_{i_1 \dots i_k}^{k*}.$$

$$M_{\sigma(i_1) \dots \sigma(i_k)}^k = M_{i_1 \dots i_k}^k.$$

*Proof.* 1. The symmetry  $S^l$  is a diagonal matrix satisfying  $S^l \mathbf{e}_l = -\mathbf{e}_l$  and  $S^l \mathbf{e}_q = \mathbf{e}_q$  for  $q \neq l$ . Hence, replacing  $S$  by  $S^l$  in (7.3.91), it holds

$$\mathcal{X}_{i_1 \dots i_k}^{k*} = (-1)^{\delta_{i_1 l} + \dots + \delta_{i_k l}} \mathcal{X}_{i_1 \dots i_k}^{k*},$$

which implies the result.

2. Applying (7.3.91) to the symmetry  $S^{l,m}$  yields the result for  $\sigma = \tau$  where  $\tau$  is the transposition exchanging  $l$  and  $m$ . Since this holds for any transposition, this implies the statement for any permutation  $\sigma$ .  $\square$

Let us illustrate how the previous corollary translates for the tensors  $M^2$  and  $M^4$ :

- if  $Y$  is symmetric with respect to the cell axes  $(\mathbf{e}_l)_{1 \leq l \leq d}$ , then only the coefficients of the form  $M_{ii}^2$ ,  $M_{iijj}^4$ ,  $M_{iiii}^4$  with  $1 \leq i, j \leq d$  and  $i \neq j$  are non zeros (in particular  $M^2$  is diagonal).
- if in addition  $Y$  is symmetric with respect to the hyperplane orthogonal to  $\mathbf{e}_l - \mathbf{e}_m$ , then these coefficients do not depend on the chosen distinct indices  $i$  and  $j$ :  $M^2$  is a multiple of the identity and  $M^4$  reduces to two effective coefficients.

## 7.4 HIGH ORDER HOMOGENIZATION FOR THE PERFORATED ELASTICITY SYSTEM

This section extends the previous analysis to the elasticity system. For a periodic, smooth, *vectorial* right-hand side  $\mathbf{f} \in C^\infty(D, \mathbb{R}^d)$ , we consider now  $\mathbf{u}_\varepsilon$  to be the solution of the elasticity system

$$\begin{cases} -\operatorname{div}(A \nabla \mathbf{u}_\varepsilon) = \mathbf{f} \text{ in } D_\varepsilon \\ \mathbf{u}_\varepsilon = 0 \text{ on } \partial\omega_\varepsilon \\ \mathbf{u}_\varepsilon \text{ is } D \text{ periodic.} \end{cases} \quad (7.4.1)$$

The displacement  $\mathbf{u}_\varepsilon$  is assumed to be zero on  $\partial\omega_\varepsilon$ : physically, (7.4.1) models a mechanical structure clamped with many screws.

In (7.4.1),  $A := (A_{ijkl})_{1 \leq i, j, k, l \leq d}$  is a 4th order tensor written in standard notation, which means that for any  $d \times d$  matrix  $P := (P_{kl})_{1 \leq k, l \leq d}$ ,  $AP$  is the  $d \times d$  matrix defined by  $(AP)_{ij} := A_{ijkl} P_{kl}$ . Hence, it holds

$$A \nabla \mathbf{u} = A_{ijkl} \partial_l u_k$$

for a vector field  $\mathbf{u} := (u_k)_{1 \leq k \leq d} \in \mathbb{R}^d$ . Note that this notation differs slightly from our notational conventions of section 7.2.4, whereby  $A \equiv A^2$  would rather be considered a second order matrix valued

tensor, see (7.4.40) below. The fourth order tensor  $A$  is assumed symmetric, positive definite in the following sense:

$$\begin{aligned} \forall P, Q \in \mathbb{R}^{d \times d}, P : AQ = AP : Q \\ \exists \nu > 0, \forall P \in \mathbb{R}^{d \times d}, P : AP \geq \nu P : P, \end{aligned}$$

where “:” denotes the Frobenius scalar product on  $d \times d$  matrices:  $P : Q = P_{ij}Q_{ij}$ .

In a more classical setting,  $A$  can be considered to be the Hooke’s tensor for isotropic homogeneous materials, defined for two given Lamé constants  $\lambda, \mu > 0$  as follows:

$$A \nabla \mathbf{v} = \mu(\nabla \mathbf{v} + \nabla \mathbf{v}^T) + \lambda \text{Tr}(\nabla \mathbf{v}), \tag{7.4.2}$$

although we do not rely on this assumption in what follows.

Hereafter, we reproduce and extend the analysis of section 7.3 to the elasticity system (7.4.1). The main difference with the analysis of section 7.3 lies in the vectorial nature of the problem. This reflects in the following substantial changes with respect to the scalar context:

- tensors  $\mathcal{X}^k, M^k, N^k$  become matrix valued (e.g.  $\mathcal{X}_{i_1 \dots i_k}^k$  is a  $d \times d$  matrix for any set of indices  $i_1 \dots i_k$ );
- odd order tensors  $\mathcal{X}^{2k+1*}$  and  $M^{2k+1}$  do not, in general, vanish anymore (we provide numerical evidence for this fact in section 7.4.5 below) but are antisymmetric matrix valued. However, these odd order tensors vanish in the case where the obstacle  $\eta T$  is symmetric with respect to the cell axes;
- even order tensors  $\mathcal{X}^{2k*}$  and  $M^{2k}$  are symmetric matrix valued;
- low volume fraction asymptotics remain valid, and it turns out that *strange* odd order operators  $\varepsilon^{2k+1} M^{2k+1} \cdot \nabla^{2k+1}$  disappear as  $\eta \rightarrow 0$ .

### 7.4.1 Formal infinite order two-scale asymptotics and matrix valued tensors $\mathcal{X}^k$

Again, we consider the traditional two-scale ansatz

$$\mathbf{u}_\varepsilon = \sum_{i=0}^{+\infty} \varepsilon^{i+2} \mathbf{u}_i(x, x/\varepsilon) \tag{7.4.3}$$

where the unknowns  $\mathbf{u}_i(x, y)$  are now vector functions  $\mathbf{u}_i : D \times Y \rightarrow \mathbb{R}^d D$  and  $P$ -periodic. Similarly as before and this is not compulsory, we assume that  $\mathbf{f}$  itself can depend on  $\varepsilon$  in the form of a formal power series in  $\varepsilon$  featuring only non-oscillating terms:

$$\forall x \in D, \mathbf{f}(x) = \sum_{i=0}^{+\infty} \varepsilon^i \mathbf{f}^i(x). \tag{7.4.4}$$

Inserting formally the ansatz (7.4.3) into the elasticity system (7.4.1) yields the cascade of equations

$$\begin{cases} \mathcal{A}_{yy} \mathbf{u}_{i+2} = -\mathcal{A}_{xy} \mathbf{u}_{i+1} - \mathcal{A}_{xx} \mathbf{u}_i + \mathbf{f}_{i+2}, & \forall i \geq -2 \\ \mathbf{u}_{-2}(x, y) = \mathbf{u}_{-1}(x, y) = 0, \\ \mathbf{u}_i(x, \cdot) = 0 \text{ on } \partial(\eta T) \forall i \geq 0, \\ \mathbf{u}_i(x, \cdot) \text{ is } P\text{-periodic } \forall i \geq 0, \end{cases} \tag{7.4.5}$$

where the operators  $\mathcal{A}_{xx}, \mathcal{A}_{xy}$  and  $\mathcal{A}_{yy}$  are defined by

$$\mathcal{A}_{yy} := -\text{div}_y(A \nabla_y \cdot), \quad \mathcal{A}_{xy} := -\text{div}_x(A \nabla_y \cdot) - \text{div}_y(A \nabla_x \cdot), \quad \mathcal{A}_{xx} := -\text{div}_x(A \nabla_x \cdot).$$

We introduce an additional differential operator  $\partial_l^A$  acting on differentiable vector fields  $\mathbf{v}$  as follows:

$$\partial_l^A \mathbf{v} := \frac{1}{2}(A \nabla \mathbf{v} \cdot \mathbf{e}_l + \text{div}(A[\mathbf{v} \mathbf{e}_l^T])). \tag{7.4.6}$$



The operator  $\partial_l^A$  behaves like a partial derivative in the direction  $\mathbf{e}_l$ . It can be verified that  $\partial_l^A$  is an antisymmetric operator: if  $\mathbf{v}$  and  $\mathbf{w}$  are two smooth  $P$ -periodic functions vanishing on the hole  $\partial(\eta T)$ , then it holds

$$\begin{aligned} \int_Y \mathbf{v} \cdot \partial_l^A \mathbf{w} dy &= \frac{1}{2} \int_Y (\mathbf{v} \cdot A \nabla \mathbf{w} - \mathbf{w} \cdot A \nabla \mathbf{v}) \cdot \mathbf{e}_l dy \\ &= - \int_Y \mathbf{w} \cdot \partial_l^A \mathbf{v} dy. \end{aligned}$$

The system (7.4.5) can be solved by introducing matrix valued tensors  $\mathcal{X}^k : Y \rightarrow \mathbb{R}^{d \times d}$ : they are defined from their column vectors  $(\mathcal{X}_j^k)_{1 \leq j \leq d}$  by a recurrence analogous to (7.3.15):

$$\left\{ \begin{array}{l} \mathcal{A}_{yy} \mathcal{X}_j^0 = \mathbf{e}_j \text{ in } Y, \\ \mathcal{A}_{yy} \mathcal{X}_j^1 = 2\partial_l^A \mathcal{X}_j^0 \otimes \mathbf{e}_l \text{ in } Y, \\ \mathcal{A}_{yy} \mathcal{X}_j^{k+2} = 2\partial_l^A \mathcal{X}_j^{k+1} \otimes \mathbf{e}_l + (A[\mathcal{X}_j^k \mathbf{e}_l^T] \cdot \mathbf{e}_m) \mathbf{e}_l \otimes \mathbf{e}_m \text{ in } Y, \quad \forall k \geq 0, \forall 1 \leq j \leq d. \\ \mathcal{X}_j^k = 0 \text{ on } \partial(\eta T), \\ \mathcal{X}_j^k \text{ is } P\text{-periodic} \end{array} \right. \quad (7.4.7)$$

The coefficients of the matrix valued tensor  $\mathcal{X}^k := [\mathcal{X}_1^k \dots \mathcal{X}_d^k]$  are then given by

$$\mathcal{X}_{ij}^k = \mathcal{X}_j^k \cdot \mathbf{e}_i, \quad 1 \leq i, j \leq d.$$

Finally, we associate to these tensors a family of differential operators  $\mathcal{X}^k \cdot \nabla^k$  defined for any smooth vector field  $\mathbf{v} := (v_q)_{1 \leq q \leq d} \in C^\infty(D, \mathbb{R}^d)$  by  $(\mathcal{X}^k \cdot \nabla^k \mathbf{v})_p := \mathcal{X}_{i_1 \dots i_k, pq}^k \partial_{i_1 \dots i_k}^k v_q$ , following the conventions of section 7.2.4.

**Remark 7.17.** With our notational conventions, the non bold symbols  $\otimes \mathbf{e}_l$  and  $\otimes \mathbf{e}_m$  refer to the tensor part of  $\mathcal{X}^k$  and indicates the occurrence of additional indices  $l$  and  $m$ , which is consistent with the definition (7.2.22). For instance the last line of (7.4.7) must be read as

$$\mathcal{A}_{yy} \mathcal{X}_{i_1 \dots i_{k+2}, j}^{k+2} = 2\partial_{i_{k+2}}^A \mathcal{X}_{i_1 \dots i_{k+1}, j}^{k+1} + A[\mathcal{X}_{i_1 \dots i_k, j}^k \mathbf{e}_{i_{k+1}}^T] \cdot \mathbf{e}_{i_{k+2}}.$$

**Remark 7.18.** There is no coupling between  $\mathcal{X}_j^k$  and  $\mathcal{X}_l^k$  for  $j \neq l$ : the columns of  $\mathcal{X}^k$  are defined from uncorrelated recurrences.

Using a methodology similar to section 7.3, we may prove

**Proposition 7.16.** *The solutions  $\mathbf{u}_k$  of the cascade of equations (7.4.5) are given by*

$$\mathbf{u}_i(x, y) = \sum_{k=0}^i \mathcal{X}^k(y) \cdot \nabla^k \mathbf{f}^{i-k}(x), \quad (7.4.8)$$

which can be formally written as the following double series product

$$\mathbf{u}_\varepsilon(x) = \sum_{i=0}^{+\infty} \varepsilon^{i+2} \mathcal{X}^i(x/\varepsilon) \cdot \nabla^i \mathbf{f}(x) = \left( \sum_{i=0}^{+\infty} \varepsilon^{i+2} \mathcal{X}^i(x/\varepsilon) \cdot \nabla^i \right) \left( \sum_{i=0}^{+\infty} \varepsilon^i \mathbf{f}^i(x) \right), \quad x \in D_\varepsilon.$$

*Proof.* The result is obtained as in proposition 7.1, by using the following two identities which hold for a scalar function  $f \in C^\infty(D)$  and a vector field  $\mathbf{a} \in C^\infty(Y, \mathbb{R}^d)$ :

$$\begin{aligned} \mathcal{A}_{xy}(f(x)\mathbf{a}(y)) &= -2\partial_{x_j} f(x) \partial_{y_j}^A \mathbf{a}(y), \\ \mathcal{A}_{xx}(f(x)\mathbf{a}(y)) &= -\partial_{x_l x_m}^2 f(x) A[\mathbf{a}(y) \mathbf{e}_l^T] \cdot \mathbf{e}_m. \end{aligned}$$

□

Following section 7.3, we denote by  $\mathcal{X}^{i*}$  and  $\mathbf{u}_i^*$  the averages of  $\mathcal{X}^i(y)$  and  $\mathbf{u}_i(x, y)$  with respect to the  $y$  variable:

$$\mathcal{X}^{i*} := \int_Y \mathcal{X}^i(y) dy, \quad \mathbf{u}_i^*(x) := \int_Y \mathbf{u}_i(x, y) dy, \quad \forall x \in D, \quad \forall i \in \mathbb{N}. \quad (7.4.9)$$

Our main objective is the derivation of high order homogenized equations yielding an approximation of the “infinite order” homogenized solution formally defined by

$$\mathbf{u}_\varepsilon^*(x) := \sum_{i=0}^{+\infty} \varepsilon^{i+2} \mathbf{u}_i^*(x) = \sum_{i=0}^{+\infty} \varepsilon^{i+2} \mathcal{X}^{i*} \cdot \nabla^i f(x), \quad x \in D. \quad (7.4.10)$$

Our first step is to extend the identities of [proposition 7.2](#) to the present context: we show that  $\mathcal{X}^{2k*}$  and  $\mathcal{X}^{2k+1*}$  are respectively symmetric and antisymmetric matrix valued tensors.

**Proposition 7.17.** *For any  $0 \leq p \leq k$ , the following identity holds for the matrix valued tensor  $\mathcal{X}^{k*}$ :*

$$\forall 1 \leq i, j \leq d, \mathcal{X}_{ij}^{k*} = (-1)^p \int_Y (\mathcal{A}_{yy} \mathcal{X}_i^p \cdot \mathcal{X}_j^{k-p} - A[\mathcal{X}_i^{p-1} \mathbf{e}_l^T] : [\mathcal{X}_j^{k-p-1} \mathbf{e}_m^T] \otimes e_l \otimes e_m) dy, \quad (7.4.11)$$

where the convention  $\mathcal{X}_i^{-1} = 0$  is assumed. In particular, for any  $k \geq 0$ :

- $\mathcal{X}^{2k*}$  takes values in the set of  $d \times d$  symmetric matrices:

$$\forall 1 \leq i, j \leq d, \mathcal{X}_{ij}^{2k*} = \mathcal{X}_{ji}^{2k*}.$$

- $\mathcal{X}^{2k+1*}$  takes values in the set of  $d \times d$  antisymmetric matrices:

$$\forall 1 \leq i, j \leq d, \mathcal{X}_{ij}^{2k+1*} = -\mathcal{X}_{ji}^{2k+1*}.$$

Furthermore,  $\mathcal{X}^{2k*}$  and  $\mathcal{X}^{2k+1*}$  depend only on  $(\mathcal{X}_j^k)_{1 \leq j \leq d}$  and  $(\mathcal{X}_j^{k-1})_{1 \leq j \leq d}$  as implied by the following identities:

$$\mathcal{X}_{ij}^{2k*} = (-1)^k \int_Y (A \nabla \mathcal{X}_i^k : \nabla \mathcal{X}_j^k - A[\mathcal{X}_i^{k-1} \mathbf{e}_l^T] : [\mathcal{X}_j^{k-1} \mathbf{e}_m^T] \otimes e_l \otimes e_m) dy. \quad (7.4.12)$$

$$\begin{aligned} \mathcal{X}_{ij}^{2k+1*} &= (-1)^k \int_Y (\mathcal{X}_i^k \cdot A \nabla \mathcal{X}_j^k - \mathcal{X}_j^k \cdot A \nabla \mathcal{X}_i^k) \cdot e_l \otimes e_l dy \\ &\quad + (-1)^k \int_Y (A[\mathcal{X}_j^{k-1} \mathbf{e}_l^T] : [\mathcal{X}_i^k \mathbf{e}_m^T] - A[\mathcal{X}_i^{k-1} \mathbf{e}_l^T] : [\mathcal{X}_j^k \mathbf{e}_m^T]) \otimes e_l \otimes e_m dy. \end{aligned} \quad (7.4.13)$$

**Remark 7.19.** 1. There is no reason for the odd order tensors  $\mathcal{X}^{2k+1*}$  being zero for all  $k \geq 0$ . Although it does not seem straightforward to exhibit a particular shape of hole for which it could be proved that  $\mathcal{X}^{2k+1*} \neq 0$ , numerical experiments tend to confirm this conjecture (see [section 7.4.5](#)). This fact is an important difference with the scalar case.

2. We shall see in [section 7.4.4](#) below that odd order tensors  $\mathcal{X}^{2k+1*}$  and  $M^{2k+1}$  vanish under symmetry assumptions of the periodic pattern  $Y$ .
3. Identity [\(7.4.12\)](#) is similar to [\(7.3.21\)](#):  $\mathcal{X}^{2k*}$  is a difference between two positive tensors.

It can also be shown that the tensors  $\mathcal{X}_j^k$  take linearly independent vector values in  $\mathbb{R}^d$  (in particular, there must exist non zero components of  $\mathcal{X}_j^k$  for any  $1 \leq j \leq d$ ):

**Proposition 7.18.** *For any  $1 \leq j \leq d$ , the following identity holds:*

$$\forall k \geq 0, \mathcal{A}_{yy}(\partial_{i_1 \dots i_k}^k \mathcal{X}_{i_1 \dots i_k, j}^k) = (-1)^k (k+1) \mathbf{e}_j.$$

*Proof.* The result is obtained by adapting the proof of [proposition 7.3](#) and by using the following formulas which hold for any smooth vector field  $\mathcal{X} \in C^\infty(Y)$ :

$$-\partial_j \partial_j^A \mathcal{X} = \mathcal{A}_{yy} \mathcal{X}, \quad (7.4.14)$$

$$-\partial_{i_m}^2 (A[\mathcal{X} \mathbf{e}_l^T] \cdot \mathbf{e}_m) = \mathcal{A}_{yy} \mathcal{X}.$$

□

### 7.4.2 High order homogenized equations: tensors $M^k$ , $N^k$ , $\mathbb{B}_K$ and $\mathbb{D}_K$

We now derive high order homogenized equations for the perforated elasticity problem (7.4.1) following the methodology described in the beginning of section 7.3.2.

Let us start by recalling the matrix  $\mathcal{X}^{0*}$  is symmetric positive definite as a consequence of the identity

$$\mathcal{X}_{ij}^0 = \int_Y A \nabla \mathbf{x}_i^0 : \nabla \mathbf{x}_j^0 dy, \quad 1 \leq i, j \leq d,$$

and of the linear independence of the vector fields  $(\mathbf{x}_j^0)_{1 \leq j \leq d}$  (proposition 7.18).

**Proposition 7.19.** *We define a family  $(M^k)$  of tensors of order  $k$  by induction as follows:*

$$\begin{cases} M^0 = (\mathcal{X}^{0*})^{-1} \\ M^k = -(\mathcal{X}^{0*})^{-1} \sum_{p=0}^{k-1} \mathcal{X}^{k-p*} \otimes M^p \quad \forall k \geq 1. \end{cases} \quad (7.4.15)$$

Then the source terms  $\mathbf{f}_i$  (eqn. (7.4.4)) are given in terms of the averaged ansatz terms  $\mathbf{u}_i^*$  through the following identity:

$$\forall i \geq 0, \mathbf{f}_i(x) = \sum_{k=0}^i M^k \cdot \nabla^k \mathbf{u}_{i-k}^*(x). \quad (7.4.16)$$

Recognizing a Cauchy product, this formula can be rewritten formally as the following “infinite order homogenized equation” for the averaged ansatz  $\mathbf{u}_\varepsilon^*$  (eqn. (7.4.10)):

$$\sum_{i=0}^{+\infty} \varepsilon^{i-2} M^k \cdot \nabla^k \mathbf{u}_\varepsilon^* = \mathbf{f}. \quad (7.4.17)$$

The explicit formula for the tensors  $M^k$  of proposition 7.7 extends to the vectorial setting:

**Proposition 7.20.** *For any  $k \geq 1$ , the tensor  $M^k$  of (7.4.15) reads explicitly*

$$M^k = \sum_{p=1}^k (-1)^p \sum_{\substack{i_1 + \dots + i_p = k \\ 1 \leq i_1, \dots, i_p \leq k}} (\mathcal{X}^{0*})^{-1} \otimes \mathcal{X}^{i_1*} \otimes (\mathcal{X}^{0*})^{-1} \otimes \dots \otimes (\mathcal{X}^{0*})^{-1} \otimes \mathcal{X}^{i_p*} \otimes (\mathcal{X}^{0*})^{-1}. \quad (7.4.18)$$

We now introduce the tensors  $N^k$  which allow to write the criminal ansatz expressing  $\mathbf{u}_\varepsilon$  in terms of  $\mathbf{u}_\varepsilon^*$ :

**Proposition 7.21.** *Let  $N^k$  be the  $k$ -th order matrix valued tensor defined by*

$$N^k(y) := \sum_{p=0}^k \mathcal{X}^{k-p}(y) \otimes M^p, \quad y \in Y.$$

Then the terms  $\mathbf{u}_i(x, y)$  of the oscillating ansatz (7.4.3) can be written in terms of their averages  $\mathbf{u}_i^*(x)$  (eqn. (7.4.9)) as follows:

$$\forall i \geq 0, \mathbf{u}_i(x, y) = \sum_{k=0}^i N^k(y) \cdot \nabla^k \mathbf{u}_{i-k}^*(x). \quad (7.4.19)$$

Recognizing a Cauchy product, this formally rewrites in terms of the following “criminal” ansatz expressing the oscillating solution  $\mathbf{u}_\varepsilon$  in terms of its formal average  $\mathbf{u}_\varepsilon^*$ :

$$\mathbf{u}_\varepsilon(x) = \sum_{k=0}^{+\infty} \varepsilon^k N^k(x/\varepsilon) \cdot \nabla \mathbf{u}_\varepsilon^*(x), \quad x \in D_\varepsilon. \quad (7.4.20)$$

The next proposition summarizes the properties of the tensors  $N^k$  and of their columns vectors  $(\mathbf{N}_j^k)_{1 \leq j \leq d}$

**Proposition 7.22.** Let  $(\mathbf{N}_j^k)_{1 \leq j \leq d}$  be the column vectors of the matrix tensors  $N^k$ :

$$\forall 1 \leq i, j \leq d, \mathbf{N}_j^k := N^k \mathbf{e}_j.$$

The tensors  $N^k$  and vector fields  $\mathbf{N}_j^k$  satisfy:

1.  $\int_Y N^0(y) dy = I$  and  $\int_Y N^k(y) dy = 0$  for any  $k \geq 1$ ;
2. for any  $1 \leq j \leq d$ , it holds

$$\begin{cases} \mathcal{A}_{yy} \mathbf{N}_j^0 = M^0 \mathbf{e}_j \\ \mathcal{A}_{yy} \mathbf{N}_j^1 = 2\partial_l^A \mathbf{N}_j^0 \otimes e_l + M^1 \mathbf{e}_j \\ \mathcal{A}_{yy} \mathbf{N}_j^{k+2} = 2\partial_l^A \mathbf{N}_j^{k+1} \otimes e_l + (A[\mathbf{N}_j^k \mathbf{e}_l^T] \cdot \mathbf{e}_m) \otimes e_l \otimes e_m + M^{k+2} \mathbf{e}_j; \end{cases}$$

3.  $\mathcal{A}_{yy}(\partial_{i_1 \dots i_k}^k \mathbf{N}_{i_1 \dots i_k, j}^k) = (-1)^k (k+1) M^0 \mathbf{e}_j$ , for any  $k \geq 0$ ;
4. For any  $1 \leq p \leq k-1$ ,

$$M_{ij}^k = (-1)^{p+1} \int_Y (\mathcal{A}_{yy} \mathbf{N}_i^p \cdot \mathbf{N}_j^{k-p} - A[\mathbf{N}_i^{p-1} \mathbf{e}_l^T] : [\mathbf{N}_j^{k-p-1} \mathbf{e}_m^T] \otimes e_l \otimes e_m) dy.$$

In particular,  $M^{2k}$  and  $M^{2k+1}$  depend only on the tensors  $N^k$  and  $N^{k-1}$ , which depend themselves only on the first  $k+1$  tensors  $\mathcal{X}^0 \dots \mathcal{X}^k$ .

Importantly, the next proposition states that the tensors  $M^{2k}$  and  $M^{2k+1}$  are also respectively symmetric and antisymmetric matrix valued (which could also be seen from (7.4.18)):

**Corollary 7.5.** For any  $k \geq 0$ ,

- $M^{2k}$  is a symmetric matrix valued tensor, and the following identities hold:

$$M_{ij}^0 = \int_Y A \nabla \mathbf{N}_i^0 : \nabla \mathbf{N}_j^0 dy, \quad (7.4.21)$$

$$\forall k \geq 1, M_{ij}^{2k} = (-1)^{k+1} \int_Y (A \nabla \mathbf{N}_i^k : \nabla \mathbf{N}_j^k - A[\mathbf{N}_i^{k-1} \mathbf{e}_l^T] : [\mathbf{N}_j^{k-1} \mathbf{e}_m^T] \otimes e_l \otimes e_m) dy. \quad (7.4.22)$$

- $M^{2k+1}$  is an antisymmetric matrix valued tensor, and the following identities hold:

$$\begin{aligned} \forall k \geq 0, M_{ij}^{2k+1} &= (-1)^{k+1} \int_Y (\mathbf{N}_i^k \cdot A \nabla \mathbf{N}_j^k - \mathbf{N}_j^k \cdot A \nabla \mathbf{N}_i^k) \cdot \mathbf{e}_l \otimes e_l dy \\ &\quad + (-1)^{k+1} \int_Y (A[\mathbf{N}_j^{k-1} \mathbf{e}_l^T] : [\mathbf{N}_i^k \mathbf{e}_m^T] - A[\mathbf{N}_i^{k-1} \mathbf{e}_l^T] : [\mathbf{N}_j^k \mathbf{e}_m^T]) \otimes e_l \otimes e_m dy, \end{aligned}$$

where the convention  $N^{-1} = 0$  is assumed.

**Remark 7.20.** It may appear quite surprising that odd order differential operators  $\varepsilon^{2k-1} M^{2k+1} \cdot \nabla^{2k+1}$  arise in the “infinite order” homogenized equation (7.4.17), while the original operator  $-\operatorname{div}(A \nabla \cdot)$  is symmetric. In fact, there is no contradiction because the antisymmetry of  $M^{2k+1}$  compensates the one induced by odd order derivatives which makes  $M^{2k+1} \cdot \nabla^{2k+1}$  be a symmetric operator. Indeed, for two vector fields  $\mathbf{u} := (u_i)_{1 \leq i \leq d}$ ,  $\mathbf{v} = (v_i)_{1 \leq i \leq d}$ , it holds

$$\begin{aligned} \int_Y \mathbf{v} \cdot M^{2k+1} \cdot \nabla^{2k+1} \mathbf{u} dy &= \int_Y (M_{ij}^{2k+1} \cdot \nabla^{2k+1} u_j) v_i dy = \int_Y (-1)^{2k+1} (M_{ij}^{2k+1} \cdot \nabla^{2k+1} v_i) u_j dy \\ &= \int_Y (-1)^{2k+1} (-1) (M_{ji}^{2k+1} \cdot \nabla^{2k+1} v_i) u_j dy = \int_Y \mathbf{u} \cdot M^{2k+1} \cdot \nabla^{2k+1} \mathbf{v} dy. \end{aligned}$$

**Remark 7.21.** If the original tensor  $A$  acts on symmetric matrices (e.g. if  $A$  is given by the Hooke’s Law (7.4.2)), which means  $A_{ijkl} = A_{ijlk}$  for  $1 \leq i, j, k, l \leq d$ , the same property is *a priori* not necessarily verified for the homogenized tensor  $M^2$ , which would mean  $M_{i_1 i_2, p q}^2 = M_{i_1 q, p i_2}^2$  (see the formulas (7.4.18) and (7.4.13)). We interpret this physically by the fact that inhomogeneities in the obstacle may translate in  $M^2$  being sensitive to the antisymmetric part of the gradient. The same remark applies as well for the other tensors  $M^k$ : partial derivative indices cannot *a priori* be permuted with spatial indices. Such occurs however in case of symmetries of the obstacle, see the dedicated section 7.4.4.

We now derive well-posed homogenized equations of higher but finite order, which are in principle amenable to numerical computations. Following [section 7.3.3](#), recall  $\mathbf{u}_\varepsilon$  is the solution to the energy minimization problem

$$\begin{aligned} \min_{\mathbf{w} \in H^1(D_\varepsilon, \mathbb{R}^d)} \quad & J(\mathbf{u}, \mathbf{f}) := \int_D \left( \frac{1}{2} A \nabla \mathbf{u} : \nabla \mathbf{u} - \mathbf{f} \cdot \mathbf{u} \right) dx \\ \text{s.t.} \quad & \begin{cases} \mathbf{w} = 0 \text{ on } \partial\omega_\varepsilon \\ \mathbf{w} \text{ is } D\text{-periodic.} \end{cases} \end{aligned}$$

We consider truncated criminal ansatz  $\mathbf{w}_{\varepsilon, K}(\mathbf{v})$  depending on a non oscillating vector function  $\mathbf{v} \in H^{K+1}(D)$  obtained by truncation of [\(7.4.20\)](#):

$$\mathbf{w}_{\varepsilon, K}(\mathbf{v})(x) := \sum_{k=0}^K \varepsilon^k N^k(x/\varepsilon) \cdot \nabla^k \mathbf{v}(x), \quad x \in D_\varepsilon. \quad (7.4.23)$$

This allows us to consider a minimization problem for the function  $\mathbf{v}$  which is sought to approximate the formal infinite average  $\mathbf{u}_\varepsilon^*$  of [\(7.4.10\)](#):

$$\begin{aligned} \min_{\mathbf{v} \in H^{K+1}(D)} \quad & J(\mathbf{w}_{\varepsilon, K}(\mathbf{v}), \mathbf{f}) \\ \text{s.t.} \quad & \mathbf{v} \text{ is } D\text{-periodic.} \end{aligned} \quad (7.4.24)$$

Still following [\(7.3.3\)](#), we eliminate the fast variable  $x/\varepsilon$  in [\(7.4.24\)](#) by applying [lemma 7.3](#) which yields an approximate energy  $J_K^*(\mathbf{v}, \mathbf{f}, \varepsilon)$  which we now define.

For a given  $K \in \mathbb{N}$ ,  $1 \leq j \leq d$  and  $0 \leq k \leq K+1$ , we denote by  $(\widetilde{\mathbf{N}}_{p,j}^l)_{1 \leq j \leq d}$  the  $d$  vector tensors defined for  $1 \leq j \leq d$  by

$$\widetilde{\mathbf{N}}_{p,j}^l := \begin{cases} \partial_p \mathbf{N}_j^0(y) & \text{if } l = 0 \\ \partial_p \mathbf{N}_j^l(y) + \mathbf{N}_j^{l-1}(y) \otimes e_p & \text{if } 1 \leq l \leq K, \\ \mathbf{N}_j^K(y) \otimes e_p & \text{if } l = K+1. \end{cases} \quad (7.4.25)$$

Then we define a constant, bilinear matrix valued tensor  $\mathbb{B}_K^{k,l}$  of order  $k+l$ :

$$\mathbb{B}_{K,ij}^{l,m} := \int_Y A[\widetilde{\mathbf{N}}_{p,i}^l \mathbf{e}_p^T] : [\widetilde{\mathbf{N}}_{q,j}^m \mathbf{e}_q^T] dy, \quad \forall 1 \leq i, j \leq d, \quad (7.4.26)$$

where the summation over the repeated indices  $p$  and  $q$  and a tensor product in the double contraction is implicitly assumed, following the convention of [section 7.2.4](#). Recall that for a vector field  $\mathbf{v} := (v_p)_{1 \leq p \leq d}$  we denote by  $\mathbb{B}_K^{l,m} \nabla^l \mathbf{v} \nabla^m \mathbf{v}$  the contraction

$$\mathbb{B}_K^{l,m} \nabla^l \mathbf{v} \nabla^m \mathbf{v} := (\mathbb{B}_{K,pq}^{l,m})_{i_1 \dots i_l j_1 \dots j_m} \partial_{i_1 \dots i_l}^l v_p \partial_{j_1 \dots j_m}^m v_q.$$

The tensors  $\mathbb{B}^{l,m}$  allow us to formulate a new energy  $J_K^*(\mathbf{v}, \mathbf{f}, \varepsilon)$ :

$$J_K^*(\mathbf{v}, \mathbf{f}, \varepsilon) := \int_D \left( \frac{1}{2} \sum_{k,l=0}^{K+1} \varepsilon^{k+l-2} \mathbb{B}_K^{k,l} \nabla^k \mathbf{v} \nabla^l \mathbf{v} - \mathbf{f} \cdot \mathbf{v} \right) dx. \quad (7.4.27)$$

The definition of the energy  $J_K^*$  is motivated by the following result, which can be obtained by a straightforward adaptation of the proof of [proposition 7.10](#):

**Proposition 7.23.** *Assume  $\mathbf{f} \in C^\infty(D, \mathbb{R}^d)$  and  $D$ -periodic. For any smooth  $D$ -periodic vector field  $\mathbf{v} \in C^\infty(D, \mathbb{R}^d)$  and truncated ansatz  $\mathbf{w}_{\varepsilon, K}(\mathbf{v}) \in C^\infty(D_\varepsilon)$  of the form of [\(7.4.23\)](#), the following energy asymptotic holds true with  $m \geq 0$  arbitrarily high:*

$$|J(\mathbf{w}_{\varepsilon, K}(\mathbf{v}), \mathbf{f}) - J_K^*(\mathbf{v}, \mathbf{f}, \varepsilon)| \leq C_{K,m} (\|\nabla^{m+2} \mathbf{v}\|_{L^\infty(D, \mathbb{R}^d \times d)}^2 + \|\mathbf{f}\|_{H^m(D, \mathbb{R}^d)}^2).$$

Homogenized equations are obtained by considering the new minimization problem where  $J(\mathbf{w}_{\varepsilon, K}(\mathbf{v}), \mathbf{f})$  is replaced with  $J_K^*(\mathbf{v}, \mathbf{f}, \varepsilon)$  in [\(7.4.24\)](#).

**Definition 7.7.** For any  $K \in \mathbb{N}$ , we call homogenized equation of order  $2K + 2$  the Euler-Lagrange equation associated with the minimization problem

$$\min_{\substack{\mathbf{v} \in H^{K+1}(D, \mathbb{R}^d), \\ D \text{ periodic}}} J_K^*(\mathbf{v}, \mathbf{f}, \varepsilon). \quad (7.4.28)$$

This equation reads explicitly in terms of a higher order homogenized solution  $\mathbf{v}_K^* \in H^{K+1}(D, \mathbb{R}^d)$  as

$$\begin{cases} \sum_{k=0}^{2K+2} \varepsilon^{k-2} \mathbb{D}_K^k \cdot \nabla^k \mathbf{v}_K^* = \mathbf{f} \\ \mathbf{v}_K^* \text{ is } D\text{-periodic,} \end{cases} \quad (7.4.29)$$

where the constant (matrix valued) tensors  $\mathbb{D}_K^k$  are defined by the following formula for any  $0 \leq k \leq 2K + 2$ :

$$\forall 1 \leq i, j \leq d, \mathbb{D}_{K,ij}^k := \begin{cases} \sum_{K=0}^k (-1)^l \frac{1}{2} (\mathbb{B}_{K,ij}^{l,k-l} + \mathbb{B}_{K,ji}^{l,k-l}), & \text{if } k \text{ is even} \\ \sum_{K=0}^k (-1)^l \frac{1}{2} (\mathbb{B}_{K,ij}^{l,k-l} - \mathbb{B}_{K,ji}^{l,k-l}), & \text{if } k \text{ is odd,} \end{cases} \quad (7.4.30)$$

where the convention  $\mathbb{B}_K^{l,m} = 0$  whenever  $l > K + 1$  or  $m > K + 1$  is assumed.

Again, an important difference with the scalar case (eqn. (7.3.50)) is the occurrence of differential operators  $\varepsilon^{2k-1} \mathbb{D}_K^{2k+1} \cdot \nabla^{2k+1}$  of odd orders in the homogenized equation (7.4.29). The  $k$ -th order tensor  $\mathbb{D}_K^k$  is also symmetric matrix valued for even values of  $k$ , and antisymmetric otherwise, which was to be expected from remark 7.20.

