



**HAL**  
open science

# The Reduced Basis Method Applied to Aerothermal Simulations

Jean-Baptiste Wahl

► **To cite this version:**

Jean-Baptiste Wahl. The Reduced Basis Method Applied to Aerothermal Simulations. Analysis of PDEs [math.AP]. Université de Strasbourg, 2018. English. NNT: . tel-01869098v2

**HAL Id: tel-01869098**

**<https://theses.hal.science/tel-01869098v2>**

Submitted on 6 Nov 2018 (v2), last revised 19 Jun 2019 (v3)

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

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

# Thèse

INSTITUT DE  
RECHERCHE  
MATHÉMATIQUE  
AVANCÉE

UMR 7501

Strasbourg

présentée pour obtenir le grade de docteur de  
l'Université de Strasbourg  
Spécialité MATHÉMATIQUES APPLIQUÉES

**Jean-Baptiste Wahl**

**The Reduced Basis Method Applied to  
Aerothermal Simulations**

Soutenue le 13 Septembre 2018  
devant la commission d'examen

Martin Grepl, rapporteur

Anthony Nouy, rapporteur

Fabian Mangeant, examinateur

Robert Mosé, examinateur

Christophe Prud'homme, directeur de thèse

Yannick Hoarau, co-directeur de thèse

[www-irma.u-strasbg.fr](http://www-irma.u-strasbg.fr)







## UNIVERSITÉ DE STRASBOURG

École Doctorale Mathématiques, Sciences de l'Information et de l'Ingénieur  
Institut de Recherche Mathématique Avancée, UMR 7501

**THÈSE** présentée par

**Jean-Baptiste WAHL**

soutenue le 13 septembre 2018

pour obtenir le grade de Docteur de l'Université de Strasbourg

Spécialité : Mathématiques Appliquées

# The Reduced Basis Method Applied to Aerothermal Simulations

### THÈSE DIRIGÉE PAR :

Christophe PRUD'HOMME

Yannick HOARAU

Directeur de thèse, Université de Strasbourg

Co-directeur de thèse, Université de Strasbourg

### RAPPORTEURS :

Martin GREPL

Anthony NOUY

RWTH Aachen University

École Centrale Nantes

### EXAMINATEURS :

Fabien MANGEANT

Robert MOSÉ

Renault

Université de Strasbourg



*Life is like a box of chocolates.  
You never know what you're gonna get.*

FORREST'S MOTHER



*À Toi, Mam'*



---

## Acknowledgment

---

Je commencerai naturellement par remercier mon directeur de thèse Christophe Prud'homme. Tu m'as fait découvrir l'analyse numérique en première année de master et tu as su me communiquer ta passion pour cette discipline. Tu m'as ensuite très vite fait confiance pour travailler avec toi et ton équipe. Merci de m'avoir donné l'opportunité de travailler sur cette thèse. J'ai énormément appris à tes côtés et je t'en serai toujours reconnaissant.

Merci à Yannick Hoarau, mon co-directeur de thèse, pour son soutien et sa patience. Le travail sur la turbulence fut plus douloureux que prévu mais tu as toujours été de bon conseil et présent quand j'en ai eu besoin.

Merci à mes rapporteurs Anthony Nouy et Martin Grepl qui ont accepté de relire et d'apprécier mon manuscrit. Merci à Fabien Mangeant et Robert Mosé d'avoir accepté de faire parti de mon jury et de contribuer à l'évaluation de mon travail.

Je remercie également le laboratoire de l'IRMA et l'Université de Strasbourg pour m'avoir accueilli pendant ces années de thèse. Merci à tous mes collègues, et particulièrement à mes co-bureau, pour les échanges et les moments partagés.

Merci à toute l'équipe Feel++ pour m'avoir embarqué dans cette belle aventure . Merci à Stéphane et Cécile pour m'avoir mis le pied à l'étrier avec les bases réduites et pour m'avoir aidé à rentrer dans le code, sans vous j'y serai probablement encore perdu. Merci à Lorenzo pour sa bonne humeur au travail et pendant le voyage au US (on remet ça bientôt !). Merci à Guillaume pour les pauses café interminables à refaire le monde, à regarder les canards ou à râler sur tout et n'importe quoi. Merci à Céline pour son soutien indéfectible et pour la grimpe. Un énorme merci à Vincent C. pour son aide précieuse et sa bonne humeur; rares sont les développements que j'aurais pu réaliser sans ton aide ! Un autre énorme merci à Romain pour les

échanges sérieux et moins sérieux, pour les barres de rire, pour les bières, pour m'avoir supporté pendant ces années et pendant tous nos déplacements ensemble. Enfin, un merci tout particulier à Vincent H., le meilleurs nounours qu'on puisse trouver quand on bloque sur un problème !