**Remark 7.22.** Let us examine (7.4.29) at order  $K = 0$ . It reads

$$\mathbb{D}_0^2 \cdot \nabla^2 \mathbf{v}_0^* + \varepsilon^{-1} \mathbb{D}_0^1 \cdot \nabla \mathbf{v}_0^* + \varepsilon^{-2} \mathbb{D}_0^0 \mathbf{v}_0^* = \mathbf{f}$$

where one can verify that the tensors  $\mathbb{D}_0^0$ ,  $\mathbb{D}_0^1$  and  $\mathbb{D}_0^2$  are respectively given by

$$\mathbb{D}_0^0 = M^0, \quad (7.4.31)$$

$$\mathbb{D}_0^1 = M^1, \quad (7.4.32)$$

$$\mathbb{D}_{0,ij}^2 = (\mathcal{X}^{0*})_{ip}^{-1} (\mathcal{X}^{0*})_{jq}^{-1} \int_Y A[\mathcal{X}_p^0 e_l^T] : [\mathcal{X}_q^0 e_m^T] \otimes e_l \otimes e_m dy, \quad 1 \leq i, j \leq d. \quad (7.4.33)$$

The physical interpretation of the rather *very strange* term  $\varepsilon^{-1} \mathbb{D}_0^1 \cdot \nabla \mathbf{v}_0^*$  does not seem obvious to us.

Observing that the dominant tensor  $\mathbb{B}_K^{K+1, K+1}$  of (7.4.26) is symmetric and nonnegative, the well-posedness result of proposition 7.11 extends without difficulty:

**Proposition 7.24.** Assume further the dominant tensor  $\mathbb{D}_K^{2K+2} = (-1)^{K+1} \mathbb{B}_K^{K+1, K+1}$  is non-degenerate, that is there exists a constant  $\nu$  such that

$$\forall \boldsymbol{\xi} = \boldsymbol{\xi}_{i_1 \dots i_{K+1}, j} \in \mathbb{R}^{d^{K+1}} \times \mathbb{R}^d, \mathbb{B}_K^{K+1, K+1} \boldsymbol{\xi} \boldsymbol{\xi} \geq \nu |\boldsymbol{\xi}|^2. \quad (7.4.34)$$

Then there exists a unique solution  $\mathbf{v}_K^* \in H^{K+1}(D, \mathbb{R}^d)$  to the homogenized equation (7.4.29).

Finally, it remains to verify that (7.4.29) is a truncation in some sense of the infinite order homogenized equation (7.4.17):

**Proposition 7.25.** The first  $K+1$  homogenized coefficients of the homogenized equation (7.4.29) coincide with those of the formal infinite order homogenized equation (7.4.17):

$$\forall 0 \leq k \leq K, \mathbb{D}_K^k = M^k.$$

*Proof.* We follow the proof of [proposition 7.12](#). For  $0 \leq k, l \leq K$  and  $1 \leq i, j \leq d$ , the coefficient  $\mathbb{B}_{K,ij}^{l,k-l}$  is given by

$$\mathbb{B}_{K,ij}^{l,k-l} = \int_Y A[(\partial_p \mathbf{N}_i^l + \mathbf{N}_i^{l-1} \otimes e_p) \mathbf{e}_p^T] : [(\partial_q \mathbf{N}_j^{k-l} + \mathbf{N}_j^{k-l-1} \otimes e_q) \mathbf{e}_q^T] dy.$$

After integration by parts, the following identities hold:

$$\begin{aligned} \int_Y A[\partial_p \mathbf{N}_i^l \mathbf{e}_p^T] : [\partial_q \mathbf{N}_j^{k-l} \mathbf{e}_q^T] dy &= \int_Y A \nabla \mathbf{N}_i^l : \nabla \mathbf{N}_j^{k-l} dy = \int_Y \mathcal{A}_{yy} \mathbf{N}_i^l \cdot \mathbf{N}_j^{k-l} dy, \\ \int_Y A[\mathbf{N}_i^{l-1} \otimes e_p \mathbf{e}_p^T] : [\partial_q \mathbf{N}_j^{k-l} \mathbf{e}_q^T] dy &= - \int_Y \operatorname{div}(A[\mathbf{N}_i^{l-1} \mathbf{e}_p^T]) \cdot \mathbf{N}_j^{k-l} \otimes e_p dy. \end{aligned}$$

This allows to rewrite  $\mathbb{B}_{K,ij}^{l,k-l}$  as follows:

$$\mathbb{B}_{K,ij}^{l,k-l} = \int_Y (\mathcal{A}_{yy} \mathbf{N}_i^l - 2\partial_p^A \mathbf{N}_i^{l-1} \otimes e_p - A[\mathbf{N}_i^{l-2} \mathbf{e}_p^T] \cdot \mathbf{e}_q \otimes e_p \otimes e_q) \cdot \mathbf{N}_j^{k-l} dy + B_{ij}^{k,l} + B_{ij}^{k,l+1}$$

where  $B_{ij}^{k,l}$  is the  $k$ -th order tensor defined by

$$B_{ij}^{k,l} := \int_Y (A \nabla \mathbf{N}_i^{l-1} : [\mathbf{N}_j^{k-l}] + A[\mathbf{N}_i^{l-2} \mathbf{e}_p^T] : [\mathbf{N}_j^{k-l} \mathbf{e}_q^T] \otimes e_p \otimes e_q) dy.$$

Hence we obtain, recognizing a telescopic series:

$$\begin{aligned} \sum_{l=0}^k (-1)^l \mathbb{B}_{K,ij}^{l,k-l} &= \sum_{l=0}^{k-1} (-1)^l \int_Y (M^l \mathbf{e}_i) \cdot \mathbf{N}_j^{k-l} dy + \sum_{l=0}^k ((-1)^l B_{ij}^{k,l} - (-1)^{l+1} B_{ij}^{k,l+1}) \\ &= (-1)^k (M^k \mathbf{e}_i) \cdot \mathbf{e}_j + B_{ij}^{k,0} - (-1)^{k+1} B_{ij}^{k,k+1} \\ &= (-1)^k M_{ji}^k. \end{aligned}$$

From there, we obtain

$$\mathbb{D}_{ij}^k = \frac{1}{2} ((-1)^k M_{ji}^k + M_{ij}^k)$$

which concludes the proof since we know that both  $M^k$  and  $\mathbb{D}^k$  are symmetric for even  $k$  and antisymmetric otherwise.  $\square$

**Remark 7.23.** This result may not be optimal, i.e.  $\mathbb{D}_K^k = M^k$  may hold up to  $k = K + 1$  under some circumstances, because equality [\(7.4.32\)](#) is not explained by the previous proposition.

The proof of the error estimates in [proposition 7.13](#) for the scalar case relied essentially on the coercivity of the operator  $-\Delta$  on the space  $H^1(D)$ . The results extends therefore without difficulty to the vectorial context due to the coercivity of the operator  $-\operatorname{div}(A \nabla \cdot)$  on  $H^1(D, \mathbb{R}^d)$ :

**Proposition 7.26.** *Let  $\mathbf{v}_K^*$  the solution to the homogenized equation [\(7.4.29\)](#) of order  $2K + 2$ . There exists a constant  $C_K(\mathbf{f})$  independent of  $\varepsilon$  such that the following error estimates hold:*

$$\left\| \mathbf{u}_\varepsilon - \sum_{k=0}^K \varepsilon^k N^k(\cdot/\varepsilon) \cdot \nabla^k \mathbf{v}_K^* \right\|_{H^1(D, \mathbb{R}^d)} \leq C_K(\mathbf{f}) \varepsilon^{K+2}. \quad (7.4.35)$$

$$\left\| \mathbf{u}_\varepsilon - \sum_{k=0}^K \varepsilon^k N^k(\cdot/\varepsilon) \cdot \nabla^k \mathbf{v}_K^* \right\|_{L^2(D, \mathbb{R}^d)} \leq C_K(\mathbf{f}) \varepsilon^{K+3}. \quad (7.4.36)$$

### 7.4.3 Low volume fraction limits when the size of the obstacle tends to 0

We now extend the results of [section 7.3.5](#) to the vectorial context: we establish low volume fraction asymptotics for the tensors  $\mathcal{X}^{k*}$  and  $M^k$  when the size  $\eta$  of the obstacle ( $\eta T$ ) converges to zero. Again, we assume

$$d \geq 3, \quad (7.4.37)$$

although analogous results could be stated for  $d = 2$ . For  $1 \leq j \leq d$ , we introduce the solution  $\Psi_j$  to the exterior problem

$$\begin{cases} -\operatorname{div}(A\nabla\Psi_j) = 0 & \text{in } \mathbb{R}^d \setminus T \\ \Psi_j = 0 & \text{on } \partial T \\ \Psi_j \rightarrow \mathbf{e}_j & \text{at } \infty, \end{cases} \quad (7.4.38)$$

where the boundary condition  $\Psi_j \rightarrow \mathbf{e}_j$  at infinity means  $\Psi_j - \mathbf{e}_j \in \mathcal{D}^{1,2}(\mathbb{R}^d \setminus T)$ . We denote by  $\Psi^* := (\Psi_{ij}^*)_{1 \leq i, j \leq d}$  the matrix collecting the normal stress components:

$$\Psi_{ij}^* := \int_{\mathbb{R}^d \setminus T} A\nabla\Psi_i : \nabla\Psi_j dx = - \int_{\partial T} \mathbf{e}_j \cdot A\nabla\Psi_i \cdot \mathbf{n} ds, \quad (7.4.39)$$

where we recall that the normal  $\mathbf{n}$  is assumed to point *inward*  $T$ . In this section, we write  $A^2$  in order to refer to the 4-th order elasticity tensor  $A$  seen as a second order matrix valued tensor: for given derivative indices  $1 \leq i_1, i_2 \leq d$  and spatial indices  $1 \leq l, m \leq d$ ,

$$A_{i_1 i_2, l m}^2 := A_{l i_1 m i_2} \quad (7.4.40)$$

which also means

$$A_{ij}^2 = A_{iljm} e_l \otimes e_m = A[\mathbf{e}_i \mathbf{e}_l^T] : [\mathbf{e}_j \mathbf{e}_m^T] e_l \otimes e_m, \quad 1 \leq i, j \leq d. \quad (7.4.41)$$

**Proposition 7.14** extends as follows:

**Proposition 7.27.** *Assume  $d \geq 3$ . For any  $k \geq 0$ , and  $1 \leq j \leq d$ , let  $\tilde{\mathcal{X}}_j^{2k}$  and  $\tilde{\mathcal{X}}_j^{2k+1}$  be the rescaled vector tensors in  $\eta^{-1}P \setminus T$  defined by*

$$\forall x \in \eta^{-1}P \setminus T, \tilde{\mathcal{X}}_j^{2k}(x) := \eta^{(d-2)(k+1)} \mathcal{X}_j^{2k}(\eta x) \text{ and } \tilde{\mathcal{X}}_j^{2k+1}(x) := \eta^{(d-2)(k+1)} \mathcal{X}_j^{2k+1}(\eta x).$$

Then:

1. there exists a constant  $C > 0$  independent of  $\eta > 0$  such that:

$$\forall \eta > 0, \|\nabla \tilde{\mathcal{X}}_j^{2k}\|_{L^2(\eta^{-1}P \setminus T, \mathbb{R}^{d \times d})} \leq C \text{ and } \|\nabla \tilde{\mathcal{X}}_j^{2k+1}\|_{L^2(\eta^{-1}P \setminus T, \mathbb{R}^{d \times d})} \leq C, \quad (7.4.42)$$

2. the following convergences hold as  $\eta \rightarrow 0$ :

$$\tilde{\mathcal{X}}_i^{2k} \rightharpoonup c_{ij}^{2k} \Psi_j \text{ weakly in } H_{loc}^1(\mathbb{R}^d \setminus T, \mathbb{R}^d), \quad (7.4.43)$$

$$\tilde{\mathcal{X}}_i^{2k+1} \rightharpoonup 0 \text{ weakly in } H_{loc}^1(\mathbb{R}^d \setminus T, \mathbb{R}^d), \quad (7.4.44)$$

$$\mathcal{X}^{2k*} \sim \frac{1}{\eta^{(d-2)(k+1)}} c^{2k}, \quad (7.4.45)$$

$$\mathcal{X}^{2k+1*} = o\left(\frac{1}{\eta^{(d-2)(k+1)}}\right), \quad (7.4.46)$$

where  $c_{ij}^{2k}$  denote the coefficients of the  $2k$ -th order matrix valued tensor  $c^{2k} := (c_{ij}^{2k})_{1 \leq i, j \leq d}$  given by:

$$c^{2k} := (\Psi^*)^{-1} \otimes \overbrace{A^2 \otimes (\Psi^*)^{-1} \dots \otimes A^2 \otimes (\Psi^*)^{-1}}^{k \text{ times}}. \quad (7.4.47)$$

*Proof.* The proof follows the lines of that of **proposition 7.10**, so we content ourselves with highlighting only the main differences that are due to the vectorial setting. The results is proved by induction on  $k$ .

1. *Case  $2k$  for  $k = 0$ .* The vector valued tensor  $\tilde{\mathcal{X}}_i^0$  satisfies  $A_{yy} \tilde{\mathcal{X}}_i^0 = \eta^d \mathbf{e}_i$  in  $\eta^{-1}P \setminus T$ , hence for any test function  $\Phi \in H^1(\eta^{-1}P \setminus T, \mathbb{R}^d)$  which is  $\eta^{-1}P$ -periodic, it holds

$$\int_{\eta^{-1}P \setminus T} A\nabla \tilde{\mathcal{X}}_i^0 : \nabla \Phi dx = \eta^d \int_{\eta^{-1}P \setminus T} \Phi \cdot \mathbf{e}_i dx + \int_{\partial T} \Phi \cdot A\nabla \tilde{\mathcal{X}}_i^0 \cdot \mathbf{n} ds. \quad (7.4.48)$$



Substituting  $\Phi = \tilde{\mathcal{X}}_i^0$  in (7.4.48) and using lemma 7.7 yields  $\|\nabla \tilde{\mathcal{X}}_i^0\|_{L^2(\eta^{-1}P \setminus T, \mathbb{R}^{d \times d})} \leq C$ . Then (7.3.66) implies that  $\langle \tilde{\mathcal{X}}_i^0 \rangle$  is bounded in  $\mathbb{R}^d$  as  $\eta \rightarrow 0$ . Up to the extraction of a subsequence, there exists a matrix  $c^0 = (c_{ij}^0)_{1 \leq i, j \leq d}$  and vector fields  $(\widehat{\Psi}_i^0)_{1 \leq i \leq d}$  such that

$$\langle \tilde{\mathcal{X}}_i^0 \rangle \cdot e_j \rightarrow c_{ij}^0 \text{ and } \tilde{\mathcal{X}}_i^0 \rightharpoonup \widehat{\Psi}_i^0 \text{ weakly in } H_{loc}^1(\eta^{-1}P \setminus T, \mathbb{R}^d) \text{ as } \eta \rightarrow 0.$$

Passing (7.4.48) to the limit when  $\eta \rightarrow 0$  with a compactly supported  $\Phi \in C^\infty(\eta^{-1}P \setminus T, \mathbb{R}^d)$  yields then  $-\operatorname{div}(A \nabla \widehat{\Psi}_i^0) = 0$  in  $\mathbb{R}^d \setminus T$ . From (7.3.68) and the lower semi-continuity of the  $H_{loc}^1(\eta^{-1}P \setminus T, \mathbb{R}^d)$  norm, we know that  $\widehat{\Phi}_i^0 - c_{ij}^0 e_j$  belongs to  $\mathcal{D}^{1,2}(\mathbb{R}^d \setminus T, \mathbb{R}^d)$ , which readily implies  $\widehat{\Psi}_i^0 = c_{ij}^0 \Psi_j$ . Setting now  $\Phi = e_j$  in (7.4.48), using the continuity of the normal force (7.4.39) with respect to the weak convergence in  $H_{loc}^1(\mathbb{R}^d \setminus T, \mathbb{R}^d)$ , we obtain

$$-\int_{\partial T} e_j \cdot A \nabla \widehat{\Psi}_i^0 \cdot \mathbf{n} ds = c_{ip}^0 \Psi_{pj}^* = \delta_{ij}.$$

This means exactly  $c^0 = (\Psi^*)^{-1}$ , from where the results follows for  $\tilde{\mathcal{X}}_i^0$ .

2. *Case  $2k+1$  for  $k=0$ .* The tensor  $\tilde{\mathcal{X}}_i^1$  satisfies  $\mathcal{A}_{yy} \tilde{\mathcal{X}}_i^1 = 2\eta \partial_l^A \tilde{\mathcal{X}}_i^1 \otimes e_l$  which can be rewritten

$$\int_{\eta^{-1}P \setminus T} A \nabla \tilde{\mathcal{X}}_i^1 : \nabla \Phi dx = \int_{\eta^{-1}P \setminus T} 2\eta \partial_l^A \tilde{\mathcal{X}}_i^0 \cdot \Phi \otimes e_l dx + \int_{\partial T} \Phi \cdot A \nabla \tilde{\mathcal{X}}_i^1 \cdot \mathbf{n} ds \quad (7.4.49)$$

for any periodic  $\Phi$  in  $H^1(\eta^{-1}P \setminus T, \mathbb{R}^d)$ . Again, arguments analogous to those of the proof of proposition 7.14 imply  $\|\nabla \tilde{\mathcal{X}}_i^1\|_{L^2(\eta^{-1}P \setminus T, \mathbb{R}^{d \times d})} \leq C$  and  $|\langle \tilde{\mathcal{X}}_i^1 \rangle| \leq C$ , from where, up to the extraction of a subsequence, we obtain as above the existence of a matrix valued tensor  $c^1$  such that

$$\langle \tilde{\mathcal{X}}_i^1 \rangle \cdot e_j \rightarrow c_{ij}^1 \text{ and } \tilde{\mathcal{X}}_i^1 \rightharpoonup c_{ij}^1 \Psi_j \text{ weakly in } H_{loc}^1(\eta^{-1}P \setminus T, \mathbb{R}^d) \text{ as } \eta \rightarrow 0.$$

Passing (7.4.49) to the limit with  $\Phi = e_i$  implies in this case

$$c_{ip}^1 \Psi_{pj}^* = 0,$$

from where we infer  $c^1 = 0$  and the convergence results for  $\tilde{\mathcal{X}}_i^1$ .

3. *General case.* Assuming the results hold till rank  $k$ , it holds

$$\mathcal{A}_{yy} \tilde{\mathcal{X}}_i^{2k+2} = 2\eta^{d-1} \partial_l \tilde{\mathcal{X}}_i^{2k+1} \otimes e_l + \eta^d A[\tilde{\mathcal{X}}_i^{2k} e_l^T] \cdot e_m \otimes e_l \otimes e_m, \quad (7.4.50)$$

$$\mathcal{A}_{yy} \tilde{\mathcal{X}}_i^{2k+3} = 2\eta \partial_l \tilde{\mathcal{X}}_i^{2k+2} \otimes e_l + \eta^d A[\tilde{\mathcal{X}}_i^{2k+1} e_l^T] \cdot e_m \otimes e_l \otimes e_m. \quad (7.4.51)$$

Writing their associated variational formulations as above and adapting the arguments of the proof of proposition 7.14, we obtain  $\|\nabla \tilde{\mathcal{X}}_i^{2k+2}\|_{L^2(\eta^{-1}P \setminus T, \mathbb{R}^{d \times d})} \leq C$  and  $\|\nabla \tilde{\mathcal{X}}_i^{2k+3}\|_{L^2(\eta^{-1}P \setminus T, \mathbb{R}^{d \times d})} \leq C$ . Repeating the above arguments, we obtain, up to the extraction of a subsequence, the existence of matrix valued tensors  $c^{2k+1}$  and  $c^{2k+3}$  such that

$$\langle \tilde{\mathcal{X}}_i^{2k+2} \rangle \cdot e_j \rightarrow c_{ij}^{2k+2} \text{ and } \tilde{\mathcal{X}}_i^{2k+2} \rightharpoonup c_{ij}^{2k+2} \Psi_j \text{ weakly in } H_{loc}^1(\eta^{-1}P \setminus T, \mathbb{R}^d) \text{ as } \eta \rightarrow 0.$$

$$\langle \tilde{\mathcal{X}}_i^{2k+3} \rangle \cdot e_j \rightarrow c_{ij}^{2k+3} \text{ and } \tilde{\mathcal{X}}_i^{2k+3} \rightharpoonup c_{ij}^{2k+3} \Psi_j \text{ weakly in } H_{loc}^1(\eta^{-1}P \setminus T, \mathbb{R}^d) \text{ as } \eta \rightarrow 0.$$

The last step consists in passing the variational formulations for  $\tilde{\mathcal{X}}_i^{2k+2}$  and  $\tilde{\mathcal{X}}_i^{2k+3}$  to the limit with the test function  $\Phi = e_j$  in order to identify  $c_{ij}^{2k+2}$  and  $c_{ij}^{2k+3}$ . Performing this computation yields

$$c_{ip}^{2k+2} \Psi_{pj}^* = A[c_{ip}^{2k} e_p e_l^T] : [e_j e_m^T] \otimes e_l \otimes e_m,$$

$$c_{ip}^{2k+3} \Psi_{pj}^* = A[c_{ip}^{2k+1} e_p e_l^T] : [e_j e_m^T] \otimes e_l \otimes e_m.$$

Rewriting  $A$  into the matrix valued tensor  $A^2$  following the notation (7.4.40), the above two equations read  $c^{2k+2} \Psi^* = c^{2k} \otimes A^2$  and  $c^{2k+3} \Psi^* = c^{2k+1} \otimes A^2$ . The result follows by using (7.4.46) and (7.4.47) at rank  $k$ .

□

We are finally able to obtain the full asymptotics for the coefficients of the homogenized equation (7.4.29) when the size  $\eta$  of the hole converges to zero:

**Corollary 7.6.** *Assume  $d \geq 3$ . The following convergence hold for the matrix valued tensors  $M^k$  as  $\eta \rightarrow 0$ :*

$$M^0 \sim \eta^{d-2} \Psi^*, \quad (7.4.52)$$

$$M^1 = o(\eta^{d-2}), \quad (7.4.53)$$

$$M^2 \rightarrow -A^2, \quad (7.4.54)$$

$$\forall k > 1, M^{2k} = o\left(\frac{1}{\eta^{(d-2)(k-1)}}\right), \quad (7.4.55)$$

$$\forall k > 1, M^{2k+1} = o\left(\frac{1}{\eta^{(d-2)(k-1)}}\right). \quad (7.4.56)$$

*Proof.* We reproduce the proof of corollary 7.2 with small adaptations due to the vectorial context. (7.4.52) is immediate with  $M^0 = (\mathcal{X}^{0*})^{-1}$ . As for  $M^2$ , we write according to (7.4.18):

$$\begin{aligned} M^2 &= -(\mathcal{X}^{0*})^{-1} \otimes \mathcal{X}^{2*} \otimes (\mathcal{X}^{0*})^{-1} + (\mathcal{X}^{0*})^{-1} \otimes \mathcal{X}^{1*} \otimes (\mathcal{X}^{0*})^{-1} \otimes \mathcal{X}^{1*} \otimes (\mathcal{X}^{0*})^{-1} \\ &= -\frac{\eta^{2(d-2)}}{\eta^{2(d-2)}} \Psi^* \otimes c^2 \otimes \Psi^* + o\left(\frac{\eta^{3(d-2)}}{\eta^{2(d-2)}}\right) \\ &= -\Psi^* \otimes (\Psi^*)^{-1} \otimes A^2 \otimes (\Psi^*)^{-1} \otimes \Psi^* + o(\eta^{d-2}) \\ &= -A^2 + o(\eta^{d-2}). \end{aligned}$$

The identity (7.4.56) for  $M^{2k+1}$  is obtained by using (7.4.18), where we observe that, for any  $0 \leq p \leq 2k+1$  and indices  $1 \leq i_1 \dots i_p \leq 2k+1$  such that  $i_1 + \dots + i_p = 2k+1$ , there exists at least one odd index  $i_q$  with  $1 \leq q \leq p$ . Using (7.4.46), we obtain therefore

$$(\mathcal{X}^{0*})^{-1} \otimes \mathcal{X}^{i_1*} \otimes \dots \otimes \mathcal{X}^{i_p*} \otimes (\mathcal{X}^{0*})^{-1} = o\left(\frac{\eta^{(p+1)(d-2)}}{\eta^{(p+\lfloor i_1/2 \rfloor + \dots + \lfloor i_p/2 \rfloor)(d-2)}}\right) = o\left(\frac{1}{\eta^{(d-2)(k-1)}}\right),$$

which eventually implies (7.4.56). Finally, in order to obtain the result for  $M^{2k}$  with  $k > 1$ , we separate the summands of (7.4.18) into two categories. For a given  $p$  such that  $1 \leq p \leq 2k$  and indices  $1 \leq i_1, \dots, i_p \leq d$  such that  $i_1 + \dots + i_p = 2k$ , there are only two possibilities:

1. there exists at least one odd index  $i_q$ , in that case the above reasoning implies as well

$$(\mathcal{X}^{0*})^{-1} \otimes \mathcal{X}^{i_1*} \otimes \dots \otimes \mathcal{X}^{i_p*} \otimes (\mathcal{X}^{0*})^{-1} = o\left(\frac{1}{\eta^{(d-2)(k-1)}}\right);$$

2. all indices  $i_1 + \dots + i_p$  are even, in that case we may write, as  $\eta \rightarrow 0$ :

$$\begin{aligned} (\mathcal{X}^{0*})^{-1} \otimes \mathcal{X}^{i_1*} \otimes \dots \otimes \mathcal{X}^{i_p*} \otimes (\mathcal{X}^{0*})^{-1} &\sim \frac{1}{\eta^{(d-2)(k-1)}} \Psi^* \otimes c^{i_1} \otimes \dots \otimes \Psi^* \otimes c^{i_p} \otimes \Psi^* \\ &\sim \frac{1}{\eta^{(d-2)(k-1)}} \overbrace{A^2 \otimes (\Psi^*)^{-1} \otimes \dots \otimes A^2 \otimes (\Psi^*)^{-1}}^{k-1 \text{ times}} \otimes A^2. \end{aligned}$$

In both case, the asymptotics obtained for  $(\mathcal{X}^{0*})^{-1} \otimes \mathcal{X}^{i_1*} \otimes \dots \otimes \mathcal{X}^{i_p*} \otimes (\mathcal{X}^{0*})^{-1}$  do not depend on the choice of indices  $1 \leq i_1, \dots, i_p \leq d$ . This allows to eventually obtain, by isolating even indices  $i_1 := 2j_1, \dots, i_p := 2j_p$  in the formula (7.4.18):

$$\begin{aligned} M^{2k} &= \frac{1}{\eta^{(d-2)(k-1)}} \overbrace{A^2 \otimes (\Psi^*)^{-1} \otimes \dots \otimes A^2 \otimes (\Psi^*)^{-1}}^{k-1 \text{ times}} \otimes A^2 \left( \sum_{p=1}^{2k} (-1)^p \sum_{\substack{2j_1 + \dots + 2j_p = 2k \\ 1 \leq j_1, \dots, j_p \leq k}} 1 \right) \\ &\quad + o\left(\frac{1}{\eta^{(d-2)(k-1)}}\right). \end{aligned}$$

The result follows by using the identity (7.3.86). □

**Remark 7.24.** We have retrieved formally in the previous corollary the classical three regimes of homogenized equations [11, 105] depending on the hole size  $a_\varepsilon = \eta\varepsilon$ :

1. if  $1 \geq a_\varepsilon \gg \varepsilon^{d/(d-2)}$ , that is  $\eta$  goes to zero at a smaller rate than  $\varepsilon^{2/(d-2)}$ , then the zeroth order term  $\varepsilon^{-2}M^0 \sim \varepsilon^{-2}\eta^{d-2}\Psi^*$  is dominant in the infinite order homogenized equation (7.4.29);
2. if  $a_\varepsilon = \varepsilon^{d/(d-2)}$ , the coefficients of the infinite order homogenized equation (7.4.29) converge to those of

$$-\operatorname{div}(A\nabla\mathbf{u}^*) + \Psi^*\mathbf{u}^* = \mathbf{f},$$

which corresponds to a “Brinkman” regime with *strange* reaction term  $\Psi^*\mathbf{u}^*$ . Interestingly, the *very strange* first order term  $\varepsilon^{-1}M^1 \cdot \nabla$  disappears because  $M^1\varepsilon = o(\varepsilon^1)$ , although this term dominates the second order term  $M^2 \cdot \nabla$  for a fixed size  $\eta > 0$ . All other contributions  $\varepsilon^k M^k \cdot \nabla^k$  for  $k > 2$  vanish equally, at rates  $o(1)$  and  $o(\varepsilon)$  for respectively even and odd values of  $k$ ;

3. if  $a_\varepsilon = o(\varepsilon^{d/(d-2)})$ , then the hole are “too small” to be seen by the homogenized model and the coefficients of (7.4.29) converge to those of the elasticity problem (7.4.1) in the homogeneous domain  $D$  (without holes):

$$-\operatorname{div}(A\nabla\mathbf{u}^*) = \mathbf{f}.$$

#### 7.4.4 Simplifications for the tensors $\mathcal{X}^{k*}$ and $M^k$ in case of symmetries

This last section generalizes the results of section 7.3.6 to the elasticity context: our main result is corollary 7.7 where we obtain symmetry properties for the tensors  $\mathcal{X}^{k*}$  and  $M^k$  in case where the obstacle ( $\eta T$ ) is symmetric with respect to the cell axes.

Throughout this part, it is assumed that the elasticity tensor  $A$  commute with orthogonal symmetries  $S := (S_{ij})_{1 \leq i, j \leq d}$  in the following sense:

$$A(SMS) = S(AM)S. \quad (7.4.57)$$

This property is notably satisfied by the Hooke’s law (7.4.2). We use below the following technical lemma:

**Lemma 7.8.** *Under the assumption (7.4.57), the following identities hold for any smooth vector field  $\mathcal{X}$ :*

$$\mathcal{A}_{yy}(S\mathcal{X} \circ S) = S(\mathcal{A}_{yy}\mathcal{X}) \circ S, \quad (7.4.58)$$

$$\partial_i^A(S\mathcal{X} \circ S) = S_{ij}S(\partial_j^A\mathcal{X}) \circ S. \quad (7.4.59)$$

*Proof.* We write

$$\begin{aligned} \mathcal{A}_{yy}(S\mathcal{X} \circ S) &= -\operatorname{div}(A\nabla(S\mathcal{X} \circ S)) = -\operatorname{div}(A[S(\nabla\mathcal{X}) \circ SS]) \\ &= -S\operatorname{div}(A(\nabla\mathcal{X}) \circ S)S = -SA[(\nabla\partial_i\mathcal{X}) \circ S] \cdot S\mathbf{e}_j S_{ij} \\ &= -SA[(\nabla\partial_i\mathcal{X}) \circ S]\mathbf{e}_i = -S\operatorname{div}(A\nabla\mathcal{X}) \circ S, \end{aligned}$$

which proves (7.4.58). Equation (7.4.59) follows from

$$\begin{aligned} \partial_i^A(S\mathcal{X} \circ S) &= \frac{1}{2}(A\nabla(S\mathcal{X} \circ S) \cdot \mathbf{e}_i + \operatorname{div}(A[S\mathcal{X} \circ S\mathbf{e}_i^T])) \\ &= \frac{1}{2}(A[S(\nabla\mathcal{X}) \circ SS] \cdot \mathbf{e}_i + \operatorname{div}(SA[\mathcal{X} \circ S(\mathbf{S}\mathbf{e}_i)^T]S)) \\ &= \frac{1}{2}(SA(\nabla\mathcal{X}) \circ S \cdot S\mathbf{e}_i + SA[S_{lj}(\partial_l\mathcal{X}) \circ S(\mathbf{S}\mathbf{e}_i)^T]S \cdot \mathbf{e}_j) \\ &= \frac{1}{2}(SA(\nabla\mathcal{X}) \circ S \cdot S\mathbf{e}_i + SA[(\partial_l\mathcal{X}) \circ S(\mathbf{S}\mathbf{e}_i)^T]SS\mathbf{e}_l) \\ &= \frac{S}{2}(A(\nabla\mathcal{X}) \circ S \cdot (S_{ij}\mathbf{e}_j) + \operatorname{div}(A[\mathcal{X}(S_{ij}\mathbf{e}_j)^T]) \circ S). \end{aligned}$$

□

**Proposition 7.28.** *If the cell  $Y = P \setminus \eta T$  is invariant with respect to a symmetry  $S$ , i.e.  $S(Y) = Y$ , then the following identity holds for the vector valued tensors  $(\boldsymbol{\chi}_l^k)_{1 \leq l \leq d}$  of eqn. (7.4.7):*

$$S\boldsymbol{\chi}_{i_1 \dots i_k, l}^k \circ S = S_{i_1 j_1} \dots S_{i_k j_k} S_{lm} \boldsymbol{\chi}_{j_1 \dots j_k, m}^k. \quad (7.4.60)$$

As a consequence, the following identities hold for the constant matrix valued tensors  $\boldsymbol{\chi}^{k*}$  and  $M^k$ :

$$\boldsymbol{\chi}_{i_1 \dots i_k, lm}^{k*} = S_{i_1 j_1} \dots S_{i_k j_k} S_{lp} S_{mq} \boldsymbol{\chi}_{j_1 \dots j_k, pq}^{k*} \quad (7.4.61)$$

$$M_{i_1 \dots i_k, pq}^k = S_{i_1 j_1} \dots S_{i_k j_k} S_{lp} S_{mq} M_{j_1 \dots j_k, pq}^k. \quad (7.4.62)$$

*Proof.* We start by proving (7.4.60) by induction. It holds  $\mathcal{A}_{yy}(S\boldsymbol{\chi}_l^0 \circ S) = S\mathbf{e}_l \circ S = S_{ml}\mathbf{e}_m$ . If the cell is  $S$  symmetric, then  $S\boldsymbol{\chi}_l^0 \circ S$  also satisfies the boundary conditions of (7.4.7), we obtain therefore

$$S\boldsymbol{\chi}_l^0 \circ S = S_{lm}\boldsymbol{\chi}_m^0,$$

which is the result at rank  $k = 0$ . For  $k = 1$ , we write

$$\mathcal{A}_{yy}(S\boldsymbol{\chi}_{i_1, l}^1 \circ S) = 2S(\partial_{i_1}^A \boldsymbol{\chi}_l^0) \circ S = 2S_{i_1 j_1} \partial_{j_1}^A (S\boldsymbol{\chi}_l^0 \circ S) = 2S_{i_1 j_1} \partial_{j_1}^A (S_{lm}\boldsymbol{\chi}_m^0),$$

which implies  $S\boldsymbol{\chi}_{i_1, l}^1 \circ S = S_{i_1 j_1} S_{lm}\boldsymbol{\chi}_{j_1, m}^1$  and the result at rank  $k = 1$ . Assuming now that the result holds till rank  $k + 1$  with  $k \geq 0$ , we write

$$\mathcal{A}_{yy}(S\boldsymbol{\chi}_{i_1 \dots i_{k+2}, l}^{k+2} \circ S) = 2S(\partial_{i_{k+2}}^A \boldsymbol{\chi}_{i_1 \dots i_{k+1}, l}^{k+1}) \circ S + SA[\boldsymbol{\chi}_{i_1 \dots i_k, l}^k \circ S\mathbf{e}_{k+1}^T] \cdot \mathbf{e}_{k+2}.$$

The previous arguments and the result at rank  $k + 1$  imply the following equality for the first term:

$$2S(\partial_{i_{k+2}}^A \boldsymbol{\chi}_{i_1 \dots i_{k+1}, l}^{k+1}) \circ S = S_{i_1 j_1} \dots S_{i_{k+2} j_{k+2}} S_{lm} \partial_{j_{k+2}}^A \boldsymbol{\chi}_{j_1 \dots j_{k+1}, m}^{k+1}.$$

Furthermore, using the result at rank  $k$ , we can rewrite the second term as follows:

$$\begin{aligned} SA[\boldsymbol{\chi}_{i_1 \dots i_k, l}^k \circ S\mathbf{e}_{k+1}^T] \cdot \mathbf{e}_{k+2} &= A[S\boldsymbol{\chi}_{i_1 \dots i_k, l}^k \circ S\mathbf{e}_{i_{k+1}}^T S] \cdot S\mathbf{e}_{i_{k+2}} \\ &= A[S\boldsymbol{\chi}_{i_1 \dots i_k, l}^k \circ S(S\mathbf{e}_{i_{k+1}})^T] \cdot S\mathbf{e}_{i_{k+2}} \\ &= S_{i_1 j_1} \dots S_{i_{k+2} j_{k+2}} S_{lm} A[\boldsymbol{\chi}_{j_1 \dots j_k, m}^k \mathbf{e}_{j_{k+1}}] \cdot \mathbf{e}_{j_{k+2}}. \end{aligned}$$

The two above equalities imply the result at rank  $k + 2$ .

Eqn. (7.4.61) then follows by the following change of variables:

$$\boldsymbol{\chi}_{i_1 \dots i_k, lm}^{k*} = \int_Y \mathbf{e}_l \cdot \boldsymbol{\chi}_{i_1 \dots i_k, m}^k dy = \int_Y (S\mathbf{e}_l) \cdot (S\boldsymbol{\chi}_{i_1 \dots i_k, m}^k \circ S) dy.$$

Remarking that (7.4.61) rewrites also as

$$\boldsymbol{\chi}_{i_1 \dots i_k}^{k*} = S_{i_1 j_1} \dots S_{i_k j_k} (S\boldsymbol{\chi}_{j_1 \dots j_k}^{k*} S),$$

eqn. (7.4.62) follows by rewriting the summand of (7.4.18) as

$$\begin{aligned} (\boldsymbol{\chi}^{0*})^{-1} \otimes \boldsymbol{\chi}^{i_1*} \otimes \dots \otimes (\boldsymbol{\chi}^{0*})^{-1} \otimes \boldsymbol{\chi}^{i_p*} \otimes (\boldsymbol{\chi}^{0*})^{-1} \\ = (\boldsymbol{\chi}^{0*})^{-1} S \otimes S\boldsymbol{\chi}^{i_1*} S \otimes \dots \otimes S(\boldsymbol{\chi}^{0*})^{-1} S \otimes S\boldsymbol{\chi}^{i_p*} S \otimes S(\boldsymbol{\chi}^{0*})^{-1} \end{aligned} \quad (7.4.63)$$

before using identity (7.4.61).  $\square$

Following the proof of corollary 7.4, we obtain simplifications for the tensors  $\boldsymbol{\chi}^{2k*}$  and  $M^{2k}$  as well as the vanishing of the tensors  $\boldsymbol{\chi}^{2k+1*}$  and  $M^{2k+1}$  in case the obstacle  $\eta T$  is symmetric with respect to the cell axes. Recall (7.3.87) and (7.3.88) for the definition of  $S^l$  and  $S^{l,m}$ .