Merci à ma famille et mes amis qui m'ont aidé et supporté dans cette aventure pendant toutes ces années. Un énorme merci à mes parents qui n'ont jamais perdu espoir ! Vous m'avez toujours soutenu dans mes projets et fais confiance pour réussir. Votre foi dans mes capacités m'a fait avancé pendant toutes ces années. Cette thèse est aussi pour vous et surtout pour toi Maman.

Enfin, le plus gros des merci pour mon Pou et ma Biquette. Cette thèse c'est aussi la votre. Je l'ai écrite, mais on a l'a vécu et souffert ensemble. C'est en vous que j'ai trouvé la force d'aller jusqu'au bout.



---

# Contents

---

<b>Notations</b>	<b>xv</b>
<b>Introduction</b>	<b>1</b>
<b>1 Preliminary Notions</b>	<b>9</b>
1 Function Spaces . . . . .	9
2 The Finite Element Method . . . . .	12
<b>I Aerothermal Simulations</b>	<b>17</b>
<b>2 Modeling of an Aerothermal Problem</b>	<b>19</b>
1 The Navier Stokes Equations and the Boussinesq Approximation . .	20
2 Finite Element Approach . . . . .	23
3 Solution Strategy . . . . .	26
4 Validation . . . . .	32
<b>3 Stabilization Methods</b>	<b>45</b>

1	Streamline Diffusion Methods . . . . .	47
2	Stabilization Parameter . . . . .	50
3	Shock Capturing Method . . . . .	55
4	Validation . . . . .	57
<b>4</b>	<b>Turbulence Modeling</b>	<b>75</b>
1	Reynolds Averaged Navier-Stokes Method . . . . .	76
2	Spalart Allmaras Model . . . . .	79
3	$k - \omega$ SST model . . . . .	80
4	Validation: 2D Naca0012 Airfoil . . . . .	84
<b>II</b>	<b>Model Order Reduction</b>	<b>87</b>
<b>5</b>	<b>The Reduced Basis Method</b>	<b>89</b>
1	Method and Notations : Application to Coercive Problems . . . . .	90
2	Reduced Basis for Saddle Point Problems . . . . .	98
3	Numerical Applications : RB for Stokes Problem . . . . .	105
<b>6</b>	<b>Non-Affine and Non-Linear Problems</b>	<b>113</b>
1	The Empirical Interpolation Method . . . . .	114
2	EIM for Discrete Operators . . . . .	117
3	Simultaneous Empirical Interpolation and Reduced Basis Method for Non-linear Problems . . . . .	121
4	Preliminary Results . . . . .	122
<b>7</b>	<b>The Reduced Basis Method Applied to Aerothermal Problems</b>	<b>129</b>
1	Reduction of the Aerothermal Model . . . . .	130
2	OPUS Test-Case: Application to the Cooling of Electronic Components	139

<i>CONTENTS</i>	xiii
3 CHORUS Test-Case: Application to Aerothermal Simulations in an Airplane . . . . .	150
<b>III Contribution to Feel++ Library</b>	<b>161</b>
<b>8 Aerothermal Framework</b>	<b>163</b>
1 The Navier Stokes Solver . . . . .	164
2 The Advection Diffusion Solver . . . . .	165
3 Turbulence Models . . . . .	169
<b>9 Reduced Basis Framework</b>	<b>175</b>
1 Reduced Basis by Block . . . . .	177
2 Reduced Basis for Aerothermal Problems . . . . .	183
<b>10 EIM Framework for Discrete Operators</b>	<b>187</b>
1 Design of the Implementation . . . . .	187
2 Construction of the Interpolation Space . . . . .	189
<b>Conclusion</b>	<b>191</b>
<b>Appendices</b>	<b>199</b>
1 Adaptive Time Stepping and Benchmarks . . . . .	199
2 Nitsche’s Method for Slip Condition . . . . .	210
3 Offline/Online Strategy for Residual Evaluation in the CRB . . . . .	215
4 Detailed Affine Decomposition for RB with Stokes Flow . . . . .	218
<b>Résumé de Thèse en Français</b>	<b>221</b>
1 Modélisation d’un problème d’aérothermie . . . . .	226
2 Réduction d’un problème d’aérothermie . . . . .	233