**Corollary 7.7.** *1. If the cell  $Y$  is symmetric with respect to all cell axes  $\mathbf{e}_l$ , i.e.  $S^l(Y) = Y$  for any  $1 \leq l \leq d$ , then*

$$\boldsymbol{\chi}_{i_1 \dots i_k, pq}^{k*} = 0 \text{ and } M_{i_1 \dots i_k, pq}^k = 0$$

*whenever a given integer  $1 \leq r \leq d$  occurs with an odd multiplicity in the indices  $i_1 \dots i_k, pq$ .*

*In particular, this implies  $\boldsymbol{\chi}^{2k+1*} = 0$  and  $M^{2k+1} = 0$ .*

2. If the cell  $Y$  is symmetric with respect to all diagonal axes orthogonal to  $(e_l - e_m)$ , i.e.  $S^{l,m}(Y) = Y$  for any  $1 \leq l < m \leq d$ , then for any permutation  $\sigma \in \mathfrak{S}_d$ ,

$$\begin{aligned}\mathcal{X}_{\sigma(i_1)\dots\sigma(i_k),\sigma(p)\sigma(q)}^{k*} &= \mathcal{X}_{i_1\dots i_k,pq}^{k*}, \\ M_{\sigma(i_1)\dots\sigma(i_k),\sigma(p)\sigma(q)}^k &= M_{i_1\dots i_k,pq}^k.\end{aligned}$$

Let us illustrate how the previous properties translate for the tensors  $M^2$  and  $M^4$ :

- if the cell  $Y$  is symmetric with respect to all cell axes  $(e_l)_{1 \leq l \leq d}$ , only the coefficients of the form

$$M_{ii,jj}^2, M_{ij,ij}^2, M_{ii,ii}^2$$

with  $i \neq j$  are non zero. For  $M^4$ , only the coefficients of the form

$$M_{iijj,kk}^4, M_{iijk,jk}^4, M_{iiii,jj}^4, M_{iijj,ii}^4, M_{iiij,ij}^4, M_{iiii,ii}^4$$

are non zero with distinct integers  $i, j, k$ .

- If in addition the obstacle is symmetric with respect to all diagonal axes, then the values of the above coefficients do not depend on the choice of the distinct integers  $i, j, k$ . As a result,  $M^2$  reduces to at most three coefficients (the material is said to be orthotropic), and  $M^4$  reduces to at most 6 coefficients for  $d \geq 3$ , and to 4 coefficients for  $d = 2$ .

#### 7.4.5 Appendix: numerical evidences for the “very strange” tensors $\mathcal{X}^{1*}$ and $M^1$ being nonzero

In this appendix, we report on some numerical computations performed in 2-d in order to assess the magnitude of the tensors  $\mathcal{X}^1$  and  $M^1$  for a particular (non symmetric) shape of hole  $\eta T$ . We consider a cell  $P = [0, 1] \times [0, 1]$  perforated with a “boomerang” shaped hole (see [Figure 7.4](#)), parameterized by the following system of equations

$$\begin{cases} x(t) = 0.5 + r(\cos(t) + 2\cos(2t)) \\ y(t) = 0.5 + 2r\sin(t), \end{cases} \quad (7.4.64)$$

where  $t \in [0, 2\pi]$  and  $r = 0.15$ . The cell  $Y = P \setminus (\eta T)$  is discretized into a triangular mesh (represented on [Figure 7.4a](#)) whose maximum edge size was  $\text{hmax} = 0.007$ . The elasticity tensor is given by the Hooke’s law

$$A\nabla \mathbf{u} = \mu(\nabla \mathbf{u} + \nabla \mathbf{u}^T) + \lambda \text{div}(\mathbf{u}),$$

where the Lamé parameters are  $\mu = 0.1$  and  $\lambda = 1$ .

By using the finite element method with  $\mathbb{P}1$  finite elements, we solve the first cell problem of [\(7.4.7\)](#) determining  $\mathcal{X}_j^0$  for  $j = 1, 2$ :

$$\begin{cases} \mathcal{A}_{yy} \mathcal{X}_j^0 = e_j, \\ \mathcal{X}_j^0 \text{ is } P\text{-periodic}, \end{cases} \quad j = 1, 2.$$

The numerical solutions obtained for  $\mathcal{X}_j^0$  with  $j = 1, 2$  are represented on [Figs. 7.4](#) and [7.5](#) below. They are used to compute the constant tensors  $\mathcal{X}^{0*}$ ,  $\mathcal{X}^{1*}$ ,  $M^0$  and  $M^1$  from the formulas [\(7.4.13\)](#) and [\(7.4.18\)](#):

$$\mathcal{X}_{ij}^{0*} = \int_Y \mathcal{X}_j^0 \cdot e_i \, dy, \quad (7.4.65)$$

$$\mathcal{X}_{i,ij}^{1*} = \int_Y (\mathcal{X}_i^0 \cdot A\nabla \mathcal{X}_j^0 - \mathcal{X}_j^0 \cdot A\nabla \mathcal{X}_i^0) \cdot e_i \, dy, \quad (7.4.66)$$

$$M^0 = (\mathcal{X}^{0*})^{-1}, \quad (7.4.67)$$

$$M^1 = -(\mathcal{X}^{0*})^{-1} \mathcal{X}^{1*} (\mathcal{X}^{0*})^{-1}. \quad (7.4.68)$$

These computations are performed in [FreeFEM \[183\]](#); we obtain the following numerical values:

$$\mathcal{X}^{0*} = \begin{pmatrix} 0.145 & -4.30\text{e-}7 \\ -4.30\text{e-}7 & 0.0765 \end{pmatrix} \quad (7.4.69)$$

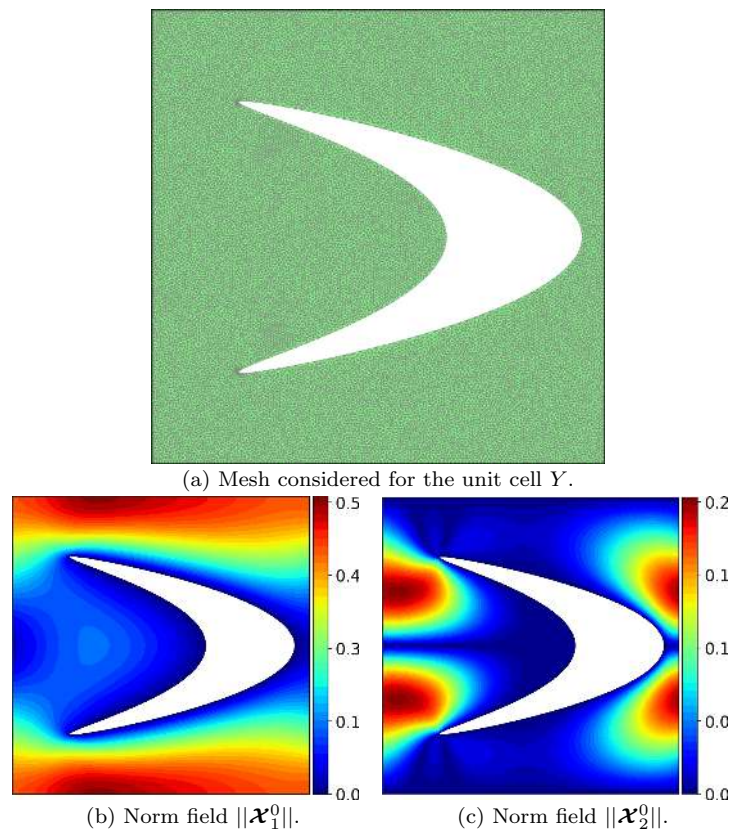


Figure 7.4: The unit cell considered for the numerical example of section 7.4.5 and the first two vector valued tensors  $\mathbf{x}_i^0$  for  $i = 1, 2$ .

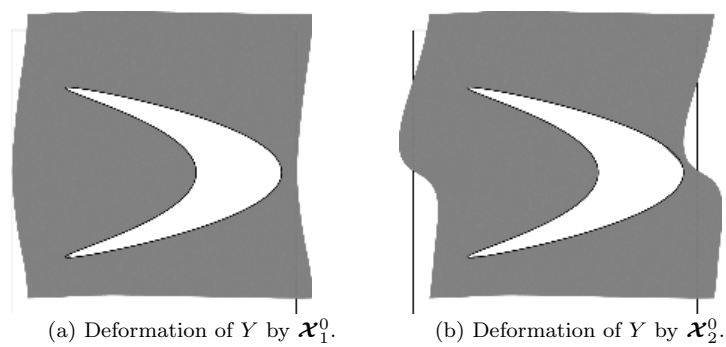


Figure 7.5: Deformation of the unit cell  $Y$  according to the first two vector valued tensors  $\mathbf{x}_i^0$  for  $i = 1, 2$ .

$$\mathcal{X}_1^{1*} = \begin{pmatrix} 1.05\text{e-}21 & 1.52\text{e-}6 \\ -1.52\text{e-}6 & -2.98\text{e-}22 \end{pmatrix}, \quad \mathcal{X}_2^{1*} = \begin{pmatrix} -7.16\text{e-}22 & 0.00237 \\ -0.00237 & -1.670\text{e-}21 \end{pmatrix}, \quad (7.4.70)$$

$$M^0 = \begin{pmatrix} 6.90 & 3.88\text{e-}05 \\ 3.88\text{e-}05 & 13.1 \end{pmatrix}, \quad (7.4.71)$$

$$M_1^1 = \begin{pmatrix} 1.074\text{e-}18 & 1.37\text{e-}4 \\ -1.37\text{e-}4 & 1.93\text{e-}18 \end{pmatrix}, \quad M_2^1 = \begin{pmatrix} 1.67\text{e-}15 & 0.214 \\ -0.214 & 3.15\text{e-}15 \end{pmatrix}, \quad (7.4.72)$$

where we have denoted by  $\mathcal{X}_i^{1*}$  and  $M_i^1$  the  $i$ -th (matrix valued) component of these tensors of order 1. These results strongly suggest that the tensor  $M^1$  is not zero (with e.g. the value  $M_{1,12}^1 = 0.214$ ). As a result, we expect the odd order differential operator  $\varepsilon^{-1}M^1 \cdot \nabla$  to have an impact when solving the high order homogenized equation (7.4.17).

## 7.5 HIGH ORDER HOMOGENIZATION FOR THE STOKES SYSTEM IN A POROUS MEDIUM

This last section is devoted to the high order homogenization of the Stokes system, which was our initial motivation:

$$\begin{cases} -\Delta \mathbf{u}_\varepsilon + \nabla p_\varepsilon = \mathbf{f} & \text{in } D_\varepsilon \\ \operatorname{div}(\mathbf{u}_\varepsilon) = 0 & \text{in } D_\varepsilon \\ \int_{D_\varepsilon} p_\varepsilon dx = 0, \\ \mathbf{u}_\varepsilon = 0 & \text{on } \partial\omega_\varepsilon \\ \mathbf{u}_\varepsilon = & \text{is } D\text{-periodic.} \end{cases} \quad (7.5.1)$$

In the sequel, we consider the following two classical assumptions for the distributions of the holes  $\omega_\varepsilon$ , following [12]:

**(H4)**  $Y \subset P$ , as a subset of the unit torus (opposite matching faces are identified) is a smooth connected set.

**(H5)** The fluid component  $D_\varepsilon$  is a smooth connected set.

**Remark 7.25.** Assumption (H5) does not necessarily imply (H4), see [7] for a counterexample. Assumption (H4) is not very restrictive and can easily be generalized to the case where the subset  $Y$  has  $m$  connected components with  $m \in \mathbb{N}$  (see section 7.5.6 for a more precise discussion). Assumption (H5) is stronger, but is also more physical. It forbids the existence of isolated fluid inclusions. Most of our derivations of homogenized tensors and homogenized equations only assume (H4). However, we rely on both assumption (H4) and (H5) when stating error bounds in section 7.5.3, because we use some technical results of [12] (namely Theorem 2.3 of this reference).

We start as previously by considering formal two-scale ansatz  $(\mathbf{u}_\varepsilon, p_\varepsilon)$  of the form

$$\mathbf{u}_\varepsilon = \sum_{i=0}^{+\infty} \varepsilon^{i+2} \mathbf{u}_i(x, x/\varepsilon), \quad p_\varepsilon(x) = \sum_{i=0}^{+\infty} \varepsilon^i (p_i^*(x) + \varepsilon p_i(x, x/\varepsilon)), \quad x \in D_\varepsilon. \quad (7.5.2)$$

where the functions  $\mathbf{u}_i(x, y)$  and  $p_i(x, y)$  are periodic in the variable  $y \in P$ . The decomposition (7.5.2) for the pressure  $p_\varepsilon$  is not the one commonly assumed in the literature, where it is more usually written only in terms of  $p_i(x, x/\varepsilon)$ . Of course, there is no loss of generality in adding the non oscillating function  $p_i^*$ , which turns to be very convenient in the sequel. In order to fix the value of  $p_i^*$  and  $p_i(x, y)$ ,  $p_i^*$  is required to have zero average with respect to the variable  $x$ , and the oscillating functions  $p_i$  are required to have zero average with respect to the variable  $y$ :

$$\int_D p_i^*(x) dx = 0, \quad \int_Y p_i(x, y) dy = 0, \quad \forall i \geq 0.$$

Assuming these conventions, our goal is to derive and analyze high order equations for the formal homogenized averages  $\mathbf{u}_\varepsilon^*$  and  $p_\varepsilon^*$  defined by

$$\mathbf{u}_\varepsilon^*(x) := \sum_{i=0}^{+\infty} \varepsilon^{i+2} \int_Y \mathbf{u}_i(x, y) dy, \quad p_\varepsilon^*(x) := \sum_{i=0}^{+\infty} \varepsilon^i p_i^*(x). \tag{7.5.3}$$

In the following, the methodologies of the previous sections 7.3 and 7.4 are applied to derive higher order homogenized equations for the perforated Stokes system (7.5.1). Most of the ingredients remain analogous to the case of the elasticity system, however additional technical difficulties occur in the proof of error estimates due to the incompressibility constraint  $\operatorname{div}(\mathbf{u}_\varepsilon) = 0$ .

This part is organized as follows. In section 7.5.1, we introduce cell problems and their solution tensors  $(\mathcal{X}^k, \alpha^k)$  which allow to identify the functions  $\mathbf{u}_i, p_i^*$  and  $p_i$  of the ansatz (7.5.2). Similarly to what occurs in the context of the elasticity system investigated in section 7.4, we prove that the averaged tensors  $\mathcal{X}^{k*}$  are symmetric matrix valued matrices for even values of  $k$ , and antisymmetric valued matrices otherwise. In section 7.5.2, we show that the formal averages  $\mathbf{u}_\varepsilon^*$  and  $p_\varepsilon^*$  are the solution of a formal, “infinite order” homogenized equation involving tensors  $M^k$ . Introducing new tensors  $N^k$  and  $\beta^k$ , we obtain “criminal” ansatz expressing  $\mathbf{u}_\varepsilon$  and  $p_\varepsilon$  in terms of the derivatives of  $\mathbf{u}_\varepsilon^*$  and  $p_\varepsilon^*$ . These tensors are used to build well-posed homogenized equations (eqn. (7.5.44)) of finite order, by using a minimization principle.

Section 7.5.3 is concerned with the error analysis of the homogenized approximations of  $(\mathbf{u}_\varepsilon, p_\varepsilon)$  generated by our procedure: our main result is stated in proposition 7.39 where we show that the solution  $(\mathbf{v}_K^*, q_K^*)$  of the homogenized equation of order  $2K + 2$  yield approximations of  $(\mathbf{u}_\varepsilon, p_\varepsilon)$  in the  $L^2(D_\varepsilon)$  norm of orders  $K + 3$  and  $K + 1$  for the velocity and the pressure respectively. In section 7.5.4, we investigate low volume fraction limits where the obstacle size  $\eta$  converges to zero: we prove the “coefficient-wise” convergence of the infinite order homogenized equation towards either of the three classical homogenized regimes depending on how  $\eta$  scales with respect to the critical size  $\varepsilon^{d/(d-2)}$ . We then establish in section 7.5.5 simplification properties for the tensors  $\mathcal{X}^{k*}$  and  $M^k$  in case the obstacle  $\eta T$  is symmetric with respect to the axes of the unit cell  $Y$ , which includes the vanishing of odd order tensors  $\mathcal{X}^{2k+1*}$  and  $M^{2k+1}$  in this case. Finally, section 7.5.6 briefly discusses the extension of our results to domains featuring multiple connected components in view of future applications to heat exchangers.

**Remark 7.26.** Following section 7.2.4, remember that the integers  $p, q \in \mathbb{N}$  may be used as a subscript or superscript indices when writing tensors in a limited number of places. Naturally, they should not to be confused with the pressure variables denoted by  $p_\varepsilon, p_i$  or  $q_K^*$ .

### 7.5.1 Formal infinite order two scale expansions: tensors $(\mathcal{X}^k, \alpha^k)$

As in the previous sections, we still assume (for the derivation of two-scale asymptotic expansions only) that the right-hand side  $\mathbf{f}$  can be formally decomposed into a power series in  $\varepsilon$ :

$$\forall x \in D, \mathbf{f}(x) = \sum_{i=0}^{+\infty} \varepsilon^i \mathbf{f}^i(x). \tag{7.5.4}$$

Inserting (7.5.2) into the Stokes system (7.5.1) yields the following cascade of equations

$$\left\{ \begin{array}{l} -\Delta_{yy} \mathbf{u}_{i+2} + \nabla_y p_{i+2} = \mathbf{f}_{i+2} - \nabla_x p_{i+2}^* - \nabla_x p_{i+1} + \Delta_{xy} \mathbf{u}_{i+1} + \Delta_{xx} \mathbf{u}_i, \quad \forall i \geq -2, \\ \operatorname{div}_y(\mathbf{u}_{i+2}) = -\operatorname{div}_x(\mathbf{u}_{i+1}), \quad \forall i \geq -2, \\ \mathbf{u}_{-2} = \mathbf{u}_{-1} = 0, \\ \int_D p_i^* dx = 0, \quad \forall i \geq 0, \\ p_{-1} = 0, \\ \mathbf{u}_i(x, \cdot) = 0 \text{ on } \partial(\eta T) \\ \mathbf{u}_i(x, \cdot) \text{ is } P\text{-periodic for any } x \in D, \forall i \geq 0, \\ \mathbf{u}_i(\cdot, y) \text{ is } D\text{-periodic for any } y \in P, \forall i \geq 0, \end{array} \right. \tag{7.5.5}$$

where the operators  $-\Delta_{yy}, -\Delta_{xy}, -\Delta_{yx}$  are defined by

$$-\Delta_{xx} = -\operatorname{div}_x(\nabla_x \cdot), \quad -\Delta_{xy} = -\operatorname{div}_x(\nabla_y \cdot) - \operatorname{div}_y(\nabla_x \cdot), \quad -\Delta_{yy} := -\operatorname{div}_y(\nabla_y \cdot).$$



We introduce a family of respectively vector valued tensors  $(\mathcal{X}_j^k(y))_{1 \leq j \leq d}$  and scalar valued tensors  $(\alpha_j^k(y))_{1 \leq j \leq d}$  defined by induction as the unique solutions in  $H^1(Y) \times L^2(Y)/\mathbb{R}$  to the following cell problems:

$$\begin{cases} -\Delta_{yy} \mathcal{X}_j^0 + \nabla_y \alpha_j^0 = \mathbf{e}_j & \text{in } Y, \\ \operatorname{div}_y(\mathcal{X}_j^0) = 0 & \text{in } Y \end{cases} \quad (7.5.6)$$

$$\begin{cases} -\Delta_{yy} \mathcal{X}_j^1 + \nabla_y \alpha_j^1 = (2\partial_l \mathcal{X}_j^0 - \alpha_j^0 \mathbf{e}_l) \otimes \mathbf{e}_l & \text{in } Y \\ \operatorname{div}_y(\mathcal{X}_j^1) = -(\mathcal{X}_j^0 - \langle \mathcal{X}_j^0 \rangle) \cdot \mathbf{e}_l \otimes \mathbf{e}_l & \text{in } Y, \end{cases} \quad (7.5.7)$$

$$\begin{cases} -\Delta_{yy} \mathcal{X}_j^{k+2} + \nabla_y \alpha_j^{k+2} = (2\partial_l \mathcal{X}_j^{k+1} - \alpha_j^{k+1} \mathbf{e}_l) \otimes \mathbf{e}_l + \mathcal{X}_j^k \otimes I & \text{in } Y \\ \operatorname{div}_y(\mathcal{X}_j^{k+2}) = -(\mathcal{X}_j^{k+1} - \langle \mathcal{X}_j^{k+1} \rangle) \cdot \mathbf{e}_l \otimes \mathbf{e}_l & \text{in } Y, \end{cases} \quad \forall k \geq 0 \quad (7.5.8)$$

where  $\langle \mathcal{X}_j^k \rangle$  denotes the cell average of the vector field  $\mathcal{X}_j^k$ :

$$\langle \mathcal{X}_j^k \rangle := \int_Y \mathcal{X}_j^k(y) dy.$$

Equations (7.5.6) to (7.5.8) are supplemented with the following boundary conditions:

$$\begin{cases} \int_Y \alpha_j^k dy = 0 \\ \mathcal{X}_j^k = 0 & \text{on } \partial(\eta T) \\ \mathcal{X}_j^k & \text{is } P\text{-periodic,} \end{cases} \quad \forall k \geq 0. \quad (7.5.9)$$

Similarly as before, we introduce the  $k$ -th order matrix valued tensors  $\mathcal{X}^k$  whose columns are the vector valued tensors  $(\mathcal{X}_j^k)$ :

$$(\mathcal{X}_{ij}^k(y))_{1 \leq i, j \leq d} := \begin{bmatrix} \mathcal{X}_1^k(y) & \dots & \mathcal{X}_d^k(y) \end{bmatrix}, \quad \forall y \in Y, \quad \forall k \geq 0.$$

We also denote by  $\alpha^k$  the  $k$ -th order vector valued tensor whose coordinates are the scalar tensors  $\alpha_j^k$ :

$$\alpha^k(y) := (\alpha_j^k(y))_{1 \leq j \leq d}, \quad \forall y \in Y, \quad \forall k \geq 0.$$

From the definition (7.5.9), the tensors  $\alpha_j^k$  are of zero average. Following the previous sections, we use a star notation to denote the averages of the tensor  $\mathcal{X}^k$  and of the functions  $\mathbf{u}_i$ :

$$\mathcal{X}^{k*} := \int_Y \mathcal{X}^k(y) dy, \quad \forall k \geq 0, \quad (7.5.10)$$

$$\mathbf{u}_i^*(x) := \int_Y \mathbf{u}_i(x, y) dy, \quad \forall x \in D, \quad \forall i \geq 0. \quad (7.5.11)$$

The tensors  $\mathcal{X}^k$  and  $\alpha^k$  allow to solve the cascade of equations (7.5.5):

**Proposition 7.29.** *Assume (H4). The solutions  $\mathbf{u}_i(x, y)$ ,  $p_i(x, y)$  of the cascade of equations (7.5.5) are given by*

$$\mathbf{u}_i(x, y) = \sum_{k=0}^i \mathcal{X}^k(y) \cdot \nabla^k (\mathbf{f}_{i-k}(x) - \nabla p_{i-k}^*(x)), \quad p_i(x, y) = \sum_{k=0}^i \alpha^k(y) \cdot \nabla^k (\mathbf{f}_{i-k}(x) - \nabla p_{i-k}^*(x)), \quad (7.5.12)$$

where the functions  $p_i^*$  are uniquely determined recursively as the solutions to the following elliptic system:

$$\begin{cases} -\operatorname{div}_x(\mathcal{X}^{0*} \nabla_x p_i^*) = \operatorname{div}_x(\mathcal{X}^{0*} \mathbf{f}_i) - \sum_{k=1}^i \operatorname{div}(\mathcal{X}^{k*} \cdot \nabla^k (\mathbf{f}_{i-k} - \nabla_x p_{i-k}^*)) & \text{in } D_\varepsilon, \\ \int_D p_i^* dx = 0 \end{cases} \quad \forall i \geq 0. \quad (7.5.13)$$

Since the average functions  $\mathbf{u}_i^*$  (eqn. (7.5.11)) are given by

$$\mathbf{u}_i^*(x) = \sum_{k=0}^i \mathcal{X}^{k*} \cdot \nabla^k (\mathbf{f}_{i-k}(x) - \nabla p_{i-k}^*(x)),$$

recognizing Cauchy products allows to rewrite identities (7.5.12) and (7.5.13) in terms of formal equality of power series:

$$\mathbf{u}_\varepsilon(x) = \sum_{i=0}^{+\infty} \varepsilon^{i+2} \mathcal{X}^i(x/\varepsilon) \cdot \nabla^i (\mathbf{f}(x) - \nabla p_\varepsilon^*(x)), \quad (7.5.14)$$

$$\mathbf{u}_\varepsilon^*(x) = \sum_{i=0}^{+\infty} \varepsilon^{i+2} \mathcal{X}^{i*} \cdot \nabla^i (\mathbf{f}(x) - \nabla p_\varepsilon^*(x)), \quad (7.5.15)$$

$$p_\varepsilon(x) = p_\varepsilon^*(x) + \sum_{i=0}^{+\infty} \varepsilon^{i+1} \alpha^i(x/\varepsilon) \cdot \nabla^i (\mathbf{f}(x) - \nabla p_\varepsilon^*(x)), \quad (7.5.16)$$

$$\operatorname{div}(\mathbf{u}_\varepsilon^*(x)) = 0. \quad (7.5.17)$$

*Proof.* The result is proved by induction. The case  $i = -1$  is straightforward thanks to the convention  $\mathbf{u}_{-1} = p_{-1} = 0$ . Assuming these results hold till rank  $i + 1$  with  $i \geq -2$ , we compute, replacing (7.5.12) into (7.5.5):

$$\left\{ \begin{aligned} & (-\Delta_{yy} \mathbf{u}_{i+2} + \nabla_y p_{i+2})(x, y) = (f_{i+2,j}(x) - \partial_j p_{i+2}^*(x)) \mathbf{e}_j \\ & + \sum_{k=0}^{i+1} (-\alpha_j^k(y) \mathbf{e}_l \otimes \mathbf{e}_l + 2\partial_l \mathcal{X}_j^k(y) \otimes \mathbf{e}_l) \cdot \nabla^{k+1} (f_{i+1-k,j}(x) - \partial_j p_{i+1-k}^*(x)) \\ & + \sum_{k=0}^i (\mathcal{X}_j^k(y) \otimes I) \cdot \nabla^{k+2} (f_{i-k,j}(x) - \partial_j p_{i-k}^*(x)) \\ & = (f_{i+2,j}(x) - \partial_j p_{i+2}^*(x)) \mathbf{e}_j + (2\partial_l \mathcal{X}_j^0(y) \otimes \mathbf{e}_l - \alpha_j^0(y) \mathbf{e}_l \otimes \mathbf{e}_l) \cdot \nabla (f_{i+1,j}(x) - \partial_j p_{i+1}^*(x)) \\ & + \sum_{k=0}^i (2\partial_l \mathcal{X}_j^{k+1}(y) \otimes \mathbf{e}_l - \alpha_j^{k+1}(y) \mathbf{e}_l \otimes \mathbf{e}_l + \mathcal{X}_j^k(y) \otimes I) \cdot \nabla^{k+2} (f_{i-k,j}(x) - \partial_j p_{i-k}^*(x)) \\ & \operatorname{div}_y(\mathbf{u}_{i+2})(x, y) = - \sum_{k=0}^{i+1} (\mathcal{X}_j^k(y) \cdot \mathbf{e}_l \otimes \mathbf{e}_l) \cdot \nabla^{k+1} (f_{i+1-k,j}(x) - \partial_j p_{i+1-k}^*(x)). \end{aligned} \right. \quad (7.5.18)$$

The system (7.5.18) admits a unique solution  $(\mathbf{u}_{i+2}, p_{i+2})$  with  $\int_Y p_{i+2}(x, y) dy = 0$  if and only if the following compatibility condition (the so-called ‘‘Fredholm alternative’’) holds:

$$\forall i \geq -1, \int_Y \operatorname{div}_y(\mathbf{u}_{i+2})(x, y) dy = - \sum_{k=0}^{i+1} [\langle \mathcal{X}_j^k \rangle \cdot \mathbf{e}_l \otimes \mathbf{e}_l] \cdot \nabla^{k+1} (f_{i+1-k,j}(x) - \partial_j p_{i+1-k}^*(x)) = 0.$$

The above equation can be rewritten in terms of an equation determining  $p_{i+1}^*$  given the values of  $p_k^*$  for  $0 \leq k \leq i$ :

$$\langle \mathcal{X}_j^0 \rangle \cdot \mathbf{e}_l \partial_l (f_{i+1,j} - \partial_j p_{i+1}^*) = - \sum_{k=1}^{i+1} [\langle \mathcal{X}_j^k \rangle \cdot \mathbf{e}_l \otimes \mathbf{e}_l] \cdot \nabla^{k+1} (f_{i+1-k,j} - \partial_j p_{i+1-k}^*),$$

which is a rewriting of (7.5.13) (at order  $i + 1$  with any  $i \geq -1$ ). In particular, this identity allows to rewrite the second equation of (7.5.18) as

$$\operatorname{div}_y(\mathbf{u}_{i+2})(x, y) = - \sum_{k=0}^{i+1} [(\mathcal{X}_j^k(y) - \langle \mathcal{X}_j^k \rangle) \cdot \mathbf{e}_l \otimes \mathbf{e}_l] \cdot \nabla^{k+1} (f_{i+1-k,j}(x) - \partial_j p_{i+1-k}^*(x)). \quad (7.5.19)$$

By linearity, (7.5.18) and (7.5.19) and the cell problems (7.5.6) to (7.5.8) imply

$$\begin{aligned} \mathbf{u}_{i+2}(x, y) &= (f_{i+2,j}(x) - \partial_j p_{i+2}^*(x)) \mathbf{x}_j^0(y) + \mathbf{x}_j^1(y) \cdot \nabla (f_{i+1,j}(x) - \partial_j p_{i+1}^*(x)) \\ &\quad + \sum_{k=0}^i \mathbf{x}_j^{k+2}(y) \cdot \nabla^{k+2} (f_{i-k,j}(x) - \partial_j p_{i-k}^*(x)) \\ p_{i+2}(x, y) &= (f_{i+2,j}(x) - \partial_j p_{i+2}^*(x)) \alpha_j^0 + \alpha_j^1(y) \cdot \nabla (f_{i+1,j}(x) - \partial_j p_{i+1}^*(x)) \\ &\quad + \sum_{k=0}^i \alpha_j^{k+2}(y) \cdot \nabla^{k+2} (f_{i-k,j}(x) - \partial_j p_{i-k}^*(x)), \end{aligned}$$

which is exactly the result (7.5.12) at rank  $i + 2$ .  $\square$

**Remark 7.27.** The truncation of the series (7.5.15) at first order yields the well-known Darcy's law. The next terms of the series have been obtained in [228, 74], at least up to the order  $i = 1$ .

**Remark 7.28.** The ansatz (7.5.14) involves a bit of asymptotic *crime* because it features  $p_\varepsilon^*$  which is itself a formal power series in  $\varepsilon$  (recall its definition (7.5.3)). Although this ansatz is criminal, it remains close to the classical ansatz (7.5.2). Therefore, in the subsequent sections, the denomination ‘‘criminal ansatz for the Stokes system (7.5.1)’’ is used exclusively for referring to the more *severe* asymptotic crime (7.5.32) introduced hereafter.

The identities of propositions 7.2 and 7.17 for the average of the tensors  $\mathcal{X}^{k*}$  extend as follows in the context of the fluid system (7.5.1):

**Proposition 7.30.** *For any  $k \geq 0$  and  $0 \leq p \leq k$ , the following identity holds for the matrix valued tensor  $\mathcal{X}^{k*}$ :*

$$\forall 1 \leq i, j \leq d, \mathcal{X}_{ij}^{k*} = (-1)^p \int_Y ((-\Delta_{yy} \mathbf{x}_i^p + \nabla \alpha_i^p) \cdot \mathbf{x}_j^{k-p} + \nabla \alpha_j^{k-p} \cdot \mathbf{x}_i^p - \mathbf{x}_j^{k-p-1} \cdot \mathbf{x}_i^{p-1} \otimes I) dy \quad (7.5.20)$$

with  $\mathbf{x}_i^{-1} = 0$  by convention. In particular, for any  $k \geq 0$ :

- $\mathcal{X}^{2k*}$  takes values in the set of  $d \times d$  symmetric matrices:

$$\mathcal{X}_{ij}^{2k*} = (-1)^k \int_Y (\nabla \mathbf{x}_i^k : \nabla \mathbf{x}_j^k + \nabla \alpha_i^k \cdot \mathbf{x}_j^k + \nabla \alpha_j^k \cdot \mathbf{x}_i^k - \mathbf{x}_i^{k-1} \cdot \mathbf{x}_j^{k-1} \otimes I) dy \quad (7.5.21)$$

- $\mathcal{X}^{2k+1*}$  takes values in the set of  $d \times d$  antisymmetric matrices:

$$\begin{aligned} \mathcal{X}_{ij}^{2k+1*} &= (-1)^k \int_Y (\mathbf{x}_i^k \cdot \nabla \mathbf{x}_j^k - \mathbf{x}_j^k \cdot \nabla \mathbf{x}_i^k + \alpha_i^k \mathbf{x}_j^k - \alpha_j^k \mathbf{x}_i^k) \cdot \mathbf{e}_l \otimes \mathbf{e}_l dy \\ &\quad + (-1)^k \int_Y (\mathbf{x}_j^{k-1} \cdot \mathbf{x}_j^k - \mathbf{x}_i^{k-1} \cdot \mathbf{x}_j^k) dy. \end{aligned} \quad (7.5.22)$$

*Proof.* The result holds for  $p = 0$  because

$$\mathcal{X}_{ij}^{k*} = \int_Y \mathbf{x}_j^k \cdot \mathbf{e}_i dy = \int_Y \mathbf{x}_j^k \cdot (-\Delta_{yy} \mathbf{x}_i^0 + \nabla \alpha_i^0) dy.$$

Assume now that (7.5.20) holds till rank  $p$  with  $k > p \geq 0$ , and let us show that it implies the result at rank  $p + 1$ . We write, after an integration by parts and by using (7.5.6) to (7.5.8):

$$\begin{aligned} \mathcal{X}_{ij}^{k*} &= (-1)^p \int_Y (-\mathbf{x}_i^p \cdot \Delta \mathbf{x}_j^{k-p} - \alpha_i^p \operatorname{div}(\mathbf{x}_j^{k-p}) - \alpha_j^{k-p} \operatorname{div}(\mathbf{x}_i^p) - \mathbf{x}_j^{k-p-1} \cdot \mathbf{x}_i^{p-1} \otimes I) dy \\ &= (-1)^p \int_Y [(\mathbf{x}_i^p \cdot (2\partial_l \mathbf{x}_j^{k-p-1} - \alpha_j^{k-p-1} \mathbf{e}_l) \otimes \mathbf{e}_l + \mathbf{x}_j^{k-p-2} \otimes I - \nabla \alpha_j^{k-p}) \cdot \mathbf{x}_i^p \\ &\quad + \alpha_i^p \mathbf{x}_j^{k-p-1} \cdot \mathbf{e}_l \otimes \mathbf{e}_l + \alpha_j^{k-p} \mathbf{x}_i^{p-1} \cdot \mathbf{e}_l \otimes \mathbf{e}_l - \mathbf{x}_j^{k-p-1} \cdot \mathbf{x}_i^{p-1} \otimes I] dy \\ &= (-1)^p \int_Y [-\mathbf{x}_j^{k-p-1} \cdot ((2\partial_l \mathbf{x}_i^p - \alpha_i^p \mathbf{e}_l) \otimes \mathbf{e}_l + \mathbf{x}_i^{p-1} \otimes I) + \alpha_j^{k-p-1} \operatorname{div}(\mathbf{x}_i^{p+1}) \\ &\quad - \nabla \alpha_j^{k-p} \cdot \mathbf{x}_i^p - \alpha_j^{k-p} \operatorname{div}(\mathbf{x}_i^p) + \mathbf{x}_j^{k-p-2} \cdot \mathbf{x}_i^p \otimes I] dy \\ &= (-1)^p \int_Y [-\mathbf{x}_j^{k-p-1} \cdot (-\Delta_{yy} \mathbf{x}_i^{p+1} + \nabla \alpha_i^{p+1}) - \nabla \alpha_j^{k-p-1} \cdot \mathbf{x}_i^{p+1} + \mathbf{x}_j^{k-p-2} \cdot \mathbf{x}_i^p \otimes I] dy, \end{aligned}$$

whence (7.5.20) at rank  $p + 1$ .

The expression (7.5.21) for  $\mathcal{X}_{ij}^{2k*}$  is obtained by setting  $k \leftarrow 2k$  and  $p \leftarrow k$  in (7.5.20). The expression for  $\mathcal{X}_{ij}^{2k+1*}$  is obtained by setting  $k \leftarrow 2k + 1$  and  $p \leftarrow k$  and performing the following integration by part:

$$\begin{aligned}\mathcal{X}_{ij}^{2k+1*} &= (-1)^k \int_Y ((-\Delta_{yy} \mathbf{x}_j^{k+1} + \nabla \alpha_j^{k+1}) \cdot \mathbf{x}_i^k + \nabla \alpha_i^k \cdot \mathbf{x}_j^{k+1} - \mathbf{x}_j^k \cdot \mathbf{x}_i^{k-1} \otimes I) dy \\ &= (-1)^k \int_Y [(2\partial_l \mathbf{x}_j^k - \alpha_j^k e_l) \otimes e_l + \mathbf{x}_j^{k-1} \otimes I] \cdot \mathbf{x}_i^k + \alpha_i^k \mathbf{x}_i^k \cdot e_l \otimes e_l - \mathbf{x}_j^k \cdot \mathbf{x}_i^{k-1} \otimes I] dy,\end{aligned}$$

from where one recognizes (7.5.22).  $\square$

We end this part with a result analogous to that of proposition 7.18 which yields as a corollary the linear independence of the tensors  $(\mathcal{X}_j^k)_{1 \leq j \leq d}$  and the positive definiteness of the permeability tensor  $\mathcal{X}^{0*}$ .

**Proposition 7.31.** *The following identities hold for any  $1 \leq j \leq d$ :*

$$\forall k \geq 0, -\Delta_{yy}(\partial_{i_1 \dots i_k}^k \mathbf{x}_{i_1 \dots i_k, j}^k) + \nabla \left( \sum_{p=0}^k (-1)^p \partial_{i_1 \dots i_{k-p}}^{k-p} \alpha_{i_1 \dots i_{k-p}, j}^{k-p} \right) = (-1)^k (k+1) \mathbf{e}_j, \quad (7.5.23)$$

$$\operatorname{div}(\partial_{i_1 \dots i_k}^k \mathbf{x}_{i_1 \dots i_k, j}^k) = 0. \quad (7.5.24)$$

*Proof.* We prove the following identity by induction on  $k$ ,

$$\begin{aligned}\forall k \geq 0, -\Delta_{yy}(\partial_{i_1 \dots i_k}^k \mathbf{x}_{i_1 \dots i_k, j}^k) + \nabla(\partial_{i_1 \dots i_k}^k \alpha_{i_1 \dots i_k, j}^k) \\ = (-1)^k (k+1) \mathbf{e}_j - (-1)^k \sum_{p=0}^{k-1} (-1)^p \nabla(\partial_{i_1 \dots i_p}^p \alpha_{i_1 \dots i_p, j}^p),\end{aligned} \quad (7.5.25)$$

which is equivalent to (7.5.23). The identity clearly holds at rank  $k = 0$  and  $k = -1$  by assuming the convention  $\mathbf{x}_j^{-1} = 0$ . Assume it to be true up to rank  $k - 1$ , we write

$$\begin{aligned}-\Delta_{yy}(\partial_{i_1 \dots i_k}^k \mathbf{x}_{i_1 \dots i_k, j}^k) + \nabla(\partial_{i_1 \dots i_k}^k \alpha_{i_1 \dots i_k, j}^k) &= \partial_{i_1 \dots i_k}^k ((2\partial_{i_k} \mathbf{x}_{i_1 \dots i_{k-1}, j}^{k-1} - \alpha_{i_1 \dots i_{k-1}, j}^{k-1} \mathbf{e}_{i_k}) + \delta_{i_k i_{k-1}} \mathbf{x}_{i_1 \dots i_{k-2}, j}^{k-2}) \\ &= 2\Delta_{yy}(\partial_{i_1 \dots i_{k-1}}^{k-1} \mathbf{x}_{i_1 \dots i_{k-1}, j}^{k-1}) + \Delta_{yy}(\partial_{i_1 \dots i_{k-2}}^{k-2} \mathbf{x}_{i_1 \dots i_{k-2}, j}^{k-2}) - \nabla(\partial_{i_1 \dots i_{k-1}}^{k-1} \alpha_{i_1 \dots i_{k-1}, j}^{k-1}) \\ &= 2 \left( -(-1)^{k-1} k \mathbf{e}_j + (-1)^{k-1} \sum_{p=0}^{k-2} (-1)^p \nabla(\partial_{i_1 \dots i_p}^p \alpha_{i_1 \dots i_p, j}^p) + \nabla(\partial_{i_1 \dots i_{k-1}}^{k-1} \alpha_{i_1 \dots i_{k-1}, j}^{k-1}) \right) \\ &\quad + \left( -(-1)^{k-2} (k-1) \mathbf{e}_j + (-1)^{k-2} \sum_{p=0}^{k-3} (-1)^p \nabla(\partial_{i_1 \dots i_p}^p \alpha_{i_1 \dots i_p, j}^p) + \nabla(\partial_{i_1 \dots i_{k-2}}^{k-2} \alpha_{i_1 \dots i_{k-2}, j}^{k-2}) \right) \\ &\quad - \nabla(\partial_{i_1 \dots i_{k-1}}^{k-1} \alpha_{i_1 \dots i_{k-1}, j}^{k-1}) \\ &= (-1)^k (k+1) \mathbf{e}_j - (-1)^k \sum_{p=0}^{k-3} (-1)^p \nabla(\partial_{i_1 \dots i_p}^p \alpha_{i_1 \dots i_p, j}^p) + \underbrace{(2(-1)^{2k-3} + 1)}_{-1 = -(-1)^k (-1)^{k-2}} \nabla(\partial_{i_1 \dots i_{k-2}}^{k-2} \alpha_{i_1 \dots i_{k-2}, j}^{k-2}) \\ &\quad - \underbrace{(-1)}_{-(-1) = -(-1)^k (-1)^{k-1}} \nabla(\partial_{i_1 \dots i_{k-1}}^{k-1} \alpha_{i_1 \dots i_{k-1}, j}^{k-1}),\end{aligned}$$

from where (7.5.25) follows at rank  $k$ .

The second equality is obtained similarly by writing

$$\begin{aligned}\operatorname{div}(\partial_{i_1 \dots i_k}^k \mathbf{x}_{i_1 \dots i_k, j}^k) &= -\partial_{i_1 \dots i_k}^k (\mathbf{x}_{i_1 \dots i_{k-1}, j}^{k-1} - \langle \mathbf{x}_{i_1 \dots i_{k-1}, j}^{k-1} \rangle) \cdot \mathbf{e}_{i_k} \\ &= -\operatorname{div}(\partial_{i_1 \dots i_{k-1}}^{k-1} \mathbf{x}_{i_1 \dots i_{k-1}, j}^{k-1}).\end{aligned}$$

$\square$

**Corollary 7.8.** *Assume (H4). The family of  $k$ -th order vector valued tensors  $(\mathcal{X}_j^k)_{1 \leq j \leq d}$  is linearly independent. In particular, the matrix*

$$\mathcal{X}^{0*} = (\mathcal{X}_{ij}^{0*})_{1 \leq i, j \leq d} \text{ with } \mathcal{X}_{ij}^{0*} = \int_Y \nabla \mathbf{x}_i^0 : \nabla \mathbf{x}_j^0 dy, \quad \forall 1 \leq i, j \leq d,$$

*is symmetric positive definite.*

*Proof.* (see also [272]) Let  $(\lambda_j)_{1 \leq j \leq d}$  be some coefficients such that  $\sum_{i=1}^d \lambda_j \mathcal{X}_j^k = 0$ . Then [proposition 7.31](#) implies the existence of a periodic function  $\phi$  satisfying

$$\nabla \phi = (-1)^k (k+1) \sum_{j=1}^d \lambda_j \mathbf{e}_j. \quad (7.5.26)$$

From the fact that  $Y$  is connected, there exist constant coefficients  $(C_i)_{1 \leq i \leq d}$  such that

$$\forall x \in Y, \phi(x) = (-1)^k (k+1) \sum_{j=1}^d (\lambda_j x_j + C_j).$$

Remembering that  $\phi$  must be a  $P$ -periodic function, this is possible only if  $\lambda_j = 0$  for all  $1 \leq j \leq d$ .  $\square$

### 7.5.2 Higher order homogenized equations: tensors $M^k, N^k, \beta^k$ and $\mathbb{D}_K$

In this section, we derive an “infinite order” homogenized equation (eqn. (7.5.29) below) for the formal average  $\mathbf{u}_\varepsilon^*$ , before outlining how to truncate so as to obtain well-posed homogenized equations of finite order (eqn. (7.5.39) below). Since the ansatz (7.5.14) has the same structure than (7.4.8) in the elasticity case considered in [section 7.4](#), the construction of the tensors  $M^k$  and  $N_j^k$  follows similarly. The following definition makes sense because we recalled in [corollary 7.8](#) that the Darcy permeability tensor  $\mathcal{X}^{0*}$  is invertible.

**Proposition 7.32.** *Let  $M^k$  be the tensor of order  $k$  defined by induction as follows:*

$$\begin{cases} M^0 = (\mathcal{X}^{0*})^{-1} \\ M^k = -(\mathcal{X}^{0*})^{-1} \sum_{p=0}^{k-1} \mathcal{X}^{k-p*} \otimes M^p, \quad \forall k \geq 1. \end{cases} \quad (7.5.27)$$

The source terms  $\mathbf{f}_i$  (eqn. (7.5.4)) are given in terms of the averaged ansatz terms  $\mathbf{u}_i^*(x)$  and  $p_i^*(x)$  through the following identity:

$$\forall i \geq 0, \mathbf{f}_i(x) - \nabla p_i^*(x) = \sum_{k=0}^i M^k \cdot \nabla^k \mathbf{u}_{i-k}^*(x). \quad (7.5.28)$$

The above identity together with (7.5.13) can be rewritten formally as the following “infinite order” homogenized equation for the homogenized averages  $\mathbf{u}_\varepsilon^*$  and  $p_\varepsilon^*(x)$  of (7.5.3):

$$\begin{cases} \sum_{i=0}^{+\infty} \varepsilon^{i-2} M^i \cdot \nabla^i \mathbf{u}_\varepsilon^* + \nabla p_\varepsilon^* = \mathbf{f} & \text{in } D \\ \operatorname{div}(\mathbf{u}_\varepsilon^*) = 0 & \text{in } D \\ \int_D p_\varepsilon^* dx = 0 \\ \mathbf{u}_\varepsilon^* \text{ is } D\text{-periodic.} \end{cases} \quad (7.5.29)$$

Since the family of tensors  $(M^k)_{k \in \mathbb{N}}$  is defined by the same recurrence than (7.4.15) in the elasticity case, the explicit formula for  $M^k$  given in [proposition 7.20](#) holds without modification:

$$M^k = \sum_{p=1}^k (-1)^p \sum_{\substack{i_1 + \dots + i_p = k \\ 1 \leq i_1, \dots, i_p \leq k}} (\mathcal{X}^{0*})^{-1} \otimes \mathcal{X}^{i_1*} \otimes (\mathcal{X}^{0*})^{-1} \otimes \dots \otimes (\mathcal{X}^{0*})^{-1} \otimes \mathcal{X}^{i_p*} \otimes (\mathcal{X}^{0*})^{-1}, \quad \forall k \geq 1. \quad (7.5.30)$$

We now introduce matrix valued tensors  $N^k$  and vector valued tensors  $\beta^k$  which allow to obtain “criminal ansatz” expressing the velocity and pressure  $(\mathbf{u}_\varepsilon, p_\varepsilon)$  in terms of their formal average  $(\mathbf{u}_\varepsilon^*, p_\varepsilon^*)$ .

**Proposition 7.33.** *Let  $N^k$  and  $\beta^k$  be respectively the  $k$ -th order matrix valued and vector valued tensors defined by*

$$N^k(y) := \sum_{p=0}^k \mathcal{X}^{k-p}(y) \otimes M^p, \quad \beta^k(y) := \sum_{p=0}^k (-1)^p M^p \cdot \alpha^{k-p}(y), \quad \forall y \in Y.$$

Then the terms  $\mathbf{u}_i(x, y)$  and  $p_i(x, y)$  of the oscillating ansatz (7.5.14) and (7.5.16) can be rewritten in terms of the averages  $\mathbf{u}_i^*$  (eqn. (7.5.11)) and  $p_i^*$  as follows:

$$\forall i \geq 0, \mathbf{u}_i(x, y) = \sum_{k=0}^i N^k(y) \cdot \nabla^k \mathbf{u}_{i-k}^*(x), \quad p_i(x, y) = \sum_{k=0}^i \beta^k(y) \cdot \nabla^k \mathbf{u}_{i-k}^*(x). \quad (7.5.31)$$

These equations can be rewritten formally in terms of the following “criminal ansatz” for  $(\mathbf{u}_\varepsilon, p_\varepsilon)$ :

$$\forall x \in D_\varepsilon, \mathbf{u}_\varepsilon(x) = \sum_{i=0}^{+\infty} \varepsilon^i N^i(x/\varepsilon) \cdot \nabla^i \mathbf{u}_\varepsilon^*(x), \quad p_\varepsilon(x) = p_\varepsilon^*(x) + \sum_{i=0}^{+\infty} \varepsilon^{i-1} \beta^i(x/\varepsilon) \cdot \nabla^i \mathbf{u}_\varepsilon^*(x). \quad (7.5.32)$$

*Proof.* The proof is identical to that of proposition 7.8: the key ingredient is to substitute (7.5.28) into (7.5.12). Performing the same change of indices as in (7.3.38) yields

$$\begin{aligned} \mathbf{u}_i(x, y) &= \sum_{k=0}^i \sum_{p=0}^k (\mathcal{X}^p(y) \otimes M^{p-k}) \cdot \nabla^k \mathbf{u}_{i-k}^*(x), \\ p_i(x, y) &= \sum_{k=0}^i \sum_{p=0}^k ((M^{p-k})^T \cdot \boldsymbol{\alpha}^p(y)) \cdot \nabla^k \mathbf{u}_{i-k}^*(x). \end{aligned}$$

The result follows, where we use  $(M^{p-k})^T = (-1)^{p-k} M^{p-k}$  as stated in corollary 7.9 below.  $\square$

In what follows, we denote by  $(\mathbf{N}_j^k)_{1 \leq j \leq d}$  and by  $(\beta_j^k)_{1 \leq j \leq d}$  respectively the column vectors and the coefficients of  $N^k(y)$  and  $\beta^k(y)$ :

$$\forall 1 \leq i, j \leq d, \mathbf{N}_j^k := N^k \mathbf{e}_j \text{ and } \beta_j^k := \beta^k \cdot \mathbf{e}_j.$$

In addition, the convention  $\mathbf{N}_j^{-1} = 0$  is assumed. The proposition 7.22 and its corollary 7.5 extend as follows in the context of the Stokes system:

**Proposition 7.34.** *The  $k$ -th order tensors  $N^k$ ,  $(\mathbf{N}_j^k)_{1 \leq j \leq d}$ ,  $\beta^k$  and  $(\beta_j^k)_{1 \leq j \leq d}$  satisfy:*

- (i)  $\int_Y N^0(y) dy = I$  and  $\int_Y N^k(y) dy = 0$  for any  $k \geq 1$ ;
- (ii)  $\int_Y \beta^k(y) dy = 0$  for any  $k \geq 0$ ;
- (iii)  $\forall k \geq -2, \forall 1 \leq j \leq d$ ,

$$\begin{cases} -\Delta_{yy} \mathbf{N}_j^{k+2} + \nabla \beta_j^{k+2} = (2\partial_l \mathbf{N}_j^{k+1} - \beta_j^{k+1} \mathbf{e}_l) \otimes \mathbf{e}_l + \mathbf{N}_j^k \otimes I + M^{k+2} \mathbf{e}_j, \\ \operatorname{div}(\mathbf{N}_j^{k+2}) = -(\mathbf{N}_j^{k+1} - \langle \mathbf{N}_j^{k+1} \rangle) \cdot \mathbf{e}_l \otimes \mathbf{e}_l; \end{cases} \quad (7.5.33)$$

- (iv) for any  $k \geq 0$ ,

$$-\Delta_{yy} (\partial_{i_1 \dots i_k}^k \mathbf{N}_{i_1 \dots i_k, j}^k) + \nabla \left( \sum_{p=0}^k (-1)^p \partial_{i_1 \dots i_{k-p}}^{k-p} \beta_{i_1 \dots i_k, j}^{k-p} \right) = (-1)^{k+1} (k+1) M^0 \mathbf{e}_j; \quad (7.5.34)$$

- (v) for any  $1 \leq p \leq k-1$ ,

$$M_{ij}^k = (-1)^{p+1} \int_Y ((-\Delta_{yy} \mathbf{N}_i^p + \nabla \beta_i^p) \cdot \mathbf{N}_j^{k-p} + \nabla \beta_j^{k-p} \cdot \mathbf{N}_i^p - \mathbf{N}_i^{p-1} \cdot \mathbf{N}_j^{k-p-1} \otimes I) dy. \quad (7.5.35)$$

*Proof.* (i) and (ii) are straightforward consequences of (7.5.27).

- (iii) is obtained by writing, for  $k \geq 0$  (implicit summation on the repeated index  $j$  assumed):

$$\begin{aligned} -\Delta_{yy} \mathbf{N}_j^{k+2} + \nabla \beta_j^{k+2} &= -\Delta_{yy} \left( \sum_{p=0}^{k+2} \boldsymbol{\chi}_i^{k+2-p}(y) \otimes M_{ij}^p \right) + \nabla \left( \sum_{p=0}^{k+2} \alpha_i^{k+2-p}(y) \otimes M_{ij}^p \right) \\ &= \sum_{p=0}^k \left[ (2\partial_l \boldsymbol{\chi}_i^{k+1-p} - \alpha_i^{k+1-p} \mathbf{e}_l) \otimes \mathbf{e}_l + \boldsymbol{\chi}_i^{k-p} \otimes I \right] M_{ij}^p + (2\partial_l \boldsymbol{\chi}_i^0 - \alpha_i^0 \mathbf{e}_l) M_{ij}^{k+1} + M_{ij}^{k+2} \mathbf{e}_i \\ &= (2\partial_l \mathbf{N}_j^{k+1} - \beta_j^{k+1} \mathbf{e}_l) \otimes \mathbf{e}_l + \mathbf{N}_j^k \otimes I + M^{k+2} \mathbf{e}_j. \end{aligned}$$

$$\operatorname{div}(N_j^{k+2}) = \sum_{p=0}^{k+2} \operatorname{div}(\mathcal{X}_i^{k+2-p}) M_{ij}^p = - \sum_{p=0}^{k+1} M_{ij}^p (\mathcal{X}_i^{k+1-p} - \langle \mathcal{X}_i^{k+1-p} \rangle) \cdot e_l \otimes e_l.$$

The proof is identical for  $k = -1$  and  $k = -2$ .

(iv) is obtained by following the proof of [proposition 7.31](#) and using [\(7.5.33\)](#).

(v) is obtained by induction following the proof of [proposition 7.9](#) and by using [\(7.5.33\)](#). The case  $p = 1$  with  $k \geq 2$  is treated by writing

$$\begin{aligned} M_{ij}^k &= \int_Y N_i^0 \cdot M^k e_j dy \\ &= \int_Y N_i^0 \cdot (-\Delta N_j^k + \nabla \beta_j^k - (2\partial_l N_j^{k-1} - \beta_j^{k-1} e_l) \otimes e_l - N_j^{k-2} \otimes I) dy \\ &= \int_Y (-\Delta N_i^0 \cdot N_j^k + 2\partial_l N_i^0 \cdot N_j^{k-1} \otimes e_l + \nabla \beta_j^{k-1} \cdot N_i^1 - N_i^0 \cdot N_j^{k-2} \otimes I) dy \\ &= \int_Y (N_j^k \cdot M^0 e_i - \nabla \beta_i^0 \cdot N_j^k + (-\Delta N_i^1 + \nabla \beta_i^1 + \beta_i^0 e_l \otimes e_l - M^1 e_i) \cdot N_j^{k-1}) dy \\ &\quad + \int_Y (\nabla \beta_j^{k-1} \cdot N_i^1 - N_i^0 \cdot N_j^{k-2} \otimes I) dy \\ &= \int_Y ((-\Delta N_i^1 + \nabla \beta_i^1) \cdot N_j^{k-1} + \nabla \beta_j^{k-1} \cdot N_i^1 - N_i^0 \cdot N_j^{k-2} \otimes I) dy. \end{aligned}$$

We now assume the result holds till rank  $1 \leq p < k - 1$  and we prove it at rank  $p + 1$ :

$$\begin{aligned} M_{ij}^k &= (-1)^{p+1} \int_Y (-\Delta N_j^{k-p} \cdot N_i^p + \beta_i^p N_j^{k-p-1} \cdot e_l \otimes e_l + \nabla \beta_j^{k-p} \cdot N_i^p - N_i^{p-1} \cdot N_j^{k-p-1} \otimes I) dy \\ &= (-1)^{p+1} \int_Y (-\nabla \beta_j^{k-p} + (2\partial_l N_j^{k-p-1} - \beta_j^{k-p-1} e_l) \otimes e_l + N_j^{k-p-2} \otimes I + M^{k-p} e_j) \cdot N_i^p dy \\ &\quad + (-1)^{p+1} \int_Y (\beta_i^p N_j^{k-p-1} \cdot e_l \otimes e_l + \nabla \beta_j^{k-p} \cdot N_i^p - N_i^{p-1} \cdot N_j^{k-p-1} \otimes I) dy \\ &= (-1)^{p+1} \int_Y (-N_j^{k-p-1} \cdot ((2\partial_l N_i^p - \beta_i^p e_l) \otimes e_l + N_i^{p-1} \otimes I + M^{p+1} e_i)) dy \\ &\quad + (-1)^{p+1} \int_Y (-\nabla \beta_j^{k-p-1} \cdot N_i^{p+1} + N_j^{k-p-2} \cdot N_i^p \otimes I) dy \\ &= (-1)^{p+1} \int_Y (-N^{k-p-1} \cdot (-\Delta N_i^{p+1} + \nabla \beta_i^{p+1}) - \nabla \beta_j^{k-p-1} \cdot N_i^{p+1} + N_j^{k-p-2} \cdot N_i^p \otimes I) dy. \end{aligned}$$

□

**Corollary 7.9.** For any  $k \geq 0$ ,

- $M^{2k}$  is a symmetric matrix valued tensor, and the following identities hold:

$$M_{ij}^0 = \int_Y \nabla N_i^0 : \nabla N_j^0 dy,$$

$$\forall k \geq 1, M_{ij}^{2k} = (-1)^{k+1} \int_Y (\nabla N_i^k : \nabla N_j^k + \nabla \beta_i^k \cdot N_j^k + \nabla \beta_j^k \cdot N_i^k - N_i^{k-1} \cdot N_j^{k-1} \otimes I) dy.$$

- $M^{2k+1}$  is an antisymmetric matrix valued tensor, and the following identities hold:

$$\begin{aligned} \forall k \geq 0, M_{ij}^{2k+1} &= (-1)^{k+1} \int_Y (N_i^k \cdot \nabla N_j^k - N_j^k \cdot \nabla N_i^k + \beta_i^k N_j^k - \beta_j^k N_i^k) \cdot e_l \otimes e_l dy \\ &\quad + (-1)^{k+1} \int_Y (N_j^{k-1} \cdot N_i^k - N_i^{k-1} \cdot N_j^k) \otimes I dy. \end{aligned}$$

We now use the criminal ansatz to derive well-posed homogenized equations of (finite) order  $2K + 2$ . The formal identities (7.5.32) lead us to introduce the following truncated ansatz  $\mathbf{w}_{\varepsilon,K}(\mathbf{v})$  and  $q_{\varepsilon,K}(\mathbf{v}, \phi)$  for the velocity and pressure:

$$\forall \mathbf{v} \in H^{K+1}(D, \mathbb{R}^d), \mathbf{w}_{\varepsilon,K}(\mathbf{v})(x) := \sum_{k=0}^K \varepsilon^k N^k(x/\varepsilon) \cdot \nabla^k \mathbf{v}(x), \quad x \in D_\varepsilon, \quad (7.5.36)$$

$$\forall q \in L^2(D), \forall \mathbf{v} \in H^{K+1}(D), q_{\varepsilon,K}(\mathbf{v}, \phi)(x) := \phi + \sum_{i=0}^K \varepsilon^{i-1} \beta^i(x/\varepsilon) \cdot \nabla^i \mathbf{v}(x), \quad (7.5.37)$$

where  $\mathbf{v}$  and  $\phi$  are functions that are sought to approximate the homogenized averages  $\mathbf{u}_\varepsilon^*$  and  $p_\varepsilon^*$  respectively. Following the methodology of sections 7.3.3 and 7.4.2, we build well-posed homogenized equations of higher—but finite—order from a minimization principle.

Recall that the solution  $\mathbf{u}_\varepsilon$  to the Stokes problem (7.5.1) is the unique minimizer of the constrained problem

$$\begin{aligned} \mathbf{u}_\varepsilon = \arg \min_{\mathbf{w} \in H^1(D_\varepsilon, \mathbb{R}^d)} J(\mathbf{w}, \mathbf{f}) &:= \int_D \left( \frac{1}{2} \nabla \mathbf{w} : \nabla \mathbf{w} - \mathbf{f} \cdot \mathbf{w} \right) dy \\ \text{s.t.} \begin{cases} \operatorname{div}(\mathbf{w}) = 0 \text{ in } D_\varepsilon, \\ \mathbf{w} = 0, \text{ on } \partial\omega_\varepsilon, \\ \mathbf{w} \text{ is } D\text{-periodic.} \end{cases} \end{aligned} \quad (7.5.38)$$

We consider the following minimization problem for the function  $\mathbf{v} \in H^{K+1}(D, \mathbb{R}^d)$  sought to approximate  $\mathbf{u}_\varepsilon^*$ :

$$\begin{aligned} \min_{\mathbf{v} \in H^{K+1}(D, \mathbb{R}^d)} J(\mathbf{w}_{\varepsilon,K}(\mathbf{v}), \mathbf{f}) \\ \text{s.t.} \begin{cases} \operatorname{div}(\mathbf{v}) = 0 \text{ in } D, \\ \mathbf{v} \text{ is } D\text{-periodic.} \end{cases} \end{aligned} \quad (7.5.39)$$

Applying lemma 7.3 to (7.5.39) in order to pass to the limit in the terms of  $J(\mathbf{w}_{\varepsilon,K}(\mathbf{v}), \mathbf{f})$  which depend on the oscillating variable  $x/\varepsilon$ , we obtain as in proposition 7.23, the existence of a functional  $J_K^*$  such that

$$\forall \mathbf{v} \in C^\infty(D, \mathbb{R}^d), J(\mathbf{w}_{\varepsilon,K}(\mathbf{v}), \mathbf{f}) = J_K^*(\mathbf{v}, \mathbf{f}, \varepsilon) + o(\varepsilon^m)$$

holds with  $m \in \mathbb{N}$  arbitrarily large. Following the derivations of section 7.4.2, the approximate energy  $J_K^*(\mathbf{v}, \mathbf{f}, \varepsilon)$  reads

$$J_K^*(\mathbf{v}, \mathbf{f}, \varepsilon) := \int_D \left( \frac{1}{2} \sum_{l,m=0}^{l+1} \varepsilon^{l+m-2} \mathbb{B}_K^{l,m} \nabla^l \mathbf{v} \nabla^m \mathbf{v} - \mathbf{f} \cdot \mathbf{v} \right) dx, \quad (7.5.40)$$

and the constant, matrix valued tensors  $\mathbb{B}_K^{l,m}$  of order  $l + m$  are defined by the formula

$$\mathbb{B}_{K,ij}^{l,m} := \int_Y \widetilde{\mathbf{N}}_{p,i}^l \cdot \widetilde{\mathbf{N}}_{p,j}^m dy, \quad 1 \leq i, j \leq d, \quad 0 \leq l, m \leq l+1. \quad (7.5.41)$$

The vector valued tensors  $\widetilde{\mathbf{N}}_{p,j}^l$  are still given as in (7.4.25) by

$$\widetilde{\mathbf{N}}_{p,j}^l(y) := \begin{cases} \partial_p \mathbf{N}_j^0(y) & \text{if } l = 0 \\ \partial_p \mathbf{N}_j^l(y) + \mathbf{N}_j^{l-1}(y) \otimes e_p & \text{if } 1 \leq l \leq K, \quad 1 \leq j, p \leq d. \\ \mathbf{N}_j^K(y) \otimes e_p & \text{if } l = K + 1. \end{cases} \quad (7.5.42)$$

Replacing  $J(\mathbf{w}_{\varepsilon,K}(\mathbf{v}), \mathbf{f})$  by  $J_K^*$  in (7.5.39) allows us to derive a well-posed homogenized equation of order  $2K + 2$  for the Stokes problem:



**Definition 7.8.** For any  $K \in \mathbb{N}$ , we call homogenized equation of order  $2K + 2$  the Euler-Lagrange equation associated with the minimization problem

$$\begin{aligned} \min \quad & J_K^*(\mathbf{v}, \mathbf{f}, \varepsilon) \\ \text{s.t.} \quad & \begin{cases} \mathbf{v} \in H^{K+1}(D, \mathbb{R}^d), \\ \operatorname{div}(\mathbf{v}) = 0 \text{ in } D, \\ \mathbf{v} \text{ is } D\text{-periodic.} \end{cases} \end{aligned} \quad (7.5.43)$$

This equation reads explicitly in terms of a higher order homogenized solution  $\mathbf{v}_K^*$  and a higher order homogenized pressure  $q_K^* \in L^2(D)$  as:

$$\begin{cases} \sum_{k=0}^{2K+2} \varepsilon^{k-2} \mathbb{D}_K^k \cdot \nabla^k \mathbf{v}_K^* + \nabla q_K^* = \mathbf{f}, \\ \operatorname{div}(\mathbf{v}_K^*) = 0, \\ \int_D q_K^* dx = 0, \\ \mathbf{v}_K^* \text{ is } D\text{-periodic,} \end{cases} \quad (7.5.44)$$

where the constant (matrix valued) tensors  $\mathbb{D}_K^k$  are defined by the following formula for any  $k \geq 0$ :

$$\forall 1 \leq i, j \leq d, \mathbb{D}_{K,ij}^k := \begin{cases} \sum_{l=0}^k (-1)^l \frac{1}{2} (\mathbb{B}_{K,ij}^{l,k-l} + \mathbb{B}_{K,ji}^{l,k-l}), & \text{if } k \text{ is even} \\ \sum_{l=0}^k (-1)^l \frac{1}{2} (\mathbb{B}_{K,ij}^{l,k-l} - \mathbb{B}_{K,ji}^{l,k-l}), & \text{if } k \text{ is odd,} \end{cases} \quad (7.5.45)$$

where the convention  $\mathbb{B}_K^{l,m} = 0$  for  $l > K + 1$  or  $m > K + 1$  is assumed.

Remembering  $\mathbb{B}_K^{K+1,K+1}$  is nonnegative according to (7.5.41), adapting the proof of proposition 7.11 yields that the bilinear form associated with the energy (7.5.40) is symmetric and positive definite, under the additional non-degeneracy assumption

$$\exists \nu > 0, \forall \boldsymbol{\xi} = \boldsymbol{\xi}_{i_1 \dots i_{K+1}, j} \in \mathbb{R}^{d^{K+1}} \times \mathbb{R}^d, \mathbb{B}_K^{K+1,K+1} \boldsymbol{\xi} \boldsymbol{\xi} \geq \nu |\boldsymbol{\xi}|^2. \quad (7.5.46)$$

From there, standard theory for saddle point problems based on the zero divergence constraint (see e.g. the textbooks [304, 301, 171, 144]) yields existence and uniqueness of a solution for (7.5.44):

**Proposition 7.35.** *Assume the dominant tensor  $\mathbb{D}_K^{2K+2} = (-1)^{K+1} \mathbb{B}_K^{K+1,K+1}$  is non-degenerate, that is there exists a constant  $\nu > 0$  such that (7.5.46) holds. Then there exists a unique velocity and pressure couple  $(\mathbf{v}_K^*, q_K^*) \in H^{K+1}(D, \mathbb{R}^d) \times L^2(D)/\mathbb{R}$  solving the higher order homogenized equation (7.5.44).*

Finally, our next proposition (which we shall use in the proof of our error estimates in the next section) states that (7.5.44) is indeed a ‘‘truncation’’ of the infinite order homogenized equation (7.5.29).

**Proposition 7.36.** *The first  $K+1$  homogenized coefficients of the homogenized equation (7.4.29) coincide with those of the formal infinite order homogenized equation (7.4.17):*

$$\forall 0 \leq k \leq K, \mathbb{D}_K^k = M^k.$$

*Proof.* We follow the proof of proposition 7.12. For  $0 \leq k, l \leq K$  and  $1 \leq i, j \leq d$ , the coefficient  $\mathbb{B}_{K,ij}^{l,k-l}$  is given by

$$\begin{aligned} \mathbb{B}_{K,ij}^{l,k-l} &= \int_Y (\partial_p \mathbf{N}_i^l + \mathbf{N}_i^{l-1} \otimes e_p) \cdot (\partial_p \mathbf{N}_j^{k-l} + \mathbf{N}_j^{k-l-1} \otimes e_p) dy \\ &= \int_Y (-\Delta_{yy} \mathbf{N}_i^l - 2\partial_p \mathbf{N}_i^{l-1} \otimes e_p - \mathbf{N}_i^{l-2} \otimes I) \cdot \mathbf{N}_j^{k-l} dy \\ &\quad + \int_Y (\partial_p \mathbf{N}_i^{l-1} \cdot \mathbf{N}_j^{k-l} \otimes e_p + \partial_p \mathbf{N}_i^l \otimes \mathbf{N}_j^{k-l-1} \otimes e_p) dy \\ &\quad + \int_Y (\mathbf{N}_i^{l-2} \cdot \mathbf{N}_j^{k-l} \otimes I + \mathbf{N}_i^{l-1} \cdot \mathbf{N}_j^{k-l-1} \otimes I) dy. \end{aligned}$$

Using the identities (7.5.33) and performing an integration by part, we may write

$$\begin{aligned} \int_Y (-\Delta_{yy} N_i^l - 2\partial_p N_i^{l-1} \otimes e_p - N_i^{l-2} \otimes I) \cdot N_j^{k-l} dy &= \int_Y (-\nabla \beta_i^l - \beta_i^{l-1} e_p \otimes e_p + M^l e_i) \cdot N_j^{k-l} dy \\ &= \int_Y [M^l e_i \cdot N_j^{k-l} - (\beta_i^l N_j^{k-l-1} + \beta_i^{l-1} N_j^{k-l}) \cdot e_p \otimes e_p] dy. \end{aligned}$$

This allows to rewrite  $\mathbb{B}_{K,ij}^{l,k-l}$  as follows:

$$\mathbb{B}_{K,ij}^{l,k-l} = \int_Y M^l e_i \cdot N_j^{k-l} dy + B_{ij}^{k,l} + B_{ij}^{k,l+1}$$

where  $B^{k,l}$  is the  $k$ -th order tensor defined by

$$B^{k,l} := \int_Y [(\partial_p N_i^{l-1} \cdot N_j^{k-l} - \beta_i^{l-1} N_j^{k-l} \cdot e_p) \otimes e_p + N_i^{l-2} \cdot N_j^{k-l} \otimes I] dy.$$

Recognizing a telescopic series, this allows to obtain

$$\sum_{l=0}^k (-1)^l \mathbb{B}_{K,ij}^{l,k-l} = (-1)^k M_{ji}^k$$

from where the result follows as in the proof of proposition 7.25.  $\square$

### 7.5.3 Error estimates and justifications of the higher order homogenization process

This section is devoted to establishing error estimates for the approximation of the oscillating solution  $(\mathbf{u}_\varepsilon, p_\varepsilon)$  by means of the truncated two-scale ansatz (7.5.36) and (7.5.37) and of the higher order homogenized equation (7.5.44). We obtain error bounds for the classical and criminal truncated ansatz (eqn. (7.5.14), (7.5.16) and (7.5.32) respectively) in propositions 7.37 and 7.38. We then prove approximation results for the solution of the homogenized equation of order  $2K+2$  (eqn. (7.5.39)) in proposition 7.39, which is our main result.

The mathematical analysis of the present situation is more difficult than the one of the perforated scalar and elasticity systems of the previous sections 7.3 and 7.4 because of the zero-divergence constraint (see e.g. [301, 304, 171, 12] for an overview of topics related to this matter); as a consequence, several of the previous arguments need to be adapted. Therefore, we start by stating a few technical results (more particularly, corollary 7.11) that we need in our proofs.

In this section, we use the following notation when referring to Sobolev spaces of  $D$ -periodic functions:

$$H_{per}^1(D_\varepsilon, \mathbb{R}^d) := \{\mathbf{v} \in H^1(D_\varepsilon, \mathbb{R}^d) \mid \mathbf{v} \text{ is } D\text{-periodic}\},$$

$$H_{per}^1(D, \mathbb{R}^d) := \{\mathbf{v} \in H^1(D, \mathbb{R}^d) \mid \mathbf{v} \text{ is } D\text{-periodic}\}.$$

The following lemma is due to [152], it establishes the existence of a continuous right inverse for the divergence  $B_\varepsilon$ —so-called a Bogovskii's operator—with a bound explicit in  $\varepsilon$  on the uniform continuity constant:

**Lemma 7.9.** *Assume (H4) and (H5). There exists a linear operator  $B_\varepsilon : L^2(D_\varepsilon) \rightarrow H_{per}^1(D_\varepsilon, \mathbb{R}^d)$  and a constant  $C$  independent of  $\varepsilon$  satisfying, for any  $\phi \in L^2(D_\varepsilon)$  such that  $\int_{D_\varepsilon} \phi dx = 0$ :*

(i)  $\operatorname{div}(B_\varepsilon \phi) = \phi$  in  $D_\varepsilon$ ,

(ii)  $B_\varepsilon \phi = 0$  on  $\partial\omega_\varepsilon$ ,

(iii)  $\|\nabla(B_\varepsilon \phi)\|_{L^2(D_\varepsilon, \mathbb{R}^{d \times d})} \leq C\varepsilon^{-1} \|\phi\|_{L^2(D_\varepsilon)}$ , for a constant  $C > 0$  independent of  $\phi$  and  $\varepsilon$ .

*Proof.* We follow the proof of [152], Lemma 2.1., which is also inspired by [14], Lemma 2.2.4. Let  $\phi \in L^2(D_\varepsilon)$  and consider the extension  $\tilde{\phi} \in L^2(D_\varepsilon)$  defined by

$$\tilde{\phi} = \begin{cases} \phi & \text{in } D_\varepsilon \\ 0 & \text{in } \omega_\varepsilon. \end{cases}$$

Let  $B : L^2(D) \rightarrow H_{per}^1(D, \mathbb{R}^d)$  be the “classical” Bogovskii’s operator in the non perforated domain  $D$ : there exists a constant  $C > 0$  independent of  $\varepsilon$  and  $\phi$  such that

$$\operatorname{div}(B\tilde{\phi}) = \tilde{\phi} \text{ in } D \quad (7.5.47)$$

$$\|\nabla(B\tilde{\phi})\|_{L^2(D, \mathbb{R}^{d \times d})} \leq C\|\tilde{\phi}\|_{L^2(D)} = C\|\phi\|_{L^2(D_\varepsilon)}. \quad (7.5.48)$$

From [12] (Theorem 2.3) (by adapting the proof to the periodicity with respect to  $D$ ), there exists a linear *restriction* operator  $q_\varepsilon : H_{per}^1(D, \mathbb{R}^d) \rightarrow H_{per}^1(D_\varepsilon, \mathbb{R}^d)$  satisfying the following properties:

1.  $\forall \mathbf{v} \in H_{per}^1(D, \mathbb{R}^d)$ ,  $q_\varepsilon \mathbf{v} = 0$  on  $\partial\omega_\varepsilon$ ;
2.  $\forall \mathbf{v} \in H_{per}^1(D, \mathbb{R}^d)$ ,  $\mathbf{v} = 0$  in  $\omega_\varepsilon \Rightarrow q_\varepsilon \mathbf{v} = \mathbf{v}$  in  $D_\varepsilon$ ;
3.  $\forall \mathbf{v} \in H_{per}^1(D, \mathbb{R}^d)$ ,  $\|q_\varepsilon \mathbf{v}\|_{L^2(D_\varepsilon, \mathbb{R}^d)} + \varepsilon\|\nabla(R_\varepsilon \mathbf{v})\|_{L^2(D_\varepsilon, \mathbb{R}^{d \times d})} \leq C(\|\mathbf{v}\|_{L^2(D, \mathbb{R}^d)} + \varepsilon\|\nabla \mathbf{v}\|_{L^2(D, \mathbb{R}^{d \times d})})$  for a constant  $C > 0$  independent of  $\varepsilon$ ;
4.  $\forall \mathbf{v} \in H_{per}^1(D, \mathbb{R}^d)$ ,  $\operatorname{div}(\mathbf{v}) = 0$  in  $D \Rightarrow \operatorname{div}(q_\varepsilon \mathbf{v}) = 0$  in  $D_\varepsilon$ .

Following an observation of [152] (see the proof of Lemma 2.1) and looking closely to the construction of the operator  $R_\varepsilon$  in [12], it can be shown that  $R_\varepsilon$  satisfies in fact the following property, which is stronger than (4):

$$(4^*) \quad \forall \mathbf{v} \in H_{per}^1(D, \mathbb{R}^d), \operatorname{div}(\mathbf{v}) = 0 \text{ in } \omega_\varepsilon \Rightarrow \operatorname{div}(q_\varepsilon \mathbf{v}) = \operatorname{div}(\mathbf{v}) \text{ in } D_\varepsilon.$$

Following [152], we set

$$B_\varepsilon \phi := q_\varepsilon(B\tilde{\phi}).$$

The resulting operator  $B_\varepsilon$  satisfies (i) to (iii). Indeed:

- (i)  $\operatorname{div}(B\tilde{\phi}) = 0$  in  $\omega_\varepsilon$  yields, according to (4\*):

$$\operatorname{div}(q_\varepsilon(B\tilde{\phi})) = \operatorname{div}(B\tilde{\phi}) = \tilde{\phi} = \phi \text{ in } D_\varepsilon;$$

- (ii) is a straightforward consequence of the property (2);

- (iii) is obtained by writing

$$\begin{aligned} \varepsilon\|\nabla(B_\varepsilon \phi)\|_{L^2(D_\varepsilon, \mathbb{R}^{d \times d})} &= \varepsilon\|\nabla(q_\varepsilon(B\tilde{\phi}))\|_{L^2(D_\varepsilon, \mathbb{R}^{d \times d})} \\ &\leq C(\|B\tilde{\phi}\|_{L^2(D, \mathbb{R}^d)} + \varepsilon\|\nabla(B\tilde{\phi})\|_{L^2(D, \mathbb{R}^{d \times d})}) \\ &\leq C\|\nabla(B\tilde{\phi})\|_{L^2(D, \mathbb{R}^{d \times d})} \quad (\text{from the Poincaré inequality in } D) \\ &\leq C\|\phi\|_{L^2(D)} \quad (\text{from (7.5.48)}). \end{aligned}$$

□

Let us remark that the previous lemma implies an estimation of the inf-sup constant associated with the perforated Stokes problem in terms of  $\varepsilon$ :

**Corollary 7.10** (inf-sup constant for the perforated Stokes problem). *There exists a constant  $C > 0$  independent of  $\varepsilon$  such that*

$$\inf_{\substack{\phi \in L^2(D_\varepsilon) \\ \int_{D_\varepsilon} \phi dx = 0}} \sup_{\substack{\mathbf{v} \in H_{per}^1(D_\varepsilon, \mathbb{R}^d) \\ \mathbf{v} = 0 \text{ on } \partial\omega_\varepsilon}} \int_{D_\varepsilon} \frac{\phi \operatorname{div}(\mathbf{v})}{\|\phi\|_{L^2(D_\varepsilon)} \|\nabla \mathbf{v}\|_{L^2(D_\varepsilon, \mathbb{R}^{d \times d})}} dx \geq C\varepsilon.$$

*Proof.* The result is obtained by writing, for any  $\phi \in L^2(D_\varepsilon)$  satisfying  $\int_{D_\varepsilon} \phi dx = 0$ :

$$\begin{aligned} \sup_{\substack{\mathbf{v} \in H_{per}^1(D_\varepsilon, \mathbb{R}^d) \\ \mathbf{v} = 0 \text{ on } \partial\omega_\varepsilon}} \int_{D_\varepsilon} \frac{\phi \operatorname{div}(\mathbf{v})}{\|\phi\|_{L^2(D_\varepsilon)} \|\nabla \mathbf{v}\|_{L^2(D_\varepsilon, \mathbb{R}^{d \times d})}} dx &\geq \int_{D_\varepsilon} \frac{\phi \operatorname{div}(B_\varepsilon \phi)}{\|\phi\|_{L^2(D)} \|\nabla(B_\varepsilon \phi)\|_{L^2(D_\varepsilon, \mathbb{R}^{d \times d})}} dx \\ &\geq \int_{D_\varepsilon} \frac{\phi^2}{\|\phi\|_{L^2(D_\varepsilon)} C\varepsilon^{-1} \|\phi\|_{L^2(D_\varepsilon)}} dx = \frac{1}{C}\varepsilon, \end{aligned}$$

where the constant  $C$  above is that of the point (iii) of lemma 7.9. □

**Corollary 7.11.** *Assume (H4) and (H5). For any  $\mathbf{h} \in L^2(D_\varepsilon, \mathbb{R}^d)$  and  $g \in L^2(D_\varepsilon)$  satisfying  $\int_{D_\varepsilon} g dx = 0$ , let  $(\mathbf{v}, \phi) \in H^1(D_\varepsilon, \mathbb{R}^d) \times L^2(D_\varepsilon)$  be the unique solution to the Stokes problem*

$$\left\{ \begin{array}{l} -\Delta \mathbf{v} + \nabla \phi = \mathbf{h} \text{ in } D_\varepsilon \\ \operatorname{div}(\mathbf{v}) = g \text{ in } D_\varepsilon \\ \int_{D_\varepsilon} \phi dx = 0 \\ \mathbf{v} = 0 \text{ on } \partial\omega_\varepsilon \\ \mathbf{v} \text{ is } D\text{-periodic.} \end{array} \right. \quad (7.5.49)$$

There exists a constant  $C$  independent of  $\varepsilon$ ,  $\mathbf{h}$  and  $g$  such that

$$\|\nabla \mathbf{v}\|_{L^2(D_\varepsilon, \mathbb{R}^{d \times d})} + \varepsilon \|\phi\|_{L^2(D_\varepsilon)} \leq C(\varepsilon \|\mathbf{h}\|_{L^2(D_\varepsilon, \mathbb{R}^d)} + \varepsilon^{-1} \|g\|_{L^2(D_\varepsilon)}), \quad (7.5.50)$$

*Proof.* We use the operator  $B_\varepsilon$  of lemma 7.9 to lift the divergence of  $\mathbf{v}$ . Let  $\mathbf{w} := \mathbf{v} - B_\varepsilon g \in H_{per}^1(D_\varepsilon, \mathbb{R}^d)$ , then it holds:

$$\left\{ \begin{array}{l} \operatorname{div}(\mathbf{w}) = 0 \text{ in } D_\varepsilon, \\ \mathbf{w} = 0 \text{ on } \partial\omega_\varepsilon. \end{array} \right.$$

After an integration by part, we obtain:

$$\begin{aligned} \|\nabla \mathbf{w}\|_{L^2(D_\varepsilon, \mathbb{R}^{d \times d})}^2 &= \int_{D_\varepsilon} \mathbf{h} \cdot \mathbf{w} dx - \int_{D_\varepsilon} \nabla(B_\varepsilon g) : \nabla \mathbf{w} dx \\ &\leq \|\mathbf{h}\|_{L^2(D_\varepsilon, \mathbb{R}^d)} \|\mathbf{w}\|_{L^2(D_\varepsilon, \mathbb{R}^d)} + \|\nabla(B_\varepsilon g)\|_{L^2(D_\varepsilon, \mathbb{R}^{d \times d})} \|\nabla \mathbf{w}\|_{L^2(D_\varepsilon, \mathbb{R}^{d \times d})} \\ &\leq C(\varepsilon \|\mathbf{h}\|_{L^2(D_\varepsilon, \mathbb{R}^d)} + \|\nabla(B_\varepsilon g)\|_{L^2(D_\varepsilon, \mathbb{R}^{d \times d})}) \|\nabla \mathbf{w}\|_{L^2(D_\varepsilon, \mathbb{R}^{d \times d})}, \end{aligned}$$

where the last inequality is a consequence of lemma 7.2. Therefore, simplifying by  $\|\nabla \mathbf{w}\|_{L^2(D_\varepsilon, \mathbb{R}^{d \times d})}$  and using the point (iii) of lemma 7.9, we obtain

$$\begin{aligned} \|\nabla \mathbf{v}\|_{L^2(D_\varepsilon, \mathbb{R}^{d \times d})} &\leq \|\nabla \mathbf{w}\|_{L^2(D_\varepsilon, \mathbb{R}^{d \times d})} + \|\nabla(B_\varepsilon g)\|_{L^2(D_\varepsilon, \mathbb{R}^{d \times d})} \\ &\leq C(\varepsilon \|\mathbf{h}\|_{L^2(D_\varepsilon, \mathbb{R}^d)} + \varepsilon^{-1} \|g\|_{L^2(D_\varepsilon)}), \end{aligned}$$

which proves the first part of the bound (7.5.50) on  $\nabla \mathbf{v}$ . The bound on the pressure is then classically obtained thanks to a reasoning analogous to that of the one of corollary 7.10. Using  $B_\varepsilon \phi$  a test function, we obtain:

$$\begin{aligned} \|\phi\|_{L^2(D_\varepsilon)}^2 &= \int_{D_\varepsilon} \phi \operatorname{div}(B_\varepsilon \phi) dx = - \int_{D_\varepsilon} \nabla \phi \cdot B_\varepsilon \phi dx \\ &= \int_{D_\varepsilon} (-\Delta \mathbf{v} - \mathbf{h}) \cdot B_\varepsilon \phi dx = \int_{D_\varepsilon} (\nabla \mathbf{v} \cdot \nabla(B_\varepsilon \phi) - \mathbf{h} \cdot B_\varepsilon \phi) dx \\ &\leq \|\nabla \mathbf{v}\|_{L^2(D_\varepsilon, \mathbb{R}^{d \times d})} \|\nabla(B_\varepsilon \phi)\|_{L^2(D_\varepsilon, \mathbb{R}^{d \times d})} + \|\mathbf{h}\|_{L^2(D_\varepsilon, \mathbb{R}^d)} \|B_\varepsilon \phi\|_{L^2(D_\varepsilon, \mathbb{R}^d)} \\ &\leq C(\varepsilon \|\mathbf{h}\|_{L^2(D_\varepsilon, \mathbb{R}^d)} + \varepsilon^{-1} \|g\|_{L^2(D_\varepsilon)}) \|\nabla(B_\varepsilon \phi)\|_{L^2(D_\varepsilon, \mathbb{R}^{d \times d})} \\ &\leq C\varepsilon^{-1} (\varepsilon \|\mathbf{h}\|_{L^2(D_\varepsilon, \mathbb{R}^d)} + \varepsilon^{-1} \|g\|_{L^2(D_\varepsilon)}) \|\phi\|_{L^2(D_\varepsilon)}, \end{aligned}$$

whence the result.  $\square$

The previous lemma allows to establish a convergence result analogous to proposition 7.4 for the truncation of the ‘‘classical’’ ansatz (7.5.2), which are defined for any  $K \in \mathbb{N}$  by

$$\mathbf{u}_{\varepsilon, K}(x) := \sum_{i=0}^K \varepsilon^{i+2} \mathbf{u}_i(x, x/\varepsilon), \quad p_{\varepsilon, K}(x) := \sum_{i=0}^K \varepsilon^i (p_i^*(x) + \varepsilon p_i(x, x/\varepsilon)). \quad (7.5.51)$$

We recall equations (7.5.5) and (7.5.12) for the definition of  $\mathbf{u}_i$  and  $p_i$ . The next result is the extension of proposition 7.4 to the context of the Stokes system:

**Proposition 7.37.** *Assume (H4) and (H5). Let  $(\mathbf{u}_\varepsilon, p_\varepsilon) \in H_{per}^1(D_\varepsilon, \mathbb{R}^d) \times L^2(D_\varepsilon)$  the solution to (7.5.1). Let  $(\mathbf{u}_{\varepsilon, K}, p_{\varepsilon, K})$  the truncated ansatz at rank  $K$  defined by (7.5.51). There exists a constant  $C_K > 0$  independent of  $\varepsilon$  and  $\mathbf{f}$  such that*

$$\|\mathbf{u}_\varepsilon - \mathbf{u}_{\varepsilon, K}\|_{L^2(D_\varepsilon, \mathbb{R}^d)} + \varepsilon \|\nabla(\mathbf{u}_\varepsilon - \mathbf{u}_{\varepsilon, K})\|_{L^2(D_\varepsilon, \mathbb{R}^{d \times d})} \leq C_K \varepsilon^{K+3} \|\mathbf{f}\|_{H^{K+3}(D_\varepsilon, \mathbb{R}^d)}, \quad (7.5.52)$$

$$\|p_\varepsilon - p_{\varepsilon, K}\|_{L^2(D_\varepsilon)} \leq C_K \varepsilon^{K+1} \|\mathbf{f}\|_{H^{K+3}(D_\varepsilon, \mathbb{R}^d)}. \quad (7.5.53)$$

*Proof.* By using that  $\mathbf{u}_i(x, y)$  and  $p_i(x, y)$  are the solutions to the cascade of equations (7.5.5), we obtain, for any  $K' \in \mathbb{N}$ ,

$$\begin{aligned} & \forall x \in D_\varepsilon, (-\Delta \mathbf{u}_{\varepsilon, K'} + \nabla p_{\varepsilon, K'})(x) \\ &= \sum_{i=-2}^{K'-2} \varepsilon^{i+2} (-\Delta_{xx} \mathbf{u}_i - \Delta_{xy} \mathbf{u}_{i+1} - \Delta_{yy} \mathbf{u}_{i+2})(x, x/\varepsilon) + \sum_{i=-2}^{K'-2} \varepsilon^{i+2} (\nabla_x p_{i+2}^* + \nabla_x p_{i+1} + \nabla_y p_{i+2})(x, x/\varepsilon) \\ & \quad - (\varepsilon^{K'+1} \Delta_{xx} \mathbf{u}_{K'-1} + \varepsilon^{K'+2} \Delta_{xx} \mathbf{u}_{K'} + \varepsilon^{K'+1} \Delta_{xy} \mathbf{u}_{K'})(x, x/\varepsilon) + \varepsilon^{K'+1} (\nabla_x p_K)(x, x/\varepsilon) \\ &= \mathbf{f} - (\varepsilon^{K'+1} \Delta_{xx} \mathbf{u}_{K'-1} + \varepsilon^{K'+2} \Delta_{xx} \mathbf{u}_{K'} + \varepsilon^{K'+1} \Delta_{xy} \mathbf{u}_{K'} + \varepsilon^{K'+1} \nabla_x p_K)(x, x/\varepsilon). \\ & \quad \operatorname{div}(\mathbf{u}_{\varepsilon, K'})(x) = \sum_{i=-2}^{K'-2} \varepsilon^{i+3} (\operatorname{div}_x \mathbf{u}_{i+1} + \operatorname{div}_y \mathbf{u}_{i+2})(x, x/\varepsilon) + \varepsilon^{K'+2} (\operatorname{div}_x \mathbf{u}_{K'})(x, x/\varepsilon) \\ & \quad = \varepsilon^{K'+2} (\operatorname{div}_x \mathbf{u}_{K'})(x, x/\varepsilon). \end{aligned}$$

From standard regularity theory, it can be easily shown that

$$\begin{aligned} & \left\| (\varepsilon^{K'+1} \Delta_{xx} \mathbf{u}_{K'-1} + \varepsilon^{K'+2} \Delta_{xx} \mathbf{u}_{K'} + \varepsilon^{K'+1} \Delta_{xy} \mathbf{u}_{K'} + \varepsilon^{K'+1} \nabla_x p_K)(\cdot, \cdot/\varepsilon) \right\|_{L^2(D_\varepsilon, \mathbb{R}^d)} \\ & \leq C_K \varepsilon^{K'+1} \|\mathbf{f}\|_{H^{K'+2}(D, \mathbb{R}^d)}, \\ & \|\varepsilon^{K'+2} (\operatorname{div}_x \mathbf{u}_{K'})(\cdot, \cdot/\varepsilon)\|_{L^2(D_\varepsilon)} \leq C_K \varepsilon^{K'+2} \|\mathbf{f}\|_{H^{K'+1}(D, \mathbb{R}^d)}. \end{aligned}$$

Therefore, applying the result of corollary 7.11 to  $(\mathbf{v}_\varepsilon, q_\varepsilon) := (\mathbf{u}_\varepsilon - \mathbf{u}_{\varepsilon, K'}, p_\varepsilon - p_{\varepsilon, K'})$  with  $K' := K + 1$ , we infer the existence of a constant  $C_K > 0$  such that

$$\begin{aligned} & \|\mathbf{u}_\varepsilon - \mathbf{u}_{\varepsilon, K+1}\|_{H^1(D_\varepsilon, \mathbb{R}^d)} \leq C_K (\varepsilon^{K+3} \|\mathbf{f}\|_{H^{K+3}(D, \mathbb{R}^d)} + \varepsilon^{K+2} \|\mathbf{f}\|_{H^{K+2}(D, \mathbb{R}^d)}) \\ & \|p_\varepsilon - p_{\varepsilon, K+1}\|_{L^2(D_\varepsilon)} \leq C_K (\varepsilon^{K+2} \|\mathbf{f}\|_{H^{K+3}(D, \mathbb{R}^d)} + \varepsilon^{K+1} \|\mathbf{f}\|_{H^{K+2}(D, \mathbb{R}^d)}). \end{aligned}$$

Finally, let us remark that  $\mathbf{u}_{\varepsilon, K+1}$  and  $p_{\varepsilon, K+1}$  are high order corrections of  $\mathbf{u}_{\varepsilon, K}$  and  $p_{\varepsilon, K}$ :

$$\begin{aligned} & \|\mathbf{u}_{\varepsilon, K+1} - \mathbf{u}_{\varepsilon, K}\|_{H^1(D_\varepsilon, \mathbb{R}^d)} = \|\varepsilon^{K+3} \mathbf{u}_{K+1}(\cdot, \cdot/\varepsilon)\|_{H^1(D_\varepsilon, \mathbb{R}^d)} \leq C_K \varepsilon^{K+2} \|\mathbf{f}\|_{H^{K+1}(D, \mathbb{R}^d)} \\ & \|p_{\varepsilon, K+1} - p_{\varepsilon, K}\|_{L^2(D)} = \left\| \varepsilon^{K+1} (p_{K+1}^* + \varepsilon p_{K+1}(\cdot, \cdot/\varepsilon)) \right\|_{L^2(D_\varepsilon)} \leq C_K \varepsilon^{K+1} \|\mathbf{f}\|_{H^{K+1}(D, \mathbb{R}^d)}. \end{aligned}$$

The result follows by using the triangle inequality.  $\square$

**Remark 7.29.** The term  $\varepsilon^{K+1} p_K(x, x/\varepsilon)$  can be removed in the truncated ansatz  $p_{\varepsilon, K}$  of (7.5.51) because it is of order  $\varepsilon^{K+1}$  in the  $L^2$  norm.

**Remark 7.30.** As a result of the scaling  $\varepsilon^{-1}$  in corollary 7.11, we pay a factor  $\varepsilon^{-1}$  in the error induced by the non zero divergence constraint. However we are able to obtain the right order of  $\varepsilon$  in the error estimates (7.5.52) and (7.5.53) thanks to the use of higher order terms. This phenomenon is quite classical in the truncation analysis of two-scale expansions where the higher order terms of the ansatz are used to establish the estimate and removed at the end (see e.g. [74, 19]).

**Remark 7.31.** In [12], the following convergence is obtained (in a setting involving much lower regularity for  $\mathbf{f}$  and  $D$ ):

$$P_\varepsilon \rightarrow p_0^* \text{ in } L^2(D),$$

where  $P_\varepsilon$  is the extension of the pressure  $p_\varepsilon$  to the whole domain  $D$  defined as follows. Denote by  $(\mathcal{P}_\varepsilon^\varepsilon)$  and the collection of cells of size  $\varepsilon$  and by  $(Y_i^\varepsilon)$  their corresponding fluid part. The extension of the pressure  $P_\varepsilon$  is defined by:

$$P_\varepsilon = \left\{ \begin{array}{l} p_\varepsilon \text{ in } D_\varepsilon, \\ \frac{1}{|Y_i^\varepsilon|} \int_{Y_i^\varepsilon} p_\varepsilon dx \text{ in } \mathcal{P}_\varepsilon^\varepsilon \setminus \mathcal{Y}_\varepsilon^\varepsilon \text{ for each } i. \end{array} \right. \quad (7.5.54)$$

This convergence results can be retrieved formally very naturally in view of the ansatz (7.5.2). Indeed, our estimates of proposition 7.37 yield

$$\|p_\varepsilon - p_0^*\|_{L^2(D_\varepsilon)} \leq C\varepsilon$$

for a constant independent of  $C$ . Furthermore, if  $x \in Y_i^\varepsilon$ , it can also be expected, with  $q \geq 0$  as large as one likes:

$$\frac{1}{|Y_i^\varepsilon|} \int_{Y_i^\varepsilon} p_\varepsilon dx \simeq \sum_{i=0}^{+\infty} \varepsilon^i \left( \frac{1}{|Y|} \int_Y p_i^*(x) dy + \varepsilon \frac{1}{|Y|} \int_Y p_i(x, y) dy \right) + o(\varepsilon^q) = p_0^*(x) + o(\varepsilon),$$

which suggests indeed that  $P_\varepsilon \rightarrow p_0^*$  in the whole domain  $D = D_\varepsilon \cup \omega_\varepsilon$ .

We now show that  $\mathbf{u}_\varepsilon$  and  $p_\varepsilon$  can be indeed approximated by criminal expansions of the form of (7.5.14) and (7.5.16); the following result is the analogous of lemma 7.5 for the present Stokes context.

**Proposition 7.38.** *Let  $\mathbf{u}_{\varepsilon,K}^*, p_{\varepsilon,K}^*$  the averages of the truncated expansions  $\mathbf{u}_{\varepsilon,K}$  and  $p_{\varepsilon,K}$  of (7.5.51):*

$$\forall x \in D, \mathbf{u}_{\varepsilon,K}^*(x) := \sum_{k=0}^K \varepsilon^{k+2} \mathbf{u}_i^*(x), \quad p_{\varepsilon,K}^*(x) := \sum_{k=0}^K \varepsilon^k p_k^*(x).$$

There exists a constant  $C_K(\mathbf{f})$  independent of  $\varepsilon$  such that

$$\left\| \mathbf{u}_\varepsilon - \sum_{k=0}^K \varepsilon^k N^k(\cdot/\varepsilon) \cdot \nabla^k \mathbf{u}_{\varepsilon,K}^* \right\|_{H^1(D, \mathbb{R}^d)} \leq C_K(\mathbf{f}) \varepsilon^{K+2} \quad (7.5.55)$$

$$\left\| p_\varepsilon - \left( p_{\varepsilon,K}^* + \sum_{k=0}^{K-1} \varepsilon^{k-1} \beta^k(\cdot/\varepsilon) \cdot \nabla^k \mathbf{u}_{\varepsilon,K}^* \right) \right\|_{L^2(D)} \leq C_K(\mathbf{f}) \varepsilon^{K+1}. \quad (7.5.56)$$

*Proof.* The proof of the error bound for (7.5.55) follows exactly the same lines of lemma 7.5. The result is obtained by using the estimates (7.5.52) and (7.5.53) together with the identities (7.5.31).  $\square$

**Remark 7.32.** We need only  $K-1$  derivatives in the truncated criminal ansatz (7.5.56) for the pressure, because the term of highest order has a norm of order  $\varepsilon^K$  (recall  $\mathbf{u}_{\varepsilon,K}^*$  has a norm of order  $\varepsilon^2$ ).

We now turn to our main result which states that the solution  $(\mathbf{v}_K^*, q_K^*)$  of the homogenized equation (7.5.39) of order  $2K+2$  yields higher order approximations of  $(\mathbf{u}_\varepsilon, p_\varepsilon)$  with error bounds similar to those of (7.5.55) and (7.5.56). Our proof is slightly different from that of proposition 7.13 because the variational framework is different and the treatment of the divergence constraint is more delicate. In order to prove our result, we need beforehand the following regularity estimate for the solution  $(\mathbf{v}_K^*, q_K^*)$ :

**Lemma 7.10.** *The solution  $(\mathbf{v}_K^*, q_K^*)$  of (7.5.44) is smooth and for any  $m \in \mathbb{N}$ , there exists a constant  $C_m(\mathbf{f})$  depending only on  $m$  and  $\mathbf{f}$  such that*

$$\|\mathbf{v}_{K^*}\|_{H^m(D, \mathbb{R}^d)} \leq C_m(\mathbf{f}) \varepsilon^2.$$

*Proof.* This can be obtained by solving (7.5.44) explicitly with Fourier series.  $\square$

**Proposition 7.39.** *Let  $(\mathbf{v}_K^*, q_K^*)$  be the unique solution to the homogenized equation (7.5.44) of order  $2K+2$ . There exists a constant  $C_K(\mathbf{f})$  depending only on  $K$  and  $\mathbf{f}$  (and a priori on the shape of the hole  $(\eta T)$ ) such that the following error estimates hold:*

$$\begin{aligned} \left\| \mathbf{u}_\varepsilon - \sum_{k=0}^K \varepsilon^k N^k(\cdot/\varepsilon) \cdot \nabla^k \mathbf{v}_K^* \right\|_{L^2(D, \mathbb{R}^d)} &\leq C_K(\mathbf{f}) \varepsilon^{K+3}, \\ \left\| \mathbf{u}_\varepsilon - \sum_{k=0}^K \varepsilon^k N^k(\cdot/\varepsilon) \cdot \nabla^k \mathbf{v}_K^* \right\|_{H^1(D, \mathbb{R}^d)} &\leq C_K(\mathbf{f}) \varepsilon^{K+2}, \\ \left\| p_\varepsilon - \left( q_K^* + \sum_{k=0}^{K-1} \varepsilon^{k-1} \beta^k(\cdot/\varepsilon) \cdot \nabla^k \mathbf{v}_K^* \right) \right\|_{L^2(D)} &\leq C_K(\mathbf{f}) \varepsilon^{K+1}. \end{aligned}$$

*Proof.* Let us compute

$$\begin{aligned}
 -\Delta \mathbf{w}_{\varepsilon, K+1}(\mathbf{v}_K^*) &= \sum_{k=0}^K \varepsilon^{k-2} (-\Delta \mathbf{N}_j^k - 2\partial_l \mathbf{N}_j^{k-1} \otimes \mathbf{e}_l - \mathbf{N}_j^{k-2} \otimes I)(\cdot/\varepsilon) \cdot \nabla^k v_{K,j}^* \\
 &\quad - \varepsilon^{K-1} (2\partial_l \mathbf{N}_j^K \otimes \mathbf{e}_l + \mathbf{N}_j^{K-1} \otimes I)(\cdot/\varepsilon) \cdot \nabla^{K+1} v_{K,j}^* - \varepsilon^K \mathbf{N}_j^K(\cdot/\varepsilon) \otimes I \cdot \nabla^{K+2} v_{K,j}^* \\
 &\quad - \varepsilon^{K+1} \Delta(\mathbf{N}_j^{K+1}(\cdot/\varepsilon) \cdot \nabla^{K+1} v_{K,j}^*), \\
 \nabla q_{\varepsilon, K}(q_K^*, \mathbf{v}_K^*) &= \nabla q_K^* + \sum_{k=0}^K \varepsilon^{k-2} \mathbf{e}_l (\partial_l \beta_j^k + \beta_j^{k-1} \otimes \mathbf{e}_l)(\cdot/\varepsilon) \cdot \nabla^k v_{K,j}^* + \varepsilon^{K-1} \mathbf{e}_l (\beta_j^K \otimes \mathbf{e}_l)(\cdot/\varepsilon) \cdot \nabla^K v_{K,j}^*.
 \end{aligned}$$

Summing the two equations and using (7.5.33), we obtain

$$\begin{aligned}
 -\Delta \mathbf{w}_{\varepsilon, K}(\mathbf{v}_K^*) + \nabla q_{\varepsilon, K}(q_K^*, \mathbf{v}_K^*) &= \sum_{k=0}^K \varepsilon^{k-2} M^k \cdot \nabla^k \mathbf{v}_K^* + \nabla q_K^* \\
 &\quad - \varepsilon^{K-1} (2\partial_l \mathbf{N}_j^K \otimes \mathbf{e}_l + \mathbf{N}_j^{K-1} \otimes I)(\cdot/\varepsilon) \cdot \nabla^{K+1} v_{K,j}^* - \varepsilon^K \mathbf{N}_j^K(\cdot/\varepsilon) \otimes I \cdot \nabla^{K+2} v_{K,j}^* \\
 &\quad - \varepsilon^{K+1} \Delta(\mathbf{N}_j^{K+1}(\cdot/\varepsilon) \cdot \nabla^{K+1} v_{K,j}^*). \quad (7.5.57)
 \end{aligned}$$

We now use proposition 7.36 and (7.5.44) to rewrite

$$\sum_{k=0}^K \varepsilon^{k-2} M^k \cdot \nabla^k \mathbf{v}_K^* + \nabla q_K^* = \mathbf{f} - \sum_{k=K+1}^{2K+2} \varepsilon^{k-2} \mathbb{D}_K^k \cdot \nabla^k \mathbf{v}_K^*. \quad (7.5.58)$$

Let us also compute the divergence of  $\mathbf{w}_{\varepsilon, K+1}(\mathbf{v}_K^*)$ :

$$\begin{aligned}
 \operatorname{div}(\mathbf{w}_{\varepsilon, K+1}(\mathbf{v}_K^*)) &= \sum_{k=0}^{K+1} \varepsilon^{k-1} (\operatorname{div}(\mathbf{N}_j^k)(\cdot/\varepsilon) \cdot \nabla^k v_{K,j}^* + (\mathbf{N}_j^{k-1}(\cdot/\varepsilon) \cdot \mathbf{e}_l \otimes \mathbf{e}_l) \cdot \nabla^k v_{K,j}^*) \\
 &\quad + \varepsilon^{K+1} \mathbf{N}_j^{K+1}(\cdot/\varepsilon) \cdot \mathbf{e}_l \otimes \nabla^{K+2} v_{K,j}^* \\
 &= \sum_{k=0}^{K+1} \varepsilon^{k-1} (-\langle \mathbf{N}_j^{k-1}(\cdot/\varepsilon) \rangle - \langle \mathbf{N}_j^{k-1} \rangle) \cdot \mathbf{e}_l \otimes \mathbf{e}_l \cdot \nabla^k v_{K,j}^* + (\mathbf{N}_j^{k-1}(\cdot/\varepsilon) \cdot \mathbf{e}_l \otimes \mathbf{e}_l) \cdot \nabla^k v_{K,j}^* \\
 &\quad + \varepsilon^{K+1} \mathbf{N}_j^{K+1}(\cdot/\varepsilon) \cdot \mathbf{e}_l \otimes \mathbf{e}_l \cdot \nabla^{K+2} v_{K,j}^* \\
 &= \operatorname{div}(\mathbf{v}_K^*) + \varepsilon^{K+1} \mathbf{N}_j^{K+1}(\cdot/\varepsilon) \cdot \mathbf{e}_l \otimes \mathbf{e}_l \cdot \nabla^{K+2} v_{K,j}^* = \varepsilon^{K+1} \mathbf{N}_j^{K+1}(\cdot/\varepsilon) \otimes \mathbf{e}_l \cdot \nabla^{K+2} v_{K,j}^*. \quad (7.5.59)
 \end{aligned}$$

We now apply corollary 7.11 to estimate  $(\mathbf{v}_\varepsilon, q_\varepsilon) := (\mathbf{u}_\varepsilon - \mathbf{w}_{\varepsilon, K+1}(\mathbf{v}_K^*), p_\varepsilon - q_{\varepsilon, K}(q_K^*, \mathbf{v}_K^*))$ . Equations (7.5.1) and (7.5.57) to (7.5.59) imply that  $(\mathbf{v}_\varepsilon, q_\varepsilon)$  solves the following Stokes system:

$$\left\{ \begin{array}{l} -\Delta \mathbf{v}_\varepsilon + \nabla q_\varepsilon = \mathbf{h}_\varepsilon \text{ in } D_\varepsilon \\ \operatorname{div}(\mathbf{v}_\varepsilon) = g_\varepsilon \in D_\varepsilon \\ \mathbf{v}_\varepsilon = 0 \text{ on } \partial\omega_\varepsilon, \\ \int_D q_\varepsilon = 0, \\ \mathbf{v}_\varepsilon \text{ is } D\text{-periodic,} \end{array} \right.$$

where the source functions  $\mathbf{h}_\varepsilon$  and  $g_\varepsilon$  are given by

$$\mathbf{h}_\varepsilon := - \sum_{k=K+1}^{2K+2} \varepsilon^{k-2} \mathbb{D}_K^k \cdot \nabla^k \mathbf{v}_K^* - \varepsilon^{K-1} (2\partial_l \mathbf{N}_j^K \otimes \mathbf{e}_l + \mathbf{N}_j^{K-1} \otimes I)(\cdot/\varepsilon) \cdot \nabla^{K+1} v_{K,j}^* \quad (7.5.60)$$

$$\begin{aligned}
 &\quad - \varepsilon^K \mathbf{N}_j^K(\cdot/\varepsilon) \otimes I \cdot \nabla^{K+2} v_{K,j}^* - \varepsilon^{K+1} \Delta(\mathbf{N}_j^{K+1}(\cdot/\varepsilon) \cdot \nabla^{K+1} v_{K,j}^*), \\
 g_\varepsilon &:= \varepsilon^{K+1} \mathbf{N}_j^{K+1}(\cdot/\varepsilon) \otimes \mathbf{e}_l \cdot \nabla^{K+2} v_{K,j}^*. \quad (7.5.61)
 \end{aligned}$$

Using the result of lemma 7.10, we infer the existence of a constant  $C_K(\mathbf{f})$  independent of  $\varepsilon$  such that

$$\|\mathbf{h}_\varepsilon\|_{L^2(D, \mathbb{R}^d)} \leq C_K(\mathbf{f}) \varepsilon^{K+1},$$

$$\|g_\varepsilon\|_{L^2(D)} \leq C_K(\mathbf{f})\varepsilon^{K+3}.$$

Therefore, the estimates of [corollary 7.11](#) yield

$$\begin{aligned} \|\mathbf{u}_\varepsilon - \mathbf{w}_{\varepsilon, K+1}(\mathbf{v}_K^*)\|_{H^1(D, \mathbb{R}^d)} &= \|\mathbf{v}_\varepsilon\|_{H^1(D, \mathbb{R}^d)} \leq C_K(\mathbf{f})\varepsilon^{K+2}, \\ \|p_\varepsilon - q_{\varepsilon, K}(\mathbf{v}_K^*, q_K^*)\|_{L^2(D)} &= \|q_\varepsilon\|_{L^2(D)} \leq C_K(\mathbf{f})\varepsilon^{K+1}. \end{aligned}$$

The final result follows from the triangle's and the following estimates:

$$\begin{aligned} \|\mathbf{w}_{\varepsilon, K+1}(\mathbf{v}_K^*) - \mathbf{w}_{\varepsilon, K}(\mathbf{v}_K^*)\|_{H^1(D, \mathbb{R}^d)} &= \|\varepsilon^{K+1} N^{K+1} \cdot \nabla^{K+1} \mathbf{v}_K^*\|_{H^1(D, \mathbb{R}^d)} \leq C_K(\mathbf{f})\varepsilon^{K+2}, \\ \|q_{\varepsilon, K}(\mathbf{v}_K^*, p_K^*) - q_{\varepsilon, K-1}(\mathbf{v}_K^*, q_K^*)\|_{L^2(D)} &= \|\varepsilon^{K-1} \beta^K \cdot \nabla^K \mathbf{v}_K^*\|_{L^2(D)} \leq C_K(\mathbf{f})\varepsilon^{K+1}. \end{aligned}$$

□

**Remark 7.33.** The above proof further highlights that there is not a unique way to derive well-posed higher order homogenized equations such as [\(7.5.44\)](#): what is important is the fact that the first  $K + 1$  coefficients of the equation satisfied by  $(\mathbf{v}_K^*, q_K^*)$  are those of the “infinite order” homogenized equation [\(7.5.29\)](#):  $\mathbb{D}_K^k = M^k$  for  $0 \leq k \leq K + 1$ . More precisely, the approximation result of [proposition 7.39](#) holds provided [lemma 7.10](#) holds and

$$\left\| \sum_{k=0}^K \varepsilon^{k-2} M^k \cdot \mathbf{v}_K^* + \nabla q_K^* - \mathbf{f} \right\| \leq C_K(\mathbf{f})\varepsilon^{K+1}.$$

The coefficients  $\mathbb{D}_K^k$  for  $K + 1 \leq k \leq 2K + 2$  are selected in such a way that [\(7.5.39\)](#) is well posed, which is the case when defining  $\mathbb{D}_K^k$  from the energy minimization principle [\(7.5.43\)](#). In the next [section 7.5.4](#), we provide evidences that [\(7.5.44\)](#) is a “well-behaved” equation when it is obtained from the minimization principle, because it contains all three classical homogenized regimes in the low-volume fraction limit where the obstacle's size  $\eta \rightarrow 0$ .

#### 7.5.4 Low volume fraction limits when the size of the obstacle tends to 0

This section is devoted to the study of the asymptotics of the tensors  $\mathcal{X}^{k*}$ ,  $M^k$  and  $\mathbb{D}_K^k$  in the low volume fraction limit where the obstacle's size vanishes, i.e.  $\eta \rightarrow 0$ . In this whole subsection, we assume again, for simplicity, that the space dimension is greater than 3:

$$d \geq 3.$$

The notation convention is that of [sections 7.3.5](#) and [7.4.3](#). We recall the inequalities in [lemma 7.7](#) that we are going to use extensively. We also need the following technical result which yields estimates of Stokes solutions in the domain  $\eta^{-1}P \setminus T$  uniform in  $\eta \rightarrow 0$ .

**Lemma 7.11.** *Consider  $\mathbf{h} \in L^2(\eta^{-1}P \setminus T, \mathbb{R}^d)$  and  $g \in L^2(\eta^{-1}P \setminus T)$  satisfying  $\int_{\eta^{-1}P \setminus T} g dx = 0$ . Let  $(\mathbf{v}, \phi) \in H^1(\eta^{-1}P \setminus T, \mathbb{R}^d) \times L^2(\eta^{-1}P \setminus T)$  be the unique solution to the following Stokes system:*

$$\begin{cases} -\Delta \mathbf{v} + \nabla \phi = \mathbf{h} & \text{in } \eta^{-1}P \setminus T \\ \operatorname{div}(\mathbf{v}) = g & \text{in } \eta^{-1}P \setminus T \\ \int_{\eta^{-1}P \setminus T} \phi dx = 0 \\ \mathbf{v} = 0 & \text{on } \partial T \\ \mathbf{v} \text{ is } \eta^{-1}P\text{-periodic.} \end{cases} \quad (7.5.62)$$

There exists a constant  $C > 0$  independent of  $\eta, \mathbf{h}$  and  $g$  such that

$$\begin{aligned} \|\nabla \mathbf{v}\|_{L^2(\eta^{-1}P \setminus T, \mathbb{R}^{d \times d})} + \|\phi\|_{L^2(\eta^{-1}P \setminus T)} \\ \leq C(\eta^{-1} \|\mathbf{h} - \langle \mathbf{h} \rangle\|_{L^2(\eta^{-1}P \setminus T, \mathbb{R}^d)} + \eta^{-d} \|\langle \mathbf{h} \rangle\| + \|g\|_{L^2(\eta^{-1}P \setminus T)}). \end{aligned} \quad (7.5.63)$$

*Proof.* From [Lemma 2.2.4](#) in [\[14\]](#), for any  $\eta > 0$ , there exists a linear “Bogovskii's” operator  $B_\eta : L^2(P \setminus (\eta T)) \rightarrow H^1(P \setminus (\eta T), \mathbb{R}^d)$  satisfying for any  $\phi \in L^2(P \setminus (\eta Y))$  such that  $\int_{P \setminus (\eta T)} \phi dy = 0$ :



- (i)  $\operatorname{div}(B_\eta\phi) = \phi$ ,
- (ii)  $B_\eta\phi = 0$  on  $\partial(\eta T)$ ,
- (iii)  $B_\eta\phi$  is  $P$ -periodic,
- (iv)  $\|\nabla(B_\eta\phi)\|_{L^2(P\setminus(\eta T), \mathbb{R}^{d \times d})} \leq C\|\phi\|_{L^2(P\setminus(\eta T))}$  for a constant  $C$  independent of  $\eta$  and  $\phi$ .

For any  $\tilde{\phi} \in L^2(\eta^{-1}P\setminus T)$  such that  $\int_{\eta^{-1}P\setminus T} \tilde{\phi} dy = 0$ , we define

$$\tilde{B}_\eta(\tilde{\phi}) := \eta^{-1}[B_\eta(\tilde{\phi}(\eta^{-1}\cdot))(\eta\cdot)].$$

The operator  $\tilde{B}_\eta : L^2(\eta^{-1}P\setminus T) \rightarrow H^1(\eta^{-1}P\setminus T, \mathbb{R}^d)$  satisfies the following properties: for any  $\tilde{\phi} \in L^2(\eta^{-1}P\setminus T)$  such that  $\int_{\eta^{-1}P\setminus T} \tilde{\phi} dx = 0$ ,

- (i)  $\operatorname{div}(\tilde{B}_\eta\tilde{\phi}) = \tilde{\phi}$  in  $\eta^{-1}P\setminus T$ ,
- (ii)  $\tilde{B}_\eta\tilde{\phi} = 0$  on  $\partial T$ ,
- (iii)  $\tilde{B}_\eta\tilde{\phi}$  is  $\eta^{-1}P$ -periodic,
- (iv)  $\|\nabla(\tilde{B}_\eta\tilde{\phi})\|_{L^2(\eta^{-1}P\setminus T, \mathbb{R}^{d \times d})} \leq C\|\tilde{\phi}\|_{L^2(\eta^{-1}P\setminus T)}$  for a constant  $C$  independent of  $\eta$  and  $\tilde{\phi}$ .

The proof follows then classically along the lines of [corollary 7.11](#). Upon an integration by parts and by using [lemma 7.7](#), it is readily obtained with  $\mathbf{w} := \mathbf{v} - \tilde{B}_\eta g$ :

$$\begin{aligned} \|\nabla\mathbf{w}\|_{L^2(\eta^{-1}P\setminus T, \mathbb{R}^{d \times d})}^2 &= \int_{\eta^{-1}P\setminus T} \mathbf{h} \cdot \mathbf{w} dy \\ &= \int_{\eta^{-1}P\setminus T} (\mathbf{h} - \langle \mathbf{h} \rangle) \cdot (\mathbf{w} - \langle \mathbf{w} \rangle) dy + \int_{\eta^{-1}P\setminus T} \langle \mathbf{h} \rangle \cdot \langle \mathbf{w} \rangle dy \\ &\leq C(\|\mathbf{h} - \langle \mathbf{h} \rangle\|_{L^2(\eta^{-1}P\setminus T, \mathbb{R}^d)} \|\mathbf{w} - \langle \mathbf{w} \rangle\|_{L^2(\eta^{-1}P\setminus T, \mathbb{R}^d)} + \eta^{-d} |\langle \mathbf{h} \rangle| |\langle \mathbf{w} \rangle|) \\ &\leq C(\eta^{-1} \|\mathbf{h} - \langle \mathbf{h} \rangle\|_{L^2(\eta^{-1}P\setminus T, \mathbb{R}^d)} + \eta^{-d} |\langle \mathbf{h} \rangle|) \|\nabla\mathbf{w}\|_{L^2(\eta^{-1}P\setminus T, \mathbb{R}^{d \times d})}. \end{aligned} \tag{7.5.64}$$

for a constant  $C > 0$  independent of  $\eta$  and  $\mathbf{h}$ . This implies

$$\begin{aligned} \|\nabla\mathbf{v}\|_{L^2(\eta^{-1}P\setminus T)} &\leq \|\nabla\mathbf{w}\|_{L^2(\eta^{-1}P\setminus T, \mathbb{R}^{d \times d})} + \|\nabla(\tilde{b}_\eta g)\|_{L^2(\eta^{-1}P\setminus T, \mathbb{R}^{d \times d})} \\ &\leq C(\|\nabla\mathbf{w}\|_{L^2(\eta^{-1}P\setminus T, \mathbb{R}^{d \times d})} + \|g\|_{L^2(\eta^{-1}P\setminus T)}), \end{aligned}$$

whence the bound on  $\|\nabla\mathbf{v}\|_{L^2(\eta^{-1}P\setminus T, \mathbb{R}^{d \times d})}$  by using [\(7.5.64\)](#). The bound for the pressure is obtained by writing

$$\int_{\eta^{-1}P\setminus T} \phi^2 dx = \int_{\eta^{-1}P\setminus T} \phi \operatorname{div}(B_\eta\phi) dx = - \int_{\eta^{-1}P\setminus T} \nabla\phi \cdot B_\eta\phi dx = \int_{\eta^{-1}P\setminus T} (\nabla\mathbf{v} : \nabla(B_\eta\phi) - \mathbf{h} \cdot (B_\eta\phi)) dx,$$

from where [\(7.5.63\)](#) follows analogously. □

For any  $1 \leq j \leq d$ , let us introduce the unique solution  $(\Psi_j, \sigma_j)$  to the exterior Stokes problem

$$\left\{ \begin{array}{l} -\Delta\Psi_j + \nabla\sigma_j = 0 \text{ in } \mathbb{R}^d \setminus T \\ \operatorname{div}(\Psi_j) = 0 \text{ in } \mathbb{R}^d \setminus T \\ \Psi_j = 0 \text{ on } \partial T \\ \Psi_j \rightarrow \mathbf{e}_j \text{ at } \infty \\ \sigma_j \in L^2(\mathbb{R}^d \setminus T). \end{array} \right. \tag{7.5.65}$$

The boundary condition  $\Psi_j \rightarrow \mathbf{e}_j$  at infinity must be understood in the sense that  $\Psi_j - \mathbf{e}_j$  belongs to the Deny-Lions space  $\mathcal{D}^{1,2}(\mathbb{R}^d \setminus T, \mathbb{R}^d)$  (see [definition 7.6](#)). Similarly, the pressures  $\sigma_j$  are uniquely determined by the condition  $\sigma_j \in L^2(\mathbb{R}^d \setminus T)$  (see e.g. Lemma 1.1, chapter V. of [\[163\]](#)). We denote by  $\Psi^* := (\Psi_{ij}^*)_{1 \leq i, j \leq d}$  the matrix collecting the drag force components:

$$\Psi_{ij}^* := \int_{\mathbb{R}^d \setminus T} \nabla\Psi_i : \nabla\Psi_j dx = - \int_{\partial T} \mathbf{e}_j \cdot (\nabla\Psi_i - \sigma_i I) \cdot \mathbf{n} ds, \tag{7.5.66}$$

where the normal  $\mathbf{n}$  is pointing *inward*  $T$ . The asymptotics of  $\boldsymbol{\chi}^0$  and  $\boldsymbol{\chi}^{0*}$  have been obtained in of Theorem 3.1 in [13]. The following proposition extends propositions 7.14 and 7.27 to the Stokes system: it improves the result of [13] by providing asymptotics for the whole family of tensors  $(\boldsymbol{\chi}^k)_{k \in \mathbb{N}}$  and  $(\boldsymbol{\chi}^{k*})_{k \in \mathbb{N}}$ .

**Proposition 7.40.** *Assume  $d \geq 3$ . For any  $k \geq 0$  and  $1 \leq j \leq d$ , denote by  $(\tilde{\boldsymbol{\chi}}_j^{2k}, \tilde{\alpha}_j^{2k})$  and  $(\tilde{\boldsymbol{\chi}}_j^{2k+1}, \tilde{\alpha}_j^{2k+1})$  the rescaled tensors in  $\eta^{-1}P \setminus T$  defined as follows:*

$$\forall x \in \eta^{-1}P \setminus T, \begin{cases} \tilde{\boldsymbol{\chi}}_j^{2k}(x) := \eta^{(d-2)(k+1)} \boldsymbol{\chi}_j^{2k}(\eta x) \\ \tilde{\alpha}_j^{2k}(x) := \eta^{(d-2)(k+1)-1} \alpha_j^{2k}(\eta x) \end{cases}, \quad \begin{cases} \tilde{\boldsymbol{\chi}}_j^{2k+1}(x) := \eta^{(d-2)(k+1)} \boldsymbol{\chi}_j^{2k+1}(\eta x) \\ \tilde{\alpha}_j^{2k+1}(x) := \eta^{(d-2)(k+1)-1} \alpha_j^{2k+1}(\eta x). \end{cases}$$

Then:

1. there exists a constant  $C$  independent of  $\eta > 0$  such that

$$\forall \eta > 0, \|\nabla \tilde{\boldsymbol{\chi}}_j^{2k}\|_{L^2(\eta^{-1}P \setminus T, \mathbb{R}^{d \times d})} + \|\tilde{\alpha}_j^{2k}\|_{L^2(\eta^{-1}P \setminus T)} \leq C,$$

$$\forall \eta > 0, \|\nabla \tilde{\boldsymbol{\chi}}_j^{2k+1}\|_{L^2(\eta^{-1}P \setminus T, \mathbb{R}^{d \times d})} + \|\tilde{\alpha}_j^{2k+1}\|_{L^2(\eta^{-1}P \setminus T)} \leq C;$$

2. the following convergences hold as  $\eta \rightarrow 0$ :

$$(\tilde{\boldsymbol{\chi}}_i^{2k}, \tilde{\alpha}_i^{2k}) \rightharpoonup (c_{ij}^{2k} \boldsymbol{\Psi}_j, c_{ij}^{2k} \sigma_j) \quad \text{weakly in } H_{loc}^1(\mathbb{R}^d \setminus T, \mathbb{R}^d) \times L_{loc}^2(\mathbb{R}^d \setminus T), \quad (7.5.67)$$

$$(\tilde{\boldsymbol{\chi}}_i^{2k+1}, \tilde{\alpha}_i^{2k+1}) \rightharpoonup (0, 0) \quad \text{weakly in } H_{loc}^1(\mathbb{R}^d \setminus T, \mathbb{R}^d) \times L_{loc}^2(\mathbb{R}^d \setminus T), \quad (7.5.68)$$

$$\boldsymbol{\chi}^{2k*} \sim \frac{1}{\eta^{(d-2)(k+1)}} c^{2k}, \quad (7.5.69)$$

$$\boldsymbol{\chi}^{2k+1*} = o\left(\frac{1}{\eta^{(d-2)(k+1)}}\right), \quad (7.5.70)$$

where  $c_{ij}^{2k}$  denote the coefficients of the  $2k$ -th order matrix valued tensor  $c^{2k} := (c_{ij}^{2k})_{1 \leq i, j \leq d}$  given by

$$c^{2k} := (\Psi^*)^{-(k+1)} I^{2k} \quad \text{with } I^{2k} = \overbrace{I \otimes I \otimes \dots \otimes I}^{k \text{ times}}.$$

*Proof.* Following proposition 7.27, the result is proved by induction on  $k$ .

1. Case  $2k$  with  $k = 0$ . The tensor  $(\tilde{\boldsymbol{\chi}}_i^0, \tilde{\alpha}_i^0)$  satisfies

$$\begin{cases} -\Delta \tilde{\boldsymbol{\chi}}_i^0 + \nabla \tilde{\alpha}_i^0 = \eta^d \mathbf{e}_i \text{ in } \mathbb{R}^d \setminus T \\ \operatorname{div}(\tilde{\boldsymbol{\chi}}_i^0) = 0 \text{ in } \mathbb{R}^d \setminus T, \end{cases} \quad (7.5.71)$$

as well as the other boundary conditions of (7.5.62). Therefore lemma 7.11 yields

$$\|\nabla \tilde{\boldsymbol{\chi}}_i^0\|_{L^2(\eta^{-1}P \setminus T, \mathbb{R}^{d \times d})} + \|\tilde{\alpha}_i^0\|_{L^2(\eta^{-1}P \setminus T)} \leq C \eta^{-d} \eta^d |\langle \mathbf{e}_i \rangle| \leq C.$$

From (7.3.66), we also obtain that  $\langle \tilde{\boldsymbol{\chi}}_i^0 \rangle$  is bounded. Hence, up to extracting a subsequence, there exists a constant matrix  $c^0 := (c_{ij}^0)_{1 \leq i, j \leq d}$ , and fields  $(\hat{\boldsymbol{\Psi}}_i^0, \hat{\sigma}_i^0)_{1 \leq i \leq d}$  such that

$$\langle \tilde{\boldsymbol{\chi}}_i^0 \rangle \cdot \mathbf{e}_j \rightarrow c_{ij}^0,$$

$$(\tilde{\boldsymbol{\chi}}_i^0, \tilde{\alpha}_i^0) \rightharpoonup (\hat{\boldsymbol{\Psi}}_i^0, \hat{\sigma}_i^0) \text{ weakly in } H_{loc}^1(\eta^{-1}P \setminus T, \mathbb{R}^d) \times L_{loc}^2(\eta^{-1}P \setminus T).$$

Multiplying (7.5.71) by a compactly supported test function  $\Phi \in C_c^\infty(\mathbb{R}^d \setminus T)$  and integrating by parts yields

$$\int_{\eta^{-1}P \setminus T} (\nabla \tilde{\boldsymbol{\chi}}_i^0 : \nabla \Phi - \tilde{\alpha}_i^0 \operatorname{div}(\Phi)) dx = \int_{\eta^{-1}P \setminus T} \eta^d \Phi \cdot \mathbf{e}_i dx.$$

Passing to the limit as  $\eta \rightarrow 0$  then implies

$$\begin{cases} -\Delta \widehat{\Psi}_i^0 + \nabla \widehat{\sigma}_i^0 = 0 & \text{in } \mathbb{R}^d \setminus T \\ \operatorname{div}(\Psi_i^0) = 0 & \text{in } \mathbb{R}^d \setminus T \\ \Psi_i^0 = 0 & \text{on } \partial T. \end{cases}$$

Furthermore, by applying the point (7.3.68) of lemma 7.7 and by using the lower semi-continuity of the Lebesgue space norms, we infer  $(\widehat{\Phi}_i^0 - c_{ij}^0 e_j, \widehat{\sigma}_i^0) \in \mathcal{D}^{1,2}(\mathbb{R}^d \setminus T) \times L^2(\mathbb{R}^d \setminus T)$ . Hence, by linearity, it is necessary that  $(\widehat{\Phi}_i^0, \widehat{\sigma}_i^0) = (c_{ij}^0 \Psi_j, c_{ij}^0 \sigma_j)$  where  $(\Psi_j, \sigma_j)$  are the solution to the exterior problem (7.5.65). Integrating (7.5.65) by parts against the test function  $\Phi = e_j$  then yields

$$0 = \eta^d \int_{\eta^{-1}P \setminus T} \delta_{ij} dx + \int_{\partial T} e_j \cdot (\nabla \widetilde{\mathcal{X}}_i^0 - \widetilde{\alpha}_i^0 I) \cdot \mathbf{n} dx.$$

Passing to the limit as  $\eta \rightarrow 0$  by using the continuity of the drag force with respect to the weak convergence and (7.5.66) yields then

$$0 = \delta_{ij} + \int_{\partial T} e_j \cdot (\nabla \widehat{\Phi}_i^0 - \widehat{\sigma}_i^0) \cdot \mathbf{n} dx = \delta_{ij} - c_{ip}^0 \Psi_{pj}^*.$$

This implies  $c^0 = (\Psi^*)^{-1}$  as claimed, and the convergence of the whole sequence by uniqueness of the limit. The asymptotic for  $\mathcal{X}^{0*}$  as  $\eta \rightarrow 0$  is obtained by a simple change of variable  $y = \eta x$ :

$$\mathcal{X}_{ij}^{0*} = e_i \cdot \int_{P \setminus (\eta T)} \mathcal{X}_j^0 dy = \eta^{2-d} \eta^d e_i \cdot \int_{\eta^{-1}P \setminus T} \widetilde{\mathcal{X}}_j^0 dy \sim \eta^{2-d} \langle \widetilde{\mathcal{X}}_j^0 \rangle \cdot e_i \sim \eta^{2-d} c_{ji}^0.$$

2. *Case  $2k+1$  with  $k=0$ .* The tensor  $(\widetilde{\mathcal{X}}_i^1, \widetilde{\alpha}_i^1)$  satisfies

$$\begin{cases} -\Delta \widetilde{\mathcal{X}}_i^1 + \nabla \widetilde{\alpha}_i^1 = \eta(2\partial_l \widetilde{\mathcal{X}}_i^0 - \widetilde{\alpha}_j^0 e_l) \otimes e_l & \text{in } \eta^{-1}P \setminus T \\ \operatorname{div}(\widetilde{\mathcal{X}}_i^1) = -\eta(\widetilde{\mathcal{X}}_j^0 - \langle \widetilde{\mathcal{X}}_j^0 \rangle) \cdot e_l \otimes e_l & \text{in } \eta^{-1}P \setminus T. \end{cases} \quad (7.5.72)$$

Applying lemma 7.11 and remarking that  $\langle 2\partial_l \widetilde{\mathcal{X}}_i^0 - \widetilde{\alpha}_j^0 e_l \rangle = 0$ , we obtain

$$\|\nabla \widetilde{\mathcal{X}}_i^1\|_{L^2(\eta^{-1}P \setminus T, \mathbb{R}^{d \times d})} + \|\widetilde{\alpha}_i^1\|_{L^2(\eta^{-1}P \setminus T)} \leq C.$$

Integrating (7.5.72) by parts against a compactly supported test function  $\Phi \in C^c(\mathbb{R}^d \setminus T)$  and passing to the limit as  $\eta \rightarrow 0$ , we obtain with similar arguments, the existence of a matrix valued tensor  $c^1 := (c_{ij}^1)_{1 \leq i, j \leq d}$  (of order 1) such that, up to the extraction of a subsequence:

$$\langle \widetilde{\mathcal{X}}_i^1 \rangle \cdot e_j \rightarrow c_{ij}^1,$$

$$(\widetilde{\mathcal{X}}_i^1, \widetilde{\alpha}_i^1) \rightharpoonup (c_{ij}^1 \Psi_j, c_{ij}^1 \sigma_j) \text{ weakly in } H_{loc}^1(\eta^{-1}P \setminus T, \mathbb{R}^d) \times L_{loc}^2(\eta^{-1}P \setminus T).$$

Integrating (7.5.72) by parts against the test function  $e_j$  and passing to the limit as  $\eta \rightarrow 0$  yields in this context

$$0 = c_{ij}^1 \Psi_{pj}^*$$

whence  $c^1 = 0$ .

3. *General case.* Assuming the result holds till rank  $k$ , we write the differential equations satisfied by the rescaled tensors:

$$\begin{cases} -\Delta \widetilde{\mathcal{X}}_i^{2k+2} + \nabla \widetilde{\alpha}_i^{2k+2} = \eta^{d-1}(\partial_l \widetilde{\mathcal{X}}_i^{2k+1} - \widetilde{\alpha}_i^{2k+1} e_l) \otimes e_l + \eta^d \widetilde{\mathcal{X}}_i^{2k} \otimes I & \text{in } \eta^{-1}P \setminus T \\ \operatorname{div}(\widetilde{\mathcal{X}}_i^{2k+2}) = -\eta^{d-1}(\widetilde{\mathcal{X}}_i^{2k+1} - \langle \widetilde{\mathcal{X}}_i^{2k+1} \rangle) \cdot e_l \otimes e_l & \text{in } \eta^{-1}P \setminus T. \end{cases} \quad (7.5.73)$$

$$\begin{cases} -\Delta \widetilde{\mathcal{X}}_i^{2k+3} + \nabla \widetilde{\alpha}_i^{2k+3} = \eta(\partial_l \widetilde{\mathcal{X}}_i^{2k+2} - \widetilde{\alpha}_i^{2k+2} e_l) \otimes e_l + \eta^d \widetilde{\mathcal{X}}_i^{2k+1} \otimes I & \text{in } \eta^{-1}P \setminus T \\ \operatorname{div}(\widetilde{\mathcal{X}}_i^{2k+3}) = -\eta(\widetilde{\mathcal{X}}_i^{2k+2} - \langle \widetilde{\mathcal{X}}_i^{2k+2} \rangle) \cdot e_l \otimes e_l & \text{in } \eta^{-1}P \setminus T. \end{cases} \quad (7.5.74)$$

Using lemma 7.7, lemma 7.11 and the point (1) of the proposition at rank  $k$ , we readily obtain

$$\begin{aligned} \|\nabla \tilde{\mathcal{X}}_i^{2k+2}\|_{L^2(\eta^{-1}P \setminus T, \mathbb{R}^{d \times d})} + \|\tilde{\alpha}_i^{2k+2}\|_{L^2(\eta^{-1}P \setminus T)} &\leq C, \\ \|\nabla \tilde{\mathcal{X}}_i^{2k+3}\|_{L^2(\eta^{-1}P \setminus T, \mathbb{R}^{d \times d})} + \|\tilde{\alpha}_i^{2k+3}\|_{L^2(\eta^{-1}P \setminus T)} &\leq C. \end{aligned}$$

Repeating the above arguments, we obtain, up to the extraction of a subsequence, the existence of matrix valued tensors  $c^{2k+2}$  and  $c^{2k+3}$  such that

$$\begin{aligned} \langle \tilde{\mathcal{X}}_i^{2k+2} \rangle \cdot \mathbf{e}_j &\rightarrow c_{ij}^{2k+2}, \text{ and } \langle \tilde{\mathcal{X}}_i^{2k+3} \rangle \cdot \mathbf{e}_j \rightarrow c_{ij}^{2k+3}, \\ (\tilde{\mathcal{X}}_i^{2k+2}, \tilde{\alpha}_i^{2k+2}) &\rightharpoonup (c_{ij}^{2k+2} \Psi_j, c_{ij}^{2k+2} \sigma_j) \text{ weakly in } H_{loc}^1(\mathbb{R}^d \setminus T, \mathbb{R}^d) \times L_{loc}^2(\mathbb{R}^d \setminus T), \\ (\tilde{\mathcal{X}}_i^{2k+3}, \tilde{\alpha}_i^{2k+3}) &\rightharpoonup (c_{ij}^{2k+3} \Psi_j, c_{ij}^{2k+3} \sigma_j) \text{ weakly in } H_{loc}^1(\mathbb{R}^d \setminus T, \mathbb{R}^d) \times L_{loc}^2(\mathbb{R}^d \setminus T). \end{aligned}$$

The last step consists in integrating (7.5.73) and (7.5.74) by part against the test function  $\mathbf{e}_j$  and to pass to the limit as  $\eta \rightarrow 0$  to identify  $c_{ij}^{2k+2}$  and  $c_{ij}^{2k+3}$ . Performing this computation as above, we obtain

$$\begin{aligned} 0 &= c_{ij}^{2k} \otimes I - c_{ip}^{2k+2} \Psi_{pj}^*, \\ 0 &= c_{ij}^{2k+1} \otimes I - c_{ip}^{2k+3} \Psi_{pj}^* \end{aligned}$$

from where we infer  $c^{2k+2} = c^{2k}(\Psi^*)^{-1} \otimes I$ ,  $c^{2k+3} = c^{2k+1}(\Psi^*)^{-1} \otimes I$ , hence the result (recall  $c^1 = 0$  from the point (2) of the proof).  $\square$

Using the identity (7.5.30), we obtain the asymptotics for the coefficients  $M^k$  of the infinite order homogenized equation (7.5.29).

**Corollary 7.12.** *Assume  $d \geq 3$ . The following convergences hold for the matrix valued tensors  $M^k$  as  $\eta \rightarrow 0$ :*

$$M^0 \sim \eta^{d-2} \Psi^*, \quad (7.5.75)$$

$$M^1 = o(\eta^{d-2}), \quad (7.5.76)$$

$$M^2 \rightarrow -I, \quad (7.5.77)$$

$$\forall k \geq 1, M^{2k} = o\left(\frac{1}{\eta^{(d-2)(k-1)}}\right), \quad (7.5.78)$$

$$\forall k \geq 1, M^{2k+1} = o\left(\frac{1}{\eta^{(d-2)(k-1)}}\right). \quad (7.5.79)$$

*Proof.* The proof is identical to that of corollary 7.6 by substituting  $A^2$  with  $I$ .  $\square$

### 7.5.5 Simplifications for the tensors $\mathcal{X}^{k*}$ and $M^k$ under symmetries

This final subsection generalizes the results of sections 7.3.6 and 7.4.4 to the Stokes system: we examine how the symmetries of the obstacle  $\eta T$  with respect to the cell axes reflect into the coefficients of the matrix valued tensors  $\mathcal{X}^{k*}$  and  $M^k$ . Our final result is stated in corollary 7.13, it is based on the following elementary lemma:

**Lemma 7.12.** *Let  $S \in \mathbb{R}^{d \times d}$  an orthogonal symmetry, i.e.  $S = S^T$  and  $SS = I$ . The following identities hold for any smooth vector field  $\mathcal{X}$  and scalar field  $\alpha$ :*

$$-\Delta(S\mathcal{X} \circ S) + \nabla(\alpha \circ S) = S(-\Delta\mathcal{X} + \nabla\alpha) \circ S, \quad (7.5.80)$$

$$\operatorname{div}(S\mathcal{X} \circ S) = \operatorname{div}(\mathcal{X}) \circ S, \quad (7.5.81)$$

$$\partial_i(S\mathcal{X} \circ S) = S_{ij}S(\partial_j\mathcal{X}) \circ S. \quad (7.5.82)$$

*Proof.* The first two identities are obtained by writing

$$\begin{aligned} -\Delta(S\mathcal{X} \circ S) + \nabla(\alpha \circ S) &= -S\partial_{ij}\mathcal{X} \circ SS_{il}S_{jl} + S(\nabla\alpha) \circ S \\ &= -S(\Delta\mathcal{X} + \nabla\alpha) \circ S, \end{aligned}$$

$$\operatorname{div}(S\mathcal{X} \circ S) = \operatorname{Tr}(\nabla(S\mathcal{X} \circ S)) = \operatorname{Tr}(S(\nabla\mathcal{X}) \circ SS) = \operatorname{Tr}((\nabla\mathcal{X}) \circ S) = \operatorname{div}(\mathcal{X}) \circ S.$$

Identity (7.5.82) is an elementary consequence of the chain rule.  $\square$

**Proposition 7.41.** *If the cell  $Y = P \setminus (\eta T)$  is invariant with respect to a symmetry  $S$ , i.e.  $S(Y) = Y$ , then the following identity holds for the tensors  $(\mathcal{X}_l^k, \alpha_l^k)$  (defined in (7.5.6) to (7.5.8)):*

$$S\mathcal{X}_{i_1 \dots i_k, l}^k \circ S = S_{i_1 j_1} \dots S_{i_k j_k} S_{lm} \mathcal{X}_{j_1 \dots j_k, m}^k, \quad (7.5.83)$$

$$\alpha_{i_1 \dots i_k, l}^k \circ S = S_{i_1 j_1} \dots S_{i_k j_k} S_{lm} \alpha_{j_1 \dots j_k, m}^k. \quad (7.5.84)$$

As a consequence, the following identities hold for the constant matrix valued tensors  $\mathcal{X}^{k*}$  and  $M^k$ :

$$\mathcal{X}_{i_1 \dots i_k, lm}^{k*} = S_{i_1 j_1} \dots S_{i_k j_k} S_{lp} S_{mq} \mathcal{X}_{j_1 \dots j_k, pq}^{k*} \quad (7.5.85)$$

$$M_{i_1 \dots i_k, lm}^k = S_{i_1 j_1} \dots S_{i_k j_k} S_{lp} S_{mq} M_{j_1 \dots j_k, pq}^k. \quad (7.5.86)$$

*Proof.* We follow the proof of proposition 7.28, that is we prove (7.5.83) and (7.5.84) by induction. Using section 7.5.5, we easily obtain

$$\begin{cases} -\Delta_{yy}(S\mathcal{X}_l^0 \circ S) + \nabla_y(\alpha_l^0 \circ S) = S\mathbf{e}_l \circ S = S\mathbf{e}_l = S_{mj}\mathbf{e}_m, \\ \operatorname{div}(S\mathcal{X}_l^0 \circ S) = 0. \end{cases}$$

Since the cell is symmetric with respect to  $S$ ,  $(S\mathcal{X}_l^0 \circ S, \alpha_l^0 \circ S)$  satisfies the same boundary conditions (7.5.9) than  $S_{mj}(\mathcal{X}_m^0, \alpha_m^0)$ , therefore these are equal and we infer (7.5.83) and (7.5.84) at rank  $k = 0$ . We then write, for a given  $1 \leq i_1 \leq d$ :

$$\begin{cases} -\Delta_{yy}(S\mathcal{X}_{i_1, l}^1 \circ S) + \nabla_y(\alpha_{i_1, l}^1 \circ S) = S(2\partial_{i_1} \mathcal{X}_l^0 - \alpha_l^0 \mathbf{e}_{i_1}) \circ S \\ \quad = S_{i_1 j_1} (2\partial_{j_1} (S\mathcal{X}_l^0 \circ S) - \alpha_l^0 \circ S \mathbf{e}_{j_1}) \\ \quad = S_{i_1 j_1} S_{lm} (2\partial_{j_1} \mathcal{X}_m^0 - \alpha_m^0 \mathbf{e}_{j_1}), \\ \operatorname{div}_y(S\mathcal{X}_{i_1, l}^1 \circ S) = -(\mathcal{X}_l^0 \circ S - \langle \mathcal{X}_l^0 \rangle) \cdot \mathbf{e}_{i_1} \\ \quad = -S_{lm} S(\mathcal{X}_m^0 - \langle \mathcal{X}_m^0 \rangle) \cdot \mathbf{e}_{i_1} \\ \quad = -S_{i_1 j_1} S_{lm} (\mathcal{X}_m^0 - \langle \mathcal{X}_m^0 \rangle) \cdot \mathbf{e}_{j_1}, \end{cases}$$

where we used  $\langle \mathcal{X}_l^0 \rangle = \langle \mathcal{X}_l^0 \circ S \rangle$ . This implies similarly (7.5.83) and (7.5.84) at rank  $k = 1$ . Assuming now the result holds till rank  $k + 1$  with  $k \geq 0$ , we prove with the same arguments that it holds at rank  $k + 2$ :

$$\begin{cases} -\Delta_{yy}(S\mathcal{X}_{i_1 \dots i_{k+2}, l}^{k+2} \circ S) + \nabla_y(\alpha_{i_1 \dots i_{k+2}, l}^{k+2} \circ S) \\ \quad = S(2\partial_{i_{k+2}} \mathcal{X}_{i_1 \dots i_{k+1}, l}^{k+1} - \alpha_{i_1 \dots i_{k+1}, l}^{k+1} \mathbf{e}_{i_{k+2}}) \circ S + S\mathcal{X}_{i_1 \dots i_k, l}^k \circ S \delta_{i_{k+1} i_{k+2}} \\ \quad = S_{i_{k+2} j_{k+2}} (2\partial_{j_{k+2}} (S\mathcal{X}_{i_1 \dots i_{k+1}, l}^{k+1} \circ S) - \alpha_{i_1 \dots i_{k+1}, l}^{k+1} \circ S \mathbf{e}_{j_{k+2}}) \\ \quad \quad + S_{i_{k+1} j_{k+1}} S_{i_{k+2} j_{k+2}} \delta_{j_{k+1} j_{k+2}} S\mathcal{X}_{i_1 \dots i_k, l}^k \circ S \\ \quad = S_{i_1 j_1} \dots S_{i_{k+2} j_{k+2}} S_{lm} [(2\partial_{j_{k+2}} \mathcal{X}_{j_1 \dots j_{k+1}, m}^{k+1} - \alpha_{j_1 \dots j_{k+1}, m}^{k+1} \mathbf{e}_{j_{k+2}}) + \delta_{j_{k+1} j_{k+2}} \mathcal{X}_{j_1 \dots j_k, m}^k] \\ \operatorname{div}_y(S\mathcal{X}_{i_1 \dots i_{k+2}, l}^{k+2} \circ S) = -(\mathcal{X}_{i_1 \dots i_{k+1}, l}^{k+1} \circ S - \langle \mathcal{X}_{i_1 \dots i_{k+1}, l}^{k+1} \rangle) \dots \mathbf{e}_{i_{k+2}} \\ \quad = -S_{i_1 j_1} \dots S_{i_{k+1} j_{k+1}} S_{lm} S(\mathcal{X}_{j_1 \dots j_{k+1}, m}^{k+1} - \langle \mathcal{X}_{j_1 \dots j_{k+1}, m}^{k+1} \rangle) \cdot \mathbf{e}_{i_{k+2}} \\ \quad = -S_{i_1 j_1} \dots S_{i_{k+2} j_{k+2}} S_{lm} (\mathcal{X}_{j_1 \dots j_{k+1}, m}^{k+1} - \langle \mathcal{X}_{j_1 \dots j_{k+1}, m}^{k+1} \rangle) \cdot \mathbf{e}_{j_{k+2}}. \end{cases}$$

The identities (7.5.85) and (7.5.86) follow as in the proof of proposition 7.28.  $\square$

**Corollary 7.13.** 1. *If the cell  $Y$  is symmetric with respect to all cell axes  $(\mathbf{e}_l)_{1 \leq l \leq d}$ , then*

$$\mathcal{X}_{i_1 \dots i_k, pq}^{k*} = 0 \text{ and } M_{i_1 \dots i_k, pq}^k = 0$$

*whenever any given integer  $1 \leq l \leq d$  occurs an odd number of times in the indices  $i_1 \dots i_k, p, q$ .*

*In particular, this implies  $\mathcal{X}^{2k+1*} = 0$  and  $M^{2k+1} = 0$ .*

2. *If the cell  $Y$  is symmetric with respect to all diagonal axes orthogonal to  $(\mathbf{e}_l - \mathbf{e}_m)$ , i.e.  $S^{l,m}(Y) = Y$  for any  $1 \leq l < m \leq d$ , then for any permutation  $\sigma \in \mathfrak{S}_d$ ,*

$$\mathcal{X}_{\sigma(i_1) \dots \sigma(i_k), \sigma(p)\sigma(q)}^{k*} = \mathcal{X}_{i_1 \dots i_k, pq}^{k*}$$

$$M_{\sigma(i_1) \dots \sigma(i_k), \sigma(p)\sigma(q)}^k = M_{i_1 \dots i_k, pq}^k$$

*Proof.* The proof is identical to that of [corollary 7.7](#).  $\square$

### 7.5.6 Appendix: extension to multicomponent fluid domains

Let us outline in this appendix how the previous methodologies could be extended to the case where the unit cell  $Y$  (a subset of the torus, i.e. opposite matching boundaries are identified) has  $m$  connected components  $(Y_l)_{1 \leq l \leq m}$  instead of one as assumed in [\(H4\)](#). Related models could be of interest for topology and shape optimization of multicomponent fluid systems, such as heat exchangers [[255](#), [303](#)].

In such a case the ansatz for the pressure  $p_\varepsilon(x)$  of [\(7.5.2\)](#) must be modified as follows:

$$\mathbf{u}_\varepsilon(x) = \sum_{i=0}^{+\infty} \varepsilon^{i+2} \mathbf{u}_i(x, x/\varepsilon), \quad p_\varepsilon(x) = \sum_{i=0}^{+\infty} \varepsilon^i \left( \sum_{l=1}^m p_i^{*,l}(x) \mathbf{1}_{Y_l}(x/\varepsilon) + \varepsilon p_i(x, x/\varepsilon) \right) \quad (7.5.87)$$

where  $m$  homogenized pressures  $p_i^{*,l}(x)$  are now involved, with  $1 \leq l \leq m$  and  $i \in \mathbb{N}$ . The uniqueness of the decomposition is obtained by imposing the average of  $p_i$  to be zero in every connected component:

$$\forall 1 \leq l \leq m, \forall x \in D, \int_{Y_l} p_i(x, y) dy = 0, \quad \int_D p_i^{*,l}(x) dx = 0.$$

This ansatz yields  $m$  homogenized pressures  $(p_\varepsilon^{*,l})_{1 \leq l \leq m}$  and velocity fields  $(\mathbf{u}_\varepsilon^{*,l})_{1 \leq l \leq m}$  for each of the corresponding connected fluid components:

$$\forall 1 \leq l \leq m, \mathbf{u}_\varepsilon^{*,l}(x) := \sum_{i=0}^{+\infty} \varepsilon^{i+2} \int_{Y_l} \mathbf{u}_i(x, y) dy, \quad p_\varepsilon^{*,l}(x) := \sum_{i=0}^{+\infty} \varepsilon^i p_i^{*,l}(x).$$

The analysis of [section 7.5.1](#) can be easily adapted to this new context. A Fredholm alternative, analogous to [\(7.5.13\)](#), yields an elliptic second order equation determining the value of the  $m$  pressures  $p_i^{*,l}$ :

$$\operatorname{div}_x(\mathcal{X}^{0*,l} \nabla_x(\mathbf{f}_i - \nabla_x p_i^{*,l})) = - \sum_{k=1}^i \operatorname{div}(\mathcal{X}^{k*,l} \cdot \nabla^k(\mathbf{f}_{i-k} - \nabla_x p_{i-k}^{*,l})), \quad \forall i \geq 0, \forall 1 \leq l \leq m,$$

where  $\mathcal{X}_{ij}^{0*,l} := \int_{Y_l} \mathcal{X}_j^0 \cdot \mathbf{e}_i dy$ . Then the velocity and pressure  $\mathbf{u}_i(x, y)$  and  $p_i(x, y)$  of the ansatz [\(7.5.87\)](#) read

$$\mathbf{u}_i(x, y) = \sum_{k=0}^i \mathcal{X}^k(y) \cdot \nabla^k(\mathbf{f}_{i-k}(x) - \mathbf{1}_{Y_l}(y) \nabla p_{i-k}^{*,l}(x))$$

$$p_i(x, y) = \sum_{k=0}^i \alpha^k(y) \cdot \nabla^k(\mathbf{f}_{i-k}(x) - \mathbf{1}_{Y_l}(y) \nabla p_{i-k}^{*,l}(x)),$$

where the summation over the repeated index  $l$  is assumed. These conclusions are not surprising: we have just verified that the homogenization theory for a fluid system with  $m$  components is that of  $m$  homogenized systems of single component (obtained e.g. by replacing  $m-1$  fluid components with solid).

## REFERENCES

- [1] N. AAGE, E. ANDREASSEN, B. S. LAZAROV, AND O. SIGMUND, *Giga-voxel computational morphogenesis for structural design*, *Nature*, 550 (2017), p. 84.
- [2] T. ABATZOGLOU, *The metric projection on  $C^2$  manifolds in Banach spaces*, *Journal of Approximation Theory*, 26 (1979), pp. 204–211.
- [3] A. ABDULLE AND T. POUCHON, *Effective models for the multidimensional wave equation in heterogeneous media over long time and numerical homogenization*, *Mathematical Models and Methods in Applied Sciences*, 26 (2016), pp. 2651–2684.
- [4] A. ABDULLE AND T. POUCHON, *Effective models and numerical homogenization for wave propagation in heterogeneous media on arbitrary timescales*, arXiv preprint arXiv:1905.09062, (2019).
- [5] P.-A. ABSIL, R. MAHONY, AND R. SEPULCHRE, *Optimization algorithms on matrix manifolds*, Princeton University Press, 2009.
- [6] P.-A. ABSIL AND J. MALICK, *Projection-like retractions on matrix manifolds*, *SIAM Journal on Optimization*, 22 (2012), pp. 135–158.
- [7] E. ACERBI, V. CHIADOPIAT, G. DAL MASO, AND D. PERCIVALE, *An extension theorem from connected sets, and homogenization in general periodic domains*, *Nonlinear Analysis: Theory, Methods & Applications*, 18 (1992), pp. 481–496.
- [8] M. F. ADAMS, H. H. BAYRAKTAR, T. M. KEAVENY, AND P. PAPADOPOULOS, *Ultrascale implicit finite element analyses in solid mechanics with over a half a billion degrees of freedom*, in *Proceedings of the 2004 ACM/IEEE conference on Supercomputing*, IEEE Computer Society, 2004, p. 34.
- [9] S.-H. AHN AND S. CHO, *Level set-based topological shape optimization of heat conduction problems considering design-dependent convection boundary*, *Numerical Heat Transfer, Part B: Fundamentals*, 58 (2010), pp. 304–322.
- [10] J. ALEXANDERSEN, O. SIGMUND, AND N. AAGE, *Large scale three-dimensional topology optimization of heat sinks cooled by natural convection*, *International Journal of Heat and Mass Transfer*, 100 (2016), pp. 876–891.
- [11] G. ALLAIRE, *Homogénéisation des équations de Stokes et de Navier-Stokes*, PhD thesis, Université Paris 6, 1989.
- [12] ———, *Homogenization of the Stokes flow in a connected porous medium*, *Asymptotic Analysis*, 2 (1989), pp. 203–222.
- [13] G. ALLAIRE, *Continuity of the Darcy’s law in the low-volume fraction limit*, *Annali della Scuola Normale Superiore di Pisa. Classe di Scienze. Serie IV*, 18 (1991), pp. 475–499.
- [14] G. ALLAIRE, *Homogenization of the Navier-Stokes equations in open sets perforated with tiny holes I. Abstract framework, a volume distribution of holes*, *Archive for Rational Mechanics and Analysis*, 113 (1991), pp. 209–259.
- [15] G. ALLAIRE, *Homogenization of the Navier-Stokes equations with a slip boundary condition*, *Communications on pure and applied mathematics*, 44 (1991), pp. 605–641.
- [16] G. ALLAIRE, *Homogenization and two-scale convergence*, *SIAM Journal on Mathematical Analysis*, 23 (1992), pp. 1482–1518.
- [17] G. ALLAIRE, *Conception optimale de structures*, vol. 58 of *Mathématiques & Applications* (Berlin), Springer-Verlag, Berlin, 2007.
- [18] G. ALLAIRE, *Shape optimization by the homogenization method*, vol. 146, Springer Science & Business Media, 2012.
- [19] G. ALLAIRE AND M. AMAR, *Boundary layer tails in periodic homogenization*, *ESAIM: Control, Optimisation and Calculus of Variations*, 4 (1999), pp. 209–243.

- [20] G. ALLAIRE, M. BRIANE, AND M. VANNINATHAN, *A comparison between two-scale asymptotic expansions and Bloch wave expansions for the homogenization of periodic structures*, SEMA journal, 73 (2016), pp. 237–259.
- [21] G. ALLAIRE, E. CANCÈS, AND J.-L. VIÉ, *Second-order shape derivatives along normal trajectories, governed by Hamilton-Jacobi equations*, Structural and Multidisciplinary Optimization, 54 (2016), pp. 1245–1266.
- [22] G. ALLAIRE AND C. CONCA, *Bloch wave homogenization and spectral asymptotic analysis*, Journal de mathématiques pures et appliquées, 77 (1998), pp. 153–208.
- [23] G. ALLAIRE, C. DAPOGNY, G. DELGADO, AND G. MICHAILIDIS, *Multi-phase structural optimization via a level set method*, ESAIM: Control, Optimisation and Calculus of Variations, 20 (2014), pp. 576–611.
- [24] G. ALLAIRE, C. DAPOGNY, AND P. FREY, *A mesh evolution algorithm based on the level set method for geometry and topology optimization*, Structural and Multidisciplinary Optimization, 48 (2013), pp. 711–715.
- [25] ———, *Shape optimization with a level set based mesh evolution method*, Computer Methods in Applied Mechanics and Engineering, 282 (2014), pp. 22–53.
- [26] G. ALLAIRE, F. DE GOURNAY, F. JOUVE, AND A.-M. TOADER, *Structural optimization using topological and shape sensitivity via a level set method*, Control and cybernetics, 34 (2005), p. 59.
- [27] G. ALLAIRE, P. GEOFFROY-DONDERS, AND O. PANTZ, *Topology optimization of modulated and oriented periodic microstructures by the homogenization method*, Computers & Mathematics with Applications, (2018).
- [28] G. ALLAIRE AND L. JAKABCIN, *Taking into account thermal residual stresses in topology optimization of structures built by additive manufacturing*, Mathematical Models and Methods in Applied Sciences, 28 (2018), pp. 2313–2366.
- [29] G. ALLAIRE, F. JOUVE, AND G. MICHAILIDIS, *Molding direction constraints in structural optimization via a level-set method*, in Variational Analysis and Aerospace Engineering, Springer, 2016, pp. 1–39.
- [30] G. ALLAIRE, F. JOUVE, AND G. MICHAILIDIS, *Thickness control in structural optimization via a level set method*, Structural and Multidisciplinary Optimization, 53 (2016), pp. 1349–1382.
- [31] G. ALLAIRE, F. JOUVE, AND A. TOADER, *A level-set method for shape optimization*, Comptes Rendus Mathématique, 334 (2002), pp. 1125–1130.
- [32] G. ALLAIRE, F. JOUVE, AND A.-M. TOADER, *Structural optimization using sensitivity analysis and a level-set method*, Journal of Computational Physics, 194 (2004), pp. 363–393.
- [33] G. ALLAIRE, A. LAMACZ, AND J. RAUCH, *Crime Pays; Homogenized Wave Equations for Long Times*, arXiv preprint arXiv:1803.09455, (2018).
- [34] G. ALLAIRE AND O. PANTZ, *Structural optimization with FreeFem++*, Structural and Multidisciplinary Optimization, 32 (2006), pp. 173–181.
- [35] G. ALLAIRE AND T. YAMADA, *Optimization of dispersive coefficients in the homogenization of the wave equation in periodic structures*, Numerische Mathematik, 140 (2018), pp. 265–326.
- [36] L. AMBROSIO, *Geometric evolution problems, distance function and viscosity solutions*, in Calculus of variations and partial differential equations, Springer, 2000, pp. 5–93.
- [37] N. AMENTA, S. CHOI, AND R. K. KOLLURI, *The power crust, unions of balls, and the medial axis transform*, Computational Geometry, 19 (2001), pp. 127–153.
- [38] P. R. AMESTOY, I. S. DUFF, J. KOSTER, AND J.-Y. L’EXCELLENT, *A Fully Asynchronous Multifrontal Solver Using Distributed Dynamic Scheduling*, SIAM Journal on Matrix Analysis and Applications, 23 (2001), pp. 15–41.



- [39] O. AMIR, N. AAGE, AND B. S. LAZAROV, *On multigrid-CG for efficient topology optimization*, Structural and Multidisciplinary Optimization, 49 (2014), pp. 815–829.
- [40] S. AMSTUTZ AND H. ANDRÄ, *A new algorithm for topology optimization using a level-set method*, Journal of Computational Physics, 216 (2006), pp. 573–588.
- [41] M. ANDERSEN, J. DAHL, AND L. VANDENBERGHE, *CVXOPT: A Python package for convex optimization*, Available at <http://cvxopt.org/>, (2012).
- [42] C. S. ANDREASEN AND O. SIGMUND, *Topology optimization of fluid–structure-interaction problems in poroelasticity*, Computer Methods in Applied Mechanics and Engineering, 258 (2013), pp. 55–62.
- [43] E. ANDREASSEN, A. CLAUSEN, M. SCHEVENELS, B. S. LAZAROV, AND O. SIGMUND, *Efficient topology optimization in MATLAB using 88 lines of code*, Structural and Multidisciplinary Optimization, 43 (2011), pp. 1–16.
- [44] S. ARGUILLÈRE, E. TRÉLAT, A. TROUVÉ, AND L. YOUNES, *Shape deformation analysis from the optimal control viewpoint*, Journal de Mathématiques Pures et Appliquées. Neuvième Série, 104 (2015), pp. 139–178.
- [45] D. ATTALI, J.-D. BOISSONNAT, AND H. EDELSBRUNNER, *Stability and computation of medial axes—a state-of-the-art report*, in Mathematical foundations of scientific visualization, computer graphics, and massive data exploration, Springer, 2009, pp. 109–125.
- [46] J.-L. AURIAULT, *On the domain of validity of brinkman’s equation*, Transport in Porous Media, 79 (2009), pp. 215–223.
- [47] J.-L. AURIAULT, C. GEINDREAU, AND C. BOUTIN, *Darcy’s law, brinkman’s law and poor separation of scales*, Poromechanics III: Biot Centennial (1905-2005) - Proceedings of the 3rd Biot Conference on Poromechanics, (2005), pp. 553–558.
- [48] H. AZEGAMI, S. K. M. SHIMODA, AND E. KATAMINE, *Irregularity of shape optimization problems and an improvement technique*, WIT Transactions on The Built Environment, 31 (1997).
- [49] H. AZEGAMI AND Z. C. WU, *Domain optimization analysis in linear elastic problems: approach using traction method*, JSME international journal. Ser. A, Mechanics and material engineering, 39 (1996), pp. 272–278.
- [50] P. AZÉRAD, *Equations de Navier-Stokes en bassin peu profond*, PhD thesis, Université de Neuchâtel, 1995.
- [51] P. AZÉRAD AND J. POUSIN, *Inégalité de Poincaré courbe pour le traitement variationnel de l’équation de transport*, Comptes rendus de l’Académie des sciences. Série 1, Mathématique, 322 (1996), pp. 721–727.
- [52] C. BACUTA, *A unified approach for Uzawa algorithms*, SIAM Journal on Numerical Analysis, 44 (2006), pp. 2633–2649.
- [53] N. BAKHVALOV AND G. PANASENKO, *Homogenisation: averaging processes in periodic media*, vol. 36 of Mathematics and its Applications (Soviet Series), Kluwer Academic Publishers Group, Dordrecht, 1989.
- [54] S. BALAY, S. ABHYANKAR, M. F. ADAMS, J. BROWN, P. BRUNE, K. BUSCHELMAN, L. DALCIN, A. DENER, V. EIJKHOUT, W. D. GROPP, D. KARPEYEV, D. KAUSHIK, M. G. KNEPLEY, D. A. MAY, L. C. MCINNES, R. T. MILLS, T. MUNSON, K. RUPP, P. SANAN, B. F. SMITH, S. ZAMPINI, H. ZHANG, AND H. ZHANG, *PETSc Web page*. <https://www.mcs.anl.gov/petsc>, 2019.
- [55] ———, *PETSc users manual*, Tech. Rep. ANL-95/11 - Revision 3.11, Argonne National Laboratory, 2019.
- [56] S. BALAY, W. D. GROPP, L. C. MCINNES, AND B. F. SMITH, *Efficient management of parallelism in object oriented numerical software libraries*, in Modern Software Tools in Scientific Computing, E. Arge, A. M. Bruaset, and H. P. Langtangen, eds., Birkhäuser Press, 1997, pp. 163–202.

- [57] N. V. BANICHUK, V. G. BELCPRIME SKIUI, AND V. V. KOBELEV, *Optimization in problems of elasticity with unknown boundaries*, Izvestiya Akademii Nauk SSSR. Mekhanika Tverdogo Tela, 19 (1984), pp. 46–52.
- [58] C. BARBAROSIE, *Shape optimization of periodic structures*, Computational Mechanics, 30 (2003), pp. 235–246.
- [59] C. BARBAROSIE AND S. LOPES, *A gradient-type algorithm for optimization with constraints*, submitted for publication, see also Pre-Print CMAF Pre-2011-001 at <http://cmaf.fc.ul.pt/preprints.html>, (2011).
- [60] C. BARBAROSIE, S. LOPES, AND A.-M. TOADER, *A gradient-type algorithm for constrained optimization with applications to multi-objective optimization of auxetic materials*, arXiv preprint arXiv:1711.04863, (2017).
- [61] G. BARLES, *Solutions de viscosité des équations de Hamilton-Jacobi*, vol. 17 of Mathématiques & Applications (Berlin), Springer-Verlag, Paris, 1994.
- [62] V. A. BARRY, B. A. GREGORY, AND N. ABUAF, *Turbine blade with enhanced cooling and profile optimization*, Nov. 9 1999. US Patent 5,980,209.
- [63] R. BECKER, V. HEUVELINE, AND R. RANNACHER, *An optimal control approach to adaptivity in computational fluid mechanics*, International journal for numerical methods in fluids, 40 (2002), pp. 105–120.
- [64] G. BELLETTINI, *Lecture notes on mean curvature flow, barriers and singular perturbations*, vol. 12 of Appunti. Scuola Normale Superiore di Pisa (Nuova Serie), Edizioni della Normale, Pisa, 2013.
- [65] M. P. BENDSOE AND N. KIKUCHI, *Generating optimal topologies in structural design using a homogenization method*, Computer methods in applied mechanics and engineering, 71 (1988), pp. 197–224.
- [66] M. P. BENDSOE AND O. SIGMUND, *Topology optimization*, Springer-Verlag, Berlin, 2003. Theory, methods and applications.
- [67] A. BENSALAH, *Une approche nouvelle de la modélisation mathématique et numérique en aéroacoustique par les équations de Goldstein et applications en aéronautique*, PhD thesis, Université Paris Saclay, 2018.
- [68] S. BERTOLUZZA, V. CHABANNES, C. PRUD’HOMME, AND M. SZOPOS, *Boundary conditions involving pressure for the Stokes problem and applications in computational hemodynamics*, Computer Methods in Applied Mechanics and Engineering, 322 (2017), pp. 58–80.
- [69] L. T. BIEGLER, *Nonlinear programming*, vol. 10 of MOS-SIAM Series on Optimization, Society for Industrial and Applied Mathematics (SIAM), 2010.
- [70] J.-F. BONNANS, J. C. GILBERT, C. LEMARÉCHAL, AND C. A. SAGASTIZÁBAL, *Numerical optimization: theoretical and practical aspects*, Springer Science & Business Media, 2006.
- [71] T. BORRVALL AND J. PETERSSON, *Large-scale topology optimization in 3D using parallel computing*, Computer methods in applied mechanics and engineering, 190 (2001), pp. 6201–6229.
- [72] ———, *Topology optimization of fluids in Stokes flow*, International Journal for Numerical Methods in Fluids, 41 (2003), pp. 77–107.
- [73] B. BOURDIN AND A. CHAMBOLLE, *Design-dependent loads in topology optimization*, ESAIM: Control, Optimisation and Calculus of Variations, 9 (2003), pp. 19–48.
- [74] A. BOURGEAT, E. MARUVSIC-PALOKA, AND A. MIKELIC, *Weak nonlinear corrections for Darcy’s law*, Mathematical Models and Methods in Applied Sciences, 6 (1996), pp. 1143–1155.
- [75] F. BOYER ET AL., *Trace theorems and spatial continuity properties for the solutions of the transport equation*, Differential and integral equations, 18 (2005), pp. 891–934.

- [76] H. BREZIS, *Functional analysis, Sobolev spaces and partial differential equations*, Springer Science & Business Media, 2010.
- [77] R. BRO AND S. DE JONG, *A fast non-negativity-constrained least squares algorithm*, *Journal of Chemometrics: A Journal of the Chemometrics Society*, 11 (1997), pp. 393–401.
- [78] D. BUCUR AND J. ZOLELIO, *Shape analysis of eigenvalues on manifolds*, in *Homogenization and applications to material sciences (Nice, 1995)*, vol. 9 of GAKUTO Internat. Ser. Math. Sci. Appl., Gakkotosho, Tokyo, 1995, pp. 67–79.
- [79] D. BUCUR AND J.-P. ZOLÉSIO, *Optimisation de forme sous contrainte capacitaire*, *Comptes Rendus de l'Académie des Sciences. Série I. Mathématique*, 318 (1994), pp. 795–800.
- [80] C. BUI, C. DAPOGNY, AND P. FREY, *An accurate anisotropic adaptation method for solving the level set advection equation*, *International Journal for Numerical Methods in Fluids*, 70 (2012), pp. 899–922.
- [81] M. BURGER, *A framework for the construction of level set methods for shape optimization and reconstruction*, *Interfaces and Free boundaries*, 5 (2003), pp. 301–329.
- [82] P. CANNARSA AND P. CARDALIAGUET, *Representation of equilibrium solutions to the table problem of growing sandpiles*, *Journal of the European Mathematical Society*, 6 (2004), pp. 435–464.
- [83] P. CANNARSA, G. DA PRATO, AND J.-P. ZOLÉSIO, *Dynamical shape control of the heat equation*, *Systems & Control Letters*, 12 (1989), pp. 103–109.
- [84] J. CÉA, *Conception optimale ou identification de formes, calcul rapide de la dérivée directionnelle de la fonction coût*, *ESAIM: Mathematical Modelling and Numerical Analysis*, 20 (1986), pp. 371–402.
- [85] J. CÉA, S. GARREAU, P. GUILLAUME, AND M. MASMOUDI, *The shape and topological optimizations connection*, *Computer Methods in Applied Mechanics and Engineering*, 188 (2000), pp. 713–726. IV WCCM, Part II (Buenos Aires, 1998).
- [86] J. CÉA, A. GIOAN, AND J. MICHEL, *Quelques résultats sur l'identification de domaines*, *Calcolo. A Quarterly on Numerical Analysis and Theory of Computation*, 10 (1973), pp. 207–232.
- [87] J. CÉA, A. GIOAN, AND J. MICHEL, *Adaptation de la méthode du gradient à un problème d'identification de domaine*, in *Computing methods in applied sciences and engineering (Proc. Internat. Sympos., Versailles, 1973)*, Part 2, 1974, pp. 391–402. *Lecture Notes in Comput. Sci.*, Vol. 11.
- [88] S. CHEN, M. Y. WANG, AND A. Q. LIU, *Shape feature control in structural topology optimization*, *Computer-Aided Design*, 40 (2008), pp. 951–962.
- [89] D. CHENAIS, *Optimal design of midsurface of shells: differentiability proof and sensitivity computation*, *Applied Mathematics and Optimization*, 16 (1987), pp. 93–133.
- [90] D. CHENAIS, J. MONNIER, AND J.-P. VILA, *Shape optimal design for a fluid-heat coupled system*, *Applied Mathematics and Computer Science*, 6 (1996), pp. 245–261. *Shape optimization and scientific computations (Warsaw, 1994)*.
- [91] D. CHENAIS, J. MONNIER, AND J. P. VILA, *Shape optimal design problem with convective and radiative heat transfer: analysis and implementation*, *Journal of Optimization Theory and Applications*, 110 (2001), pp. 75–117.
- [92] G. CHENG, Y. MEI, AND X. WANG, *A feature-based structural topology optimization method*, in *IUTAM Symposium on Topological Design Optimization of Structures, Machines and Materials*, Springer, 2006, pp. 505–514.
- [93] A. CHERKAEV, *Variational methods for structural optimization*, vol. 140, Springer Science & Business Media, 2012.
- [94] J. CHETBOUN, *Conception de formes aérodynamiques en présence d'écoulements décollés: contrôle et optimisation*, PhD thesis, École polytechnique, 2010.

- [95] C. CHICONE, *Ordinary Differential Equations with Applications*, Springer New York, 1999.
- [96] D. L. CHOPP, *Computing minimal surfaces via level set curvature flow*, Journal of Computational Physics, 106 (1993), pp. 77–91.
- [97] D. L. CHOPP, *Another look at velocity extensions in the level set method*, SIAM Journal on Scientific Computing, 31 (2009), pp. 3255–3273.
- [98] A. N. CHRISTIANSEN, J. A. BÆRENTZEN, M. NOBEL-JØRGENSEN, N. AAGE, AND O. SIGMUND, *Combined shape and topology optimization of 3D structures*, Computers & Graphics, 46 (2015), pp. 25–35.
- [99] A. N. CHRISTIANSEN, M. NOBEL-JØRGENSEN, N. AAGE, O. SIGMUND, AND J. A. BÆRENTZEN, *Topology optimization using an explicit interface representation*, Structural and Multidisciplinary Optimization, 49 (2014), pp. 387–399.
- [100] C. P. CHUKWUDOZIE, *Shape Optimization for Drag Minimization Using the Navier-Stokes Equation*, Master’s thesis, Louisiana State University, 2015.
- [101] H. CHUNG, O. AMIR, AND H. A. KIM, *Nonlinear Thermoelastic Topology Optimization with the Level-Set Method*, in AIAA Scitech 2019 Forum, 2019, p. 1470.
- [102] D. CIORANESCU, A. DAMLAMIAN, AND G. GRISO, *The periodic unfolding method in homogenization*, SIAM Journal on Mathematical Analysis, 40 (2008), pp. 1585–1620.
- [103] D. CIORANESCU AND F. MURAT, *A strange term coming from nowhere*, in Topics in the mathematical modelling of composite materials, vol. 31 of Progr. Nonlinear Differential Equations Appl., Birkhäuser Boston, Boston, MA, 1997, pp. 45–93.
- [104] P. COFFIN AND K. MAUTE, *Level set topology optimization of cooling and heating devices using a simplified convection model*, Structural and Multidisciplinary Optimization, 53 (2016), pp. 985–1003.
- [105] C. CONCA, F. MURAT, AND O. PIRONNEAU, *The Stokes and Navier-Stokes equations with boundary conditions involving the pressure*, Japanese journal of mathematics. New series, 20 (1994), pp. 279–318.
- [106] C. CONCA, C. PARES, O. PIRONNEAU, AND M. THIRIET, *Navier-Stokes equations with imposed pressure and velocity fluxes*, International journal for numerical methods in fluids, 20 (1995), pp. 267–287.
- [107] C. DAPOGNY, *Optimisation de formes, méthode des lignes de niveaux sur maillages non structurés et évolution de maillages*, PhD thesis, Université Pierre et Marie Curie-Paris VI, 2013.
- [108] C. DAPOGNY, C. DOBRZYNSKI, AND P. FREY, *Three-dimensional adaptive domain remeshing, implicit domain meshing, and applications to free and moving boundary problems*, Journal of Computational Physics, 262 (2014), pp. 358–378.
- [109] C. DAPOGNY, R. ESTEVEZ, A. FAURE, AND G. MICHAILIDIS, *Shape and topology optimization considering anisotropic features induced by additive manufacturing processes*, Hal preprint: <https://hal.archives-ouvertes.fr/hal-01660850/>, (2017).
- [110] C. DAPOGNY, A. FAURE, G. MICHAILIDIS, G. ALLAIRE, A. COUVELAS, AND R. ESTEVEZ, *Geometric constraints for shape and topology optimization in architectural design*, Computational Mechanics, 59 (2017), pp. 933–965.
- [111] C. DAPOGNY AND P. FREY, *Computation of the signed distance function to a discrete contour on adapted triangulation*, Calcolo, 49 (2012), pp. 193–219.
- [112] C. DAPOGNY, P. FREY, F. OMNÈS, AND Y. PRIVAT, *Geometrical shape optimization in fluid mechanics using FreeFem++*. HAL preprint hal-01481707, Mar. 2017.
- [113] ———, *Geometrical shape optimization in fluid mechanics using FreeFem++*, Structural and Multidisciplinary Optimization, (2017), pp. 1–28.

- [114] C. DAPOGNY, N. LEBBE, AND E. OUDET, *Optimization of the shape of regions supporting boundary conditions*. working paper or preprint, Mar. 2019.
- [115] G. DAVID AND S. SEMMES, *Uniform rectifiability and singular sets*, Annales de l'Institut Henri Poincaré (C) Non Linear Analysis, 13 (1996), pp. 383–443.
- [116] T. DBOUK, *A review about the engineering design of optimal heat transfer systems using topology optimization*, Applied Thermal Engineering, 112 (2017), pp. 841–854.
- [117] F. DE GOURNAY, *Velocity extension for the level-set method and multiple eigenvalues in shape optimization*, SIAM journal on control and optimization, 45 (2006), pp. 343–367.
- [118] E. M. DEDE, *Multiphysics topology optimization of heat transfer and fluid flow systems*, in proceedings of the COMSOL Users Conference, 2009.
- [119] M. DELFOUR AND J. ZOLÉSIO, *Shape analysis via distance functions*, J. Funct. Anal, 123 (1994), pp. 129–201.
- [120] M. DELFOUR AND J.-P. ZOLÉSIO, *Shape identification via metrics constructed from the oriented distance function*, Control and Cybernetics, 34 (2005), pp. 137–164.
- [121] M. C. DELFOUR, G. PAYRE, AND J.-P. ZOLÉSIO, *Shape optimal design of a radiating fin*, in System modelling and optimization (Copenhagen, 1983), vol. 59 of Lect. Notes Control Inf. Sci., Springer, Berlin, 1984, pp. 810–818.
- [122] M. C. DELFOUR AND J.-P. ZOLÉSIO, *Shapes and geometries: metrics, analysis, differential calculus, and optimization*, SIAM, 2011.
- [123] A. DERVIEUX, F. COURTY, M. VÁZQUEZ, AND B. KOOBUS, *Additive multilevel optimization and its application to sonic boom reduction*, in Numerical methods for scientific computing. Variational problems and applications, Internat. Center Numer. Methods Eng. (CIMNE), Barcelona, 2003, pp. 31–44.
- [124] T. K. DEY AND W. ZHAO, *Approximate medial axis as a voronoi subcomplex*, in Proceedings of the seventh ACM symposium on Solid modeling and applications, ACM, 2002, pp. 356–366.
- [125] D. A. DI PIETRO AND A. ERN, *Mathematical aspects of discontinuous Galerkin methods*, vol. 69, Springer Science & Business Media, 2011.
- [126] L. DIECI AND L. LOPEZ, *A survey of numerical methods for IVPs of ODEs with discontinuous right-hand side*, Journal of Computational and Applied Mathematics, 236 (2012), pp. 3967–3991.
- [127] J. DIEUDONNÉ, *Foundations of modern analysis*, Academic press, New York and London, 1960.
- [128] C. DILGEN, S. DILGEN, D. FUHRMAN, O. SIGMUND, AND B. LAZAROV, *Topology optimization of turbulent flows*, Computer Methods in Applied Mechanics and Engineering, 331 (2018), pp. 363–393.
- [129] S. B. DILGEN, C. B. DILGEN, D. R. FUHRMAN, O. SIGMUND, AND B. S. LAZAROV, *Density based topology optimization of turbulent flow heat transfer systems*, Structural and Multidisciplinary Optimization, 57 (2018), pp. 1905–1918.
- [130] Q. V. DINH, G. ROGÉ, C. SEVIN, AND B. STOUFFLET, *Shape optimization in computational fluid dynamics*, Revue Européenne des Éléments Finis. European Journal of Finite Elements, 5 (1996), pp. 569–594.
- [131] R. J. DI PERNA AND P.-L. LIONS, *Ordinary differential equations, transport theory and Sobolev spaces*, Inventiones mathematicae, 98 (1989), pp. 511–547.
- [132] V. DOLEAN, P. JOLIVET, AND F. NATAF, *An introduction to domain decomposition methods: algorithms, theory, and parallel implementation*, vol. 144, SIAM, 2015.
- [133] X. DUAN, Y. MA, AND R. ZHANG, *Optimal shape control of fluid flow using variational level set method*, Physics letters A, 372 (2008), pp. 1374–1379.

- [134] F. DUGAST, Y. FAVENNEC, C. JOSSET, Y. FAN, AND L. LUO, *Topology optimization of thermal fluid flows with an adjoint Lattice Boltzmann Method*, Journal of Computational Physics, 365 (2018), pp. 376–404.
- [135] P. D. DUNNING AND H. A. KIM, *Introducing the sequential linear programming level-set method for topology optimization*, Structural and Multidisciplinary Optimization, 51 (2015), pp. 631–643.
- [136] J. DUOANDIKOETXEA, *Forty years of Muckenhoupt weights*, Function Spaces and Inequalities, (2013), pp. 23–75.
- [137] T. DUPONT, M. F. WHEELER, ET AL., *A Galerkin procedure for approximating the flux on the boundary for elliptic and parabolic boundary value problems*, Revue française d’automatique, informatique, recherche opérationnelle. Analyse numérique, 8 (1974), pp. 47–59.
- [138] R. G. DURÁN AND M. A. MUSCHIETTI, *An explicit right inverse of the divergence operator which is continuous in weighted norms*, Univ., 2000.
- [139] M. C. DUTA, S. SHAHPAR, AND M. B. GILES, *Turbomachinery design optimization using automatic differentiated adjoint code*, in ASME Turbo Expo 2007: Power for Land, Sea, and Air, American Society of Mechanical Engineers, 2007, pp. 1435–1444.
- [140] P. DUYSINX, L. VAN MIEGROET, T. JACOBS, AND C. FLEURY, *Generalized shape optimization using X-FEM and level set methods*, in IUTAM Symposium on Topological Design Optimization of Structures, Machines and Materials, Springer, 2006, pp. 23–32.
- [141] A. EDELMAN, T. A. ARIAS, AND S. T. SMITH, *The geometry of algorithms with orthogonality constraints*, SIAM journal on Matrix Analysis and Applications, 20 (1998), pp. 303–353.
- [142] B. EPSTEIN, A. JAMESON, S. PEIGIN, D. ROMAN, N. HARRISON, AND J. VASSBERG, *Comparative study of three-dimensional wing drag minimization by different optimization techniques*, Journal of Aircraft, 46 (2009), pp. 526–541.
- [143] B. EREM AND D. H. BROOKS, *Differential geometric approximation of the gradient and hessian on a triangulated manifold*, in Proceedings/IEEE International Symposium on Biomedical Imaging: from nano to macro, NIH Public Access, 2011, p. 504.
- [144] A. ERN AND J.-L. GUERMOND, *Theory and practice of finite elements*, vol. 159, Springer Science & Business Media, 2013.
- [145] H. A. ESCHENAUER, V. V. KOBELLEV, AND A. SCHUMACHER, *Bubble method for topology and shape optimization of structures*, Structural optimization, 8 (1994), pp. 42–51.
- [146] L. C. EVANS, *Partial differential equations*, vol. 19 of Graduate Studies in Mathematics, American Mathematical Society, Providence, RI, second ed., 2010.
- [147] L. C. EVANS AND R. F. GARIEPY, *Measure theory and fine properties of functions*, Studies in Advanced Mathematics. CRC Press, 1992.
- [148] W. D. EVANS, *Weighted Sobolev Spaces*, Bulletin of the London Mathematical Society, 18 (1986), pp. 220–221.
- [149] R. D. FALGOUT AND U. M. YANG, *hypr: A library of high performance preconditioners*, in International Conference on Computational Science, Springer, 2002, pp. 632–641.
- [150] A. FAURE, *Optimisation de forme de matériaux et structures architecturés par la méthode des lignes de niveaux avec prise en compte des interfaces graduées*, PhD thesis, Université Grenoble Alpes, 2017.
- [151] A. FAURE, G. MICHAILIDIS, G. PARRY, N. VERMAAK, AND R. ESTEVEZ, *Design of thermoelastic multi-material structures with graded interfaces using topology optimization*, Structural and Multidisciplinary Optimization, 56 (2017), pp. 823–837.
- [152] E. FEIREISL AND Y. LU, *Homogenization of stationary Navier–Stokes equations in domains with tiny holes*, Journal of Mathematical Fluid Mechanics, 17 (2015), pp. 381–392.

- [153] F. FEPPON, G. ALLAIRE, F. BORDEU, J. CORTIAL, AND C. DAPOGNY, *Shape optimization of a coupled thermal fluid–structure problem in a level set mesh evolution framework*, SeMA Journal, (2019), pp. 1–46.
- [154] F. FEPPON, G. ALLAIRE, AND C. DAPOGNY, *A variational formulation for computing shape derivatives of geometric constraints along rays*, To appear in M2AN, (2019).
- [155] F. FEPPON, G. ALLAIRE, AND C. DAPOGNY, *Null space gradient flows for constrained optimization with applications to shape optimization*, Submitted, (2019).
- [156] F. FEPPON AND P. LERMUSIAUX, *The Extrinsic Geometry of Dynamical Systems Tracking Nonlinear Matrix Projections*, SIAM Journal on Matrix Analysis and Applications, 40 (2019), pp. 814–844.
- [157] F. FEPPON AND P. F. LERMUSIAUX, *A geometric approach to dynamical model order reduction*, SIAM Journal on Matrix Analysis and Applications, 39 (2018), pp. 510–538.
- [158] A. F. FILIPPOV, *Differential Equations with Discontinuous Righthand Sides: Control Systems*, vol. 18 of Mathematics and its applications, Springer Netherlands, 1988.
- [159] M. FIRDAOUSS, J.-L. GUERMOND, AND P. LE QUÉRÉ, *Nonlinear corrections to darcy’s law at low reynolds numbers*, Journal of Fluid Mechanics, 343 (1997), pp. 331–350.
- [160] Y. FISCHER, B. MARTEAU, AND Y. PRIVAT, *Some inverse problems around the Tokamak Tore Supra*, Communications on Pure and Applied Analysis, 11 (2012), pp. 2327–2349.
- [161] R. FLETCHER, *Practical methods of optimization*, John Wiley & Sons, 2013.
- [162] R. FRANSEN, *LES based aerothermal modeling of turbine blade cooling systems*, PhD thesis, Université de Toulouse, 2013.
- [163] G. P. GALDI, *Steady Stokes Flow in Exterior Domains*, Springer New York, New York, NY, 1994, pp. 244–303.
- [164] H. GARCKE, C. HECHT, M. HINZE, C. KAHLE, AND K. F. LAM, *Shape optimization for surface functionals in Navier-Stokes flow using a phase field approach*, Interfaces and Free Boundaries. Mathematical Analysis, Computation and Applications, 18 (2016), pp. 219–261.
- [165] S. GARREAU, P. GUILLAUME, AND M. MASMOUDI, *The topological asymptotic for PDE systems: the elasticity case*, SIAM Journal on Control and Optimization, 39 (2001), pp. 1756–1778.
- [166] P. GEOFFROY DONDERS, *Homogenization method for topology optimization of structures built with lattice materials.*, theses, Université Paris-Saclay, Dec. 2018.
- [167] C. GEUZAIN AND J.-F. REMACLE, *Gmsh: A 3-D finite element mesh generator with built-in pre-and post-processing facilities*, International journal for numerical methods in engineering, 79 (2009), pp. 1309–1331.
- [168] M. GIACOMINI, O. PANTZ, AND K. TRABELSI, *Volumetric expressions of the shape gradient of the compliance in structural shape optimization*, arXiv preprint arXiv:1701.05762, (2017).
- [169] R. GIANNONI, *Gas-liquid heat exchanger and method for its manufacture*, Oct. 7 1999. European Patent 1 098 156 A1.
- [170] F. GIBOU, R. FEDKIW, AND S. OSHER, *A review of level-set methods and some recent applications*, Journal of Computational Physics, 353 (2018), pp. 82 – 109.
- [171] V. GIRAULT AND P.-A. RAVIART, *Finite element approximation of the Navier-Stokes equations*, vol. 749, Springer Science & Business Media, 1979.
- [172] R. GLOWINSKI AND O. PIRONNEAU, *Towards the computation of minimum drag profiles in viscous laminar flow*, Applied Mathematical Modelling, 1 (1976/77), pp. 58–66.
- [173] G. H. GOLUB AND C. F. VAN LOAN, *Matrix computations*, The Johns Hopkins University Press, Baltimore, USA, (1989).

- [174] F. A. GOMES AND T. A. SENNE, *An SLP algorithm and its application to topology optimization*, Computational & Applied Mathematics, 30 (2011).
- [175] J. S. GRAY, J. T. HWANG, J. R. R. A. MARTINS, K. T. MOORE, AND B. A. NAYLOR, *OpenMDAO: an open-source framework for multidisciplinary design, analysis, and optimization*, Structural and Multidisciplinary Optimization, 59 (2019), pp. 1075–1104.
- [176] P. GRISVARD, *Elliptic problems in nonsmooth domains*, vol. 69 of Classics in Applied Mathematics, Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2011.
- [177] J. P. GROEN AND O. SIGMUND, *Homogenization-based topology optimization for high-resolution manufacturable microstructures*, International Journal for Numerical Methods in Engineering, 113 (2018), pp. 1148–1163.
- [178] J. K. GUEST, *Imposing maximum length scale in topology optimization*, Structural and Multidisciplinary Optimization, 37 (2009), pp. 463–473.
- [179] P. GUILLAUME AND M. MASMUDI, *Calcul numérique des dérivées d'ordre supérieur en conception optimale de formes*, Comptes Rendus de l'Académie des Sciences. Série I. Mathématique, 316 (1993), pp. 1091–1096.
- [180] J. HADAMARD, *Mémoire sur le problème d'analyse relatif à l'équilibre des plaques élastiques encastées*, vol. 33, Imprimerie nationale, 1908.
- [181] S. HARRIES, C. ABT, AND M. BRENNER, *Upfront CAD—Parametric modeling techniques for shape optimization*, in Advances in Evolutionary and Deterministic Methods for Design, Optimization and Control in Engineering and Sciences, Springer, 2019, pp. 191–211.
- [182] J. HASLINGER AND P. NEITTAANMÄKI, *Shape optimization in contact problems. Approximation and numerical realization*, RAIRO Modélisation Mathématique et Analyse Numérique, 21 (1987), pp. 269–291.
- [183] F. HECHT, *New development in FreeFem++*, Journal of Numerical Mathematics, 20 (2012), pp. 251–265.
- [184] A. HENROT AND M. PIERRE, *Variation et optimisation de formes: une analyse géométrique*, vol. 48, Springer Science & Business Media, 2006.
- [185] ———, *Shape variation and optimization*, vol. 28 of EMS Tracts in Mathematics, European Mathematical Society (EMS), Zürich, 2018.
- [186] A. HENROT AND Y. PRIVAT, *What is the optimal shape of a pipe?*, Archive for rational mechanics and analysis, 196 (2010), pp. 281–302.
- [187] M. HINTERMÜLLER, K. ITO, AND K. KUNISCH, *The primal-dual active set strategy as a semi-smooth Newton method*, SIAM Journal on Optimization, 13 (2002), pp. 865–888.
- [188] R. HIPTMAIR, A. PAGANINI, AND S. SARGHEINI, *Comparison of approximate shape gradients*, BIT Numerical Mathematics, 55 (2015), pp. 459–485.
- [189] G. P. HUANG, D. DOMAN, M. OPPENHEIMER, A. TIPTON, AND D. SIGTHORSSON, *Topology optimization of a fuel thermal management system*, in AIAA Aviation 2019 Forum, 2019, p. 3471.
- [190] K. ITO, H. T. TRAN, AND J. S. SCROGGS, *Mathematical issues in optimal design of a vapor transport reactor*, in Flow control (Minneapolis, MN, 1992), vol. 68 of IMA Vol. Math. Appl., Springer, New York, 1995, pp. 197–218.
- [191] A. JAMESON, *Aerodynamic design via control theory*, in Recent advances in computational fluid dynamics (Princeton, NJ, 1988), vol. 43 of Lecture Notes in Engrg., Springer, Berlin, 1989, pp. 377–401.
- [192] ———, *Computational algorithms for aerodynamic analysis and design*, Applied Numerical Mathematics. An IMACS Journal, 13 (1993), pp. 383–422.



- [193] I. G. JANG AND B. M. KWAK, *Evolutionary topology optimization using design space adjustment based on fixed grid*, International Journal for Numerical Methods in Engineering, 66 (2006), pp. 1817–1840.
- [194] N. JENKINS AND K. MAUTE, *Level set topology optimization of stationary fluid-structure interaction problems*, Structural and Multidisciplinary Optimization, 52 (2015), pp. 179–195.
- [195] M. JENSEN, *Discontinuous Galerkin methods for Friedrichs systems with irregular solutions*, PhD thesis, University of Oxford, 2005.
- [196] W. JING, *A unified homogenization approach for the Dirichlet problem in Perforated Domains*, arXiv preprint arXiv:1901.08251, (2019).
- [197] C. JOG, *Distributed-parameter optimization and topology design for non-linear thermoelasticity*, Computer Methods in Applied Mechanics and Engineering, 132 (1996), pp. 117–134.
- [198] P. JOLIVET, *FreeFEM tutorial*. <http://jolivet.perso.enseeiht.fr/FreeFem-tutorial/>, Apr 2019.
- [199] P. JOLIVET, F. HECHT, F. NATAF, AND C. PRUD’HOMME, *Scalable domain decomposition preconditioners for heterogeneous elliptic problems*, Scientific Programming, 22 (2014), pp. 157–171.
- [200] H. T. JONGEN AND O. STEIN, *On the complexity of equalizing inequalities*, Journal of Global Optimization, 27 (2003), pp. 367–374.
- [201] H. T. JONGEN AND O. STEIN, *Constrained global optimization: adaptive gradient flows*, in Frontiers in global optimization, vol. 74 of Nonconvex Optim. Appl., Kluwer Acad. Publ., Boston, MA, 2004, pp. 223–236.
- [202] S. KAMBAMPATI, C. JAUREGUI, K. MUSETH, AND H. A. KIM, *Fast level set topology optimization using a hierarchical data structure*, in 2018 Multidisciplinary Analysis and Optimization Conference, 2018, p. 3881.
- [203] D.-H. KIM, S.-H. LEE, I.-H. PARK, AND J.-H. LEE, *Derivation of a general sensitivity formula for shape optimization of 2-D magnetostatic systems by continuum approach*, IEEE Transactions on Magnetics, 38 (2002), pp. 1125–1128.
- [204] R. KIMMEL AND J. A. SETHIAN, *Computing geodesic paths on manifolds*, Proceedings of the national academy of Sciences, 95 (1998), pp. 8431–8435.
- [205] T. KONDOH, T. MATSUMORI, AND A. KAWAMOTO, *Drag minimization and lift maximization in laminar flows via topology optimization employing simple objective function expressions based on body force integration*, Structural and Multidisciplinary Optimization, 45 (2012), pp. 693–701.
- [206] A. KOSHAKJI, A. QUARTERONI, AND G. ROZZA, *Free form deformation techniques applied to 3D shape optimization problems*, Communications in Applied and Industrial Mathematics, 4 (2013), pp. e452, 26.
- [207] T. KRAINER AND B. SCHULZE, *Weighted Sobolev spaces*, Springer, 1985.
- [208] S. KREISSL, G. PINGEN, A. EVGRAFOV, AND K. MAUTE, *Topology optimization of flexible microfluidic devices*, Structural and Multidisciplinary Optimization, 42 (2010), pp. 495–516.
- [209] A. KUFNER AND B. OPIC, *How to define reasonably weighted Sobolev spaces*, Commentationes Mathematicae Universitatis Carolinae, 25 (1984), pp. 537–554.
- [210] O. A. LADYZHENSKAYA, *The mathematical theory of viscous incompressible flow*, vol. 2, Gordon and Breach New York, 1969.
- [211] A. B. LAMBE AND J. R. R. A. MARTINS, *Matrix-free aerostructural optimization of aircraft wings*, Structural and Multidisciplinary Optimization, 53 (2016), pp. 589–603.
- [212] S. LANG, *Fundamentals of differential geometry*, vol. 191, Springer Science & Business Media, 2012.

- [213] N. LEBBE, C. DAPOGNY, E. OUDET, K. HASSAN, AND A. GLIERE, *Robust shape and topology optimization of nanophotonic devices using the level set method*, Journal of Computational Physics, (2019).
- [214] L. LEIFSSON AND S. KOZIEL, *Aerodynamic shape optimization by variable-fidelity computational fluid dynamics models: a review of recent progress*, Journal of Computational Science, 10 (2015), pp. 45–54.
- [215] Y. LI AND L. NIRENBERG, *The distance function to the boundary, finsler geometry, and the singular set of viscosity solutions of some hamilton-jacobi equations*, Communications on pure and applied mathematics, 58 (2005), pp. 85–146.
- [216] J.-L. LIONS, *Perturbations singulières dans les problèmes aux limites et en contrôle optimal*, Lecture Notes in Mathematics, Vol. 323, Springer-Verlag, Berlin-New York, 1973.
- [217] J.-L. LIONS, *Some methods in the mathematical analysis of systems and their control*, Beijing, Science Press, (1981).
- [218] J.-L. LIONS AND E. MAGENES, *Non-Homogenous Boundary Value Problems and Applications*, vol. 1, Springer, New York, 01 1972.
- [219] W. G. LITVINOV, *On the optimal shape of a hydrofoil*, Journal of Optimization Theory and Applications, 85 (1995), pp. 325–345.
- [220] J. LIU AND Y. MA, *A survey of manufacturing oriented topology optimization methods*, Advances in Engineering Software, 100 (2016), pp. 161–175.
- [221] C. LUNDGAARD, J. ALEXANDERSEN, M. ZHOU, C. S. ANDREASEN, AND O. SIGMUND, *Revisiting density-based topology optimization for fluid-structure-interaction problems*, Structural and Multidisciplinary Optimization, 58 (2018), pp. 969–995.
- [222] J. LUO, Z. LUO, S. CHEN, L. TONG, AND M. Y. WANG, *A new level set method for systematic design of hinge-free compliant mechanisms*, Computer Methods in Applied Mechanics and Engineering, 198 (2008), pp. 318–331.
- [223] C. MANTEGAZZA AND A. C. MENNUCCI, *Hamilton-Jacobi Equations and Distance Functions on Riemannian Manifolds*, Applied Mathematics & Optimization, 47 (2003).
- [224] G. MARCK, M. NEMER, AND J.-L. HARION, *Topology optimization of heat and mass transfer problems: Laminar flow*, Numerical Heat Transfer, Part B: Fundamentals, 63 (2013), pp. 508–539.
- [225] G. MARCK AND Y. PRIVAT, *On some shape and topology optimization problems in conductive and convective heat transfers*, in OPTI 2014, An International Conference on Engineering and Applied Sciences Optimization, M. Papadarakakis, M. Karlaftis, and N. Lagaros, eds., June 2014, pp. 1640–1657.
- [226] N. MARCO AND A. DERVIEUX, *Multilevel parametrization for aerodynamical optimization of 3D shapes*, Finite Elements in Analysis and Design, 26 (1997), pp. 259–277.
- [227] J. MARTÍNEZ-FRUTOS AND D. HERRERO-PÉREZ, *GPU acceleration for evolutionary topology optimization of continuum structures using isosurfaces*, Computers & Structures, 182 (2017), pp. 119–136.
- [228] E. MARUVSIC-PALOKA, *Asymptotic expansion for a flow in a periodic porous medium*, Comptes Rendus de l’Académie des Sciences-Series IIB-Mechanics-Physics-Chemistry-Astronomy, 325 (1997), pp. 369–374.
- [229] M. MASMOUDI AND P. GUILLAUME, *Conception optimale de formes et applications*, in Optimisation et contrôle (Sophia-Antipolis, 1992), Cépaduès, Toulouse, 1993, pp. 201–216.
- [230] A. MAURY, G. ALLAIRE, AND F. JOUVE, *Elasto-plastic shape optimization using the level set method*, SIAM Journal on Control and Optimization, 56 (2018), pp. 556–581.
- [231] W. C. H. MCLEAN, *Strongly elliptic systems and boundary integral equations*, Cambridge university press, 2000.

- [232] L. A. M. MELLO, E. DE STURLER, G. H. PAULINO, AND E. C. N. SILVA, *Recycling Krylov subspaces for efficient large-scale electrical impedance tomography*, Computer Methods in Applied Mechanics and Engineering, 199 (2010), pp. 3101–3110.
- [233] M. MEYER, M. DESBRUN, P. SCHRÖDER, AND A. H. BARR, *Discrete differential-geometry operators for triangulated 2-manifolds*, in Visualization and mathematics III, Springer, 2003, pp. 35–57.
- [234] G. MICHAILIDIS, *Manufacturing Constraints and Multi-Phase Shape and Topology Optimization via a Level-Set Method*, PhD thesis, Ecole polytechnique, 2014.
- [235] B. MOHAMMADI AND B. MOHAMMADI, *Optimal shape design, reverse mode of automatic differentiation and turbulence*, in 35th Aerospace Sciences Meeting and Exhibit, 1997, p. 99.
- [236] B. MOHAMMADI AND O. PIRONNEAU, *Applied shape optimization for fluids*, Numerical Mathematics and Scientific Computation, Oxford University Press, Oxford, second ed., 2010.
- [237] P. MORIN, R. NOCHETTO, M. PAULETTI, AND M. VERANI, *Adaptive SQP method for shape optimization*, in Numerical Mathematics and Advanced Applications 2009, Springer, 2010, pp. 663–673.
- [238] J. MOULIN, P. JOLIVET, AND O. MARQUET, *Augmented Lagrangian preconditioner for large-scale hydrodynamic stability analysis*, Computer Methods in Applied Mechanics and Engineering, 351 (2019), pp. 718–743.
- [239] C. MULTIPHYSICS, *Comsol*, Inc., Burlington, MA, www.comsol.com, (1994).
- [240] D. J. MUNK, T. KIPOUROS, G. A. VIO, G. T. PARKS, AND G. P. STEVEN, *On the effect of fluid-structure interactions and choice of algorithm in multi-physics topology optimisation*, Finite Elements in Analysis and Design, 145 (2018), pp. 32–54.
- [241] F. MURAT AND J. SIMON, *Etude de problèmes d’optimal design*, in IFIP Technical Conference on Optimization Techniques, Springer, 1975, pp. 54–62.
- [242] ———, *Sur le contrôle par un domaine géométrique*, Publication du Laboratoire d’Analyse Numérique de l’Université Pierre et Marie Curie, (1976).
- [243] J.-C. NÉDÉLEC, *Acoustic and electromagnetic equations: integral representations for harmonic problems*, Springer Science & Business Media, 2001.
- [244] J. NOCEDAL AND S. J. WRIGHT, *Numerical optimization*, Springer Science, 35 (1999).
- [245] NOESIS SOLUTIONS, NV, *Optimus theoretical background*, Leuven, Belgium, (2010).
- [246] H. OKUMURA AND M. KAWAHARA, *Shape optimization of body located in incompressible Navier-Stokes flow based on optimal control theory*, CMES. Computer Modeling in Engineering & Sciences, 1 (2000), pp. 71–77.
- [247] S. OSHER AND R. FEDKIW, *Level set methods and dynamic implicit surfaces*, vol. 153 of Applied Mathematical Sciences, Springer-Verlag, New York, 2003.
- [248] S. OSHER AND J. A. SETHIAN, *Fronts propagating with curvature-dependent speed: algorithms based on Hamilton-Jacobi formulations*, Journal of computational physics, 79 (1988), pp. 12–49.
- [249] S. OSHER AND C.-W. SHU, *High-order essentially nonoscillatory schemes for Hamilton-Jacobi equations*, SIAM Journal on numerical analysis, 28 (1991), pp. 907–922.
- [250] S. J. OSHER AND F. SANTOSA, *Level set methods for optimization problems involving geometry and constraints. I. Frequencies of a two-density inhomogeneous drum*, J. Comput. Phys., 171 (2001), pp. 272–288.
- [251] B. PALMERIO AND A. DERVIEUX, *Identification de frontière dans le cas d’un problème de Dirichlet*, Comptes Rendus Hebdomadaires des Séances de l’Académie des Sciences. Séries A et B, 275 (1972), pp. A1111–A1113.

- [252] ———, *Une formule de Hadamard dans des problèmes d'identification de domaines*, Comptes Rendus Hebdomadaires des Séances de l'Académie des Sciences. Séries A et B, 280 (1975), pp. Aii, A1697–A1700.
- [253] O. PANTZ, *Sensibilité de l'équation de la chaleur aux sauts de conductivité*, Comptes Rendus Mathématique, 341 (2005), pp. 333–337.
- [254] O. PANTZ AND K. TRABELSI, *A post-treatment of the homogenization method for shape optimization*, SIAM Journal on Control and Optimization, 47 (2008), pp. 1380–1398.
- [255] P. PAPAZOGLU, *Topology optimization of heat exchangers*, Master's thesis, TU Delft, 2015.
- [256] P.-O. PERSSON, *Mesh generation for implicit geometries*, PhD thesis, Massachusetts Institute of Technology, 2005.
- [257] E. PEYNAUD, *Rayonnement sonore dans un écoulement subsonique complexe en régime harmonique: analyse et simulation numérique du couplage entre les phénomènes acoustiques et hydrodynamiques*, PhD thesis, Toulouse, INSA, 2013.
- [258] M. PIETROPAOLI, F. MONTOMOLI, AND A. GAYMANN, *Three-dimensional fluid topology optimization for heat transfer*, Structural and Multidisciplinary Optimization, 59 (2019), pp. 801–812.
- [259] O. PIRONNEAU, *On optimum profiles in Stokes flow*, Journal of Fluid Mechanics, 59 (1973), pp. 117–128.
- [260] ———, *Optimal shape design for elliptic systems*, in System Modeling and Optimization, Springer, 1982, pp. 42–66.
- [261] P. PLOTNIKOV AND J. SOKOŁOWSKI, *Compressible Navier-Stokes equations: theory and shape optimization*, vol. 73, Springer Science & Business Media, 2012.
- [262] P. I. PLOTNIKOV AND J. SOKOŁOWSKI, *Optimal shape control of airfoil in compressible gas flow governed by Navier-Stokes equations*, Evolution Equations and Control Theory, 2 (2013), pp. 495–516.
- [263] N. POLLINI, O. SIGMUND, C. S. ANDREASEN, AND J. ALEXANDERSEN, *A “poor man's” approach for high-resolution three-dimensional topology design for natural convection problems*, Advances in Engineering Software, 140 (2020), p. 102736.
- [264] T. N. POUCHON, *Effective models and numerical homogenization methods for long time wave propagation in heterogeneous media*, PhD thesis, EPFL, 2017.
- [265] J. RAUCH AND M. TAYLOR, *Potential and scattering theory on wildly perturbed domains*, Journal of Functional Analysis, 18 (1975), pp. 27–59.
- [266] A. REBEI, *Développement de méthode d'optimisation topologique adaptée aux écoulements en régime turbulent, application au cas des échangeurs de chaleur*, PhD thesis, Thèse de l'université PSL. Préparée à Mines ParisTech., 2019.
- [267] T. RICHTER, *Fluid-structure interactions*, vol. 118 of Lecture Notes in Computational Science and Engineering, Springer, Cham, 2017. Models, analysis and finite elements.
- [268] J. C. ROBINSON, *Infinite-dimensional dynamical systems: an introduction to dissipative parabolic PDEs and the theory of global attractors*, vol. 28, Cambridge University Press, 2001.
- [269] W. RUDIN, *Real and complex analysis*, Tata McGraw-Hill Education, 2006.
- [270] S. RUSINKIEWICZ, *Estimating curvatures and their derivatives on triangle meshes*, in 2nd International Symposium on 3D Data Processing, Visualization and Transmission, IEEE, 2004, pp. 486–493.
- [271] D. SALTZMAN, M. BICHNEVICIUS, S. LYNCH, T. W. SIMPSON, E. W. REUTZEL, C. DICKMAN, AND R. MARTUKANITZ, *Design and evaluation of an additively manufactured aircraft heat exchanger*, Applied Thermal Engineering, 138 (2018), pp. 254–263.

- [272] E. SANCHEZ-PALENCIA, *Fluid flow in porous media*, Non-homogeneous media and vibration theory, (1980), pp. 129–157.
- [273] E. SÁNCHEZ-PALENCIA, *On the asymptotics of the fluid flow past an array of fixed obstacles*, International Journal of Engineering Science, 20 (1982), pp. 1291–1301.
- [274] F. SANTOSA AND W. W. SYMES, *A dispersive effective medium for wave propagation in periodic composites*, SIAM Journal on Applied Mathematics, 51 (1991), pp. 984–1005.
- [275] K. R. SAVIERS, R. RANJAN, AND R. MAHMOUDI, *Design and validation of topology optimized heat exchangers*, in AIAA Scitech 2019 Forum, 2019, p. 1465.
- [276] T. SCHNEIDER AND R. KLEIN, *Overcoming mass losses in level set-based interface tracking schemes*, in Finite volumes for complex applications II, Hermes Sci. Publ., Paris, 1999, pp. 41–50.
- [277] J. SCHROPP AND I. SINGER, *A dynamical systems approach to constrained minimization*, Numerical functional analysis and optimization, 21 (2000), pp. 537–551.
- [278] V. H. SCHULZ, *A Riemannian view on shape optimization*, Foundations of Computational Mathematics. The Journal of the Society for the Foundations of Computational Mathematics, 14 (2014), pp. 483–501.
- [279] A. SEGAL, M. UR REHMAN, AND C. VUIK, *Preconditioners for incompressible Navier-Stokes solvers*, Numerical Mathematics: Theory, Methods and Applications, 3 (2010), pp. 245–275.
- [280] J. A. SETHIAN, *A fast marching level set method for monotonically advancing fronts*, Proceedings of the National Academy of Sciences, 93 (1996), pp. 1591–1595.
- [281] J. A. SETHIAN AND A. WIEGMANN, *Structural boundary design via level set and immersed interface methods*, Journal of computational physics, 163 (2000), pp. 489–528.
- [282] V. SHIKHMAN AND O. STEIN, *Constrained optimization: projected gradient flows*, Journal of optimization theory and applications, 140 (2009), pp. 117–130.
- [283] H. SI AND A. TETGEN, *A quality tetrahedral mesh generator and three-dimensional delaunay triangulator*, Weierstrass Institute for Applied Analysis and Stochastic, Berlin, Germany, 81 (2006).
- [284] O. SIGMUND AND K. MAUTE, *Topology optimization approaches*, Structural and Multidisciplinary Optimization, 48 (2013), pp. 1031–1055.
- [285] O. SIGMUND AND S. TORQUATO, *Design of materials with extreme thermal expansion using a three-phase topology optimization method*, Journal of the Mechanics and Physics of Solids, 45 (1997), pp. 1037–1067.
- [286] S. N. SKINNER AND H. ZARE-BEHTASH, *State-of-the-art in aerodynamic shape optimisation methods*, Applied Soft Computing, 62 (2018), pp. 933–962.
- [287] V. P. SMYSHLYAEV AND K. CHEREDNICHENKO, *On rigorous derivation of strain gradient effects in the overall behaviour of periodic heterogeneous media*, Journal of the Mechanics and Physics of Solids, 48 (2000), pp. 1325–1357.
- [288] J. SOBIESZCZANSKI-SOBIESKI AND R. T. HAFTKA, *Multidisciplinary aerospace design optimization: survey of recent developments*, Structural optimization, 14 (1997), pp. 1–23.
- [289] J. SOKOŁOWSKI AND A. ZOCHOWSKI, *Topological derivative in shape optimization*, Encyclopedia of Optimization, (2009), pp. 3908–3918.
- [290] J. SOKOŁOWSKI AND J.-P. ZOLÉSIO, *Shape sensitivity analysis of contact problem with prescribed friction*, Nonlinear Analysis. Theory, Methods & Applications. An International Multidisciplinary Journal, 12 (1988), pp. 1399–1411.
- [291] J. SOKOŁOWSKI AND J.-P. ZOLÉSIO, *Introduction to shape optimization*, vol. 16 of Springer Series in Computational Mathematics, Springer-Verlag, Berlin, 1992. Shape sensitivity analysis.

- [292] M. SONNTAG, S. SCHMIDT, AND N. R. GAUGER, *Shape derivatives for the compressible Navier-Stokes equations in variational form*, Journal of Computational and Applied Mathematics, 296 (2016), pp. 334–351.
- [293] M. SPIVAK, *A comprehensive introduction to differential geometry. Vol. III*, Publish or Perish, Inc., Wilmington, Del., second ed., 1979.
- [294] G. STAMPACCHIA, *Le problème de Dirichlet pour les équations elliptiques du second ordre à coefficients discontinus*, Université de Grenoble. Annales de l’Institut Fourier, 15 (1965), pp. 189–258.
- [295] J. STRAIN, *Semi-Lagrangian methods for level set equations*, Journal of Computational Physics, 151 (1999), pp. 498–533.
- [296] K. STURM, *Shape optimization with nonsmooth cost functions: from theory to numerics*, SIAM Journal on Control and Optimization, 54 (2016), pp. 3319–3346.
- [297] W. SUN AND Y.-X. YUAN, *Optimization theory and methods: nonlinear programming*, vol. 1, Springer Science & Business Media, 2006.
- [298] K. SVANBERG, *The method of moving asymptotes—a new method for structural optimization*, International Journal for Numerical Methods in Engineering, 24 (1987), pp. 359–373.
- [299] A. TAKEZAWA, S. NISHIWAKI, AND M. KITAMURA, *Shape and topology optimization based on the phase field method and sensitivity analysis*, Journal of Computational Physics, 229 (2010), pp. 2697–2718.
- [300] K. TANABE, *A geometric method in nonlinear programming*, Journal of Optimization Theory and Applications, 30 (1980), pp. 181–210.
- [301] L. TARTAR, *Topics in nonlinear analysis*, Publications mathématiques d’Orsay, 78 (1978).
- [302] ———, *An introduction to Sobolev spaces and interpolation spaces*, vol. 3, Springer Science & Business Media, 2007.
- [303] R. TAWK, B. GHANNAM, AND M. NEMER, *Topology optimization of heat and mass transfer problems in two fluids-one solid domains*, Numerical Heat Transfer, Part B: Fundamentals, 76 (2019), pp. 130–151.
- [304] R. TEMAM, *Navier stokes equations: Theory and numerical analysis*, vol. 45, North-Holland Publishing Company, 1977.
- [305] L. N. TREFETHEN AND D. BAU III, *Numerical linear algebra*, vol. 50, Siam, 1997.
- [306] B. O. TURESSON, *Nonlinear potential theory and weighted Sobolev spaces*, Springer, 2007.
- [307] M. G. UKKEN AND M. SIVAPRAGASAM, *Aerodynamic shape optimization of airfoils at ultra-low Reynolds numbers*, Sādhanā, 44 (2019), p. 130.
- [308] G. N. VANDERPLAATS AND F. MOSES, *Structural optimization by methods of feasible directions*, Computers & Structures, 3 (1973), pp. 739–755.
- [309] C. H. VILLANUEVA AND K. MAUTE, *CutFEM topology optimization of 3D laminar incompressible flow problems*, Computer Methods in Applied Mechanics and Engineering, 320 (2017), pp. 444–473.
- [310] A. WÄCHTER AND L. T. BIEGLER, *On the implementation of an interior-point filter line-search algorithm for large-scale nonlinear programming*, Mathematical programming, 106 (2006), pp. 25–57.
- [311] M. Y. WANG, X. WANG, AND D. GUO, *A level set method for structural topology optimization*, Computer methods in applied mechanics and engineering, 192 (2003), pp. 227–246.
- [312] K. WELKER, *Suitable Spaces for Shape Optimization*, arXiv preprint arXiv:1702.07579, (2017).
- [313] O. WIDLUND AND M. DRYJA, *An additive variant of the Schwarz alternating method for the case of many subregions*, Technical Report 339, Ultracomputer Note 131, Department of Computer Science, Courant Institute, 12 1987.

- [314] Q. XIA AND M. Y. WANG, *Topology optimization of thermoelastic structures using level set method*, Computational Mechanics, 42 (2008), pp. 837–857.
- [315] K. YAJI, T. YAMADA, S. KUBO, K. IZUI, AND S. NISHIWAKI, *A topology optimization method for a coupled thermal–fluid problem using level set boundary expressions*, International Journal of Heat and Mass Transfer, 81 (2015), pp. 878–888.
- [316] T. YAMADA, K. IZUI, S. NISHIWAKI, AND A. TAKEZAWA, *A topology optimization method based on the level set method incorporating a fictitious interface energy*, Computer Methods in Applied Mechanics and Engineering, 199 (2010), pp. 2876–2891.
- [317] H. YAMASHITA, *A differential equation approach to nonlinear programming*, Mathematical Programming, 18 (1980), pp. 155–168.
- [318] G. H. YOON, *Topological layout design of electro–fluid–thermal–compliant actuator*, Computer Methods in Applied Mechanics and Engineering, 209–212 (2012), pp. 28 – 44.
- [319] ———, *Stress-based topology optimization method for steady-state fluid–structure interaction problems*, Computer Methods in Applied Mechanics and Engineering, 278 (2014), pp. 499 – 523.
- [320] G. H. YOON, S. HEO, AND Y. Y. KIM, *Minimum thickness control at various levels for topology optimization using the wavelet method*, International journal of solids and structures, 42 (2005), pp. 5945–5970.
- [321] Y.-X. YUAN, *A review of trust region algorithms for optimization*, in ICIAM99: Proceedings of the Fourth International Congress on Industrial and Applied Mathematics, 09 1999.
- [322] M. YULIN AND W. XIAOMING, *A level set method for structural topology optimization with multi-constraints and multi-materials*, Acta Mechanica Sinica, 20 (2004), pp. 507–518.
- [323] H. ZHAO, *A fast sweeping method for eikonal equations*, Mathematics of computation, 74 (2005), pp. 603–627.
- [324] X. ZHAO, M. ZHOU, O. SIGMUND, AND C. ANDREASEN, *A "poor man's approach" to topology optimization of cooling channels based on a Darcy flow model*, International Journal of Heat and Mass Transfer, 116 (2018), pp. 1108–1123.
- [325] C. ZHUANG, Z. XIONG, AND H. DING, *A level set method for topology optimization of heat conduction problem under multiple load cases*, Computer Methods in Applied Mechanics and Engineering, 196 (2007), pp. 1074–1084.
- [326] J.-P. ZOLÉSIO, *The material derivative (or speed) method for shape optimization*, in Optimization of distributed parameter structures, Vol. II (Iowa City, Iowa, 1980), vol. 50 of NATO Adv. Study Inst. Ser. E: Appl. Sci., Nijhoff, The Hague, 1981, pp. 1089–1151.
- [327] G. ZOUTENDIJK, *Methods of feasible directions: A study in linear and non-linear programming*, Elsevier Publishing Co., Amsterdam-London-New York-Princeton, N.J., 1960.

**Titre :** Optimisation de formes de systèmes multiphysiques

**Mots clés :** Optimisation topologique, remaillage, transfert thermique convectif, interaction fluide-structure, contraintes géométriques, modèles homogénéisés d'ordres élevés.

**Résumé :** Cette thèse est consacrée à l'optimisation de la topologie et de la forme de systèmes multiphysiques motivés par des applications de l'industrie aéronautique. Nous calculons les dérivées de forme de fonctions de coût arbitraires pour un modèle fluide, thermique et mécanique faiblement couplé. Nous développons ensuite un algorithme de type gradient adapté à la résolution de problèmes d'optimisation de formes sous contraintes qui ne requiert par de réglage de paramètres non physiques. Nous introduisons ensuite une méthode variationnelle qui permet de calculer des intégrales le long de rayons sur un maillage par la résolution d'un problème variationnel qui ne requiert pas la détermination explicite de ces lignes sur la discrétisation spatiale. Cette technique nous a ainsi permis d'imposer une contrainte de non-mélange de phases pour une application à l'optimisation d'échangeurs de chaleur bitubes. Tous ces ingrédients ont été employés pour

traiter une variété de cas tests d'optimisation de formes pour des systèmes multi-physiques 2-d ou 3-d. Nous avons considéré des problèmes à une seule, deux ou bien trois physiques couplées en 2-d, et des problèmes de tailles relativement élevées en 3-d pour la mécanique, la conduction thermique, l'optimisation de profils aérodynamiques, et de la forme de systèmes en interaction fluide-structure. Un dernier chapitre d'ouverture est consacré à l'étude de modèles homogénéisés d'ordres élevés pour les équations de Stokes en milieu poreux. Ces équations d'ordres élevés englobent les trois régimes homogénéisés classiques—Stokes, Brinkman et Darcy—associés à divers rapports d'échelles pour la taille des obstacles. Elles pourraient permettre, lors de futurs travaux, de développer de nouvelles méthodes d'optimisation pour la conception de systèmes fluides caractérisés par des motifs multiéchelles, tels que les échangeurs thermiques industriels.

**Title :** Shape and topology optimization of multiphysics systems

**Keywords :** Topology optimization, remeshing, convective heat transfer, fluid-structure interaction, geometric constraints, high order homogenization.

**Abstract :** This work is devoted to shape and topology optimization of multiphysics systems motivated by aeronautic industrial applications. Shape derivatives of arbitrary objective functionals are computed for a weakly coupled thermal fluid-structure model. A novel gradient flow type algorithm is then developed for solving generic constrained shape optimization problems without the need for tuning non-physical metaparameters. Motivated by the need for enforcing non-mixing constraints in the design of liquid-liquid heat exchangers, a variational method is developed in order to simplify the numerical evaluation of geometric constraints: it allows to compute line integrals on a mesh by solving a variational problem without requiring the explicit knowledge of these lines on the spatial discretization. All these ingredients allowed us to im-

plement a variety of 2-d and 3-d multiphysics shape optimization test cases: from single, double or three physics problems in 2-d, to moderately large-scale 3-d test cases for structural design, thermal conduction, aerodynamic design and a fluid-structure interacting system. A final opening chapter derives high order homogenized equations for the Stokes system in a porous medium. These high order equations encompass the three classical homogenized regimes—namely Stokes, Brinkman and Darcy—associated with different obstacle's size scalings. They could allow, in future works, to develop new topology optimization methods for the design of fluid systems characterized by multi-scale patterns such as industrial heat exchangers.