3 Applications numériques . . . . . 239

**Bibliography** . . . . . **240**

## Acronyms

FE	: Finite Element
SUPG	: Streamline Upwind/Petrov-Galerkin
GLS	: Galerkin Least Square
RANS	: Reynolds Averaged Navier-Stokes
DNS	: Direct Numerical Simulation
SA	: Spalart-Allmaras
RB	: Reduced Basis
CRB	: Certified Reduced Basis
EIM	: Empirical Interpolation Method
SER	: Simultaneous EIM and RB

## Functional Analysis

$\Omega$	: Regular Bounded Domain	$\Omega \subset \mathbb{R}^d, d = 1, 2, 3$
$\mathcal{L}^p(\Omega)$	: Lebesgue Space on $\Omega$	section 1.1
$H^k(\Omega)$	: Sobolev Space on $\Omega$	section 1.2
$\gamma_a$	: Continuity Constant of $a$	
$\alpha_a$	: Coercivity Constant of $a$	
$\beta_{Br}$	: Brezzi inf-sup Constant	



$\beta_{Ba}$	: Babuška inf-sup Constant	
$\nabla u$	: Gradient of scalar field $u$	$\nabla u = \left( \frac{\partial u}{\partial x_1}, \frac{\partial u}{\partial x_2} \right)$
$\Delta u$	: Laplacian of a scalar field $u$	$\Delta u = \frac{\partial^2 u}{\partial x_1^2} + \frac{\partial^2 u}{\partial x_2^2}$
$Hu$	: Hessian matrix of a scalar field $u$	$Hu = \begin{pmatrix} \frac{\partial^2 u}{\partial x_1^2} & \frac{\partial^2 u}{\partial x_1 x_2} \\ \frac{\partial^2 u}{\partial x_2 x_1} & \frac{\partial^2 u}{\partial x_2^2} \end{pmatrix}$
$\nabla \cdot \mathbf{u}$	: Divergence of a vector field $\mathbf{u}$	$\nabla \cdot \mathbf{u} = \frac{\partial u_1}{\partial x_1} + \frac{\partial u_2}{\partial x_2}$
$\nabla \mathbf{u}$	: Gradient of vector field $\mathbf{u}$	$\nabla \mathbf{u} = \begin{pmatrix} \frac{\partial u_1}{\partial x_1} & \frac{\partial u_1}{\partial x_2} \\ \frac{\partial u_2}{\partial x_1} & \frac{\partial u_2}{\partial x_2} \end{pmatrix}$
$\Delta \mathbf{u}$	: Laplacian of a vector field $\mathbf{u}$	$\Delta \mathbf{u} = \begin{pmatrix} \frac{\partial^2 u_1}{\partial x_1^2} + \frac{\partial^2 u_1}{\partial x_2^2} \\ \frac{\partial^2 u_2}{\partial x_1^2} + \frac{\partial^2 u_2}{\partial x_2^2} \end{pmatrix}$
$H\mathbf{u}$	: Hessian rank-3 tensor of a vector field $\mathbf{u}$	$(H\mathbf{u})_k = \begin{pmatrix} \frac{\partial^2 u_k}{\partial x_1^2} & \frac{\partial^2 u_k}{\partial x_1 x_2} \\ \frac{\partial^2 u_k}{\partial x_2 x_1} & \frac{\partial^2 u_k}{\partial x_2^2} \end{pmatrix}$

## Finite Element Method

$\mathcal{T}_h$	: Mesh on $\Omega$ of characteristic size $h$	
$\mathbf{n}$	: Normal Vector to the Current Face	
$\mathcal{N}$	: FE space dimension	
$\mathbb{L}_{ch}^k(\Omega)$	: Continuous Lagrange FE space on $\Omega$	section 2.2
$\mathbb{TH}_{ch}^k(\Omega)$	: Taylor-Hood FE space on $\Omega$	$\mathbb{TH}_{ch}^k = [\mathbb{L}_{ch}^{k+1}]^d \times \mathbb{L}_{ch}^k$

## Aerothermal Simulation

$\rho$	: Fluid Density	$\text{kg.m}^{-3}$
$\mu$	: Fluid Dynamic Viscosity	$\text{Pa.s}$
$\nu$	: Fluid Kinematic Viscosity	$\text{m}^2.\text{s}^{-1}$
$\kappa$	: Fluid Thermal Diffusivity	$\text{m}^2.\text{s}^{-1}$
$\mathbf{T}$	: Cauchy Stress Tensor	$\mathbf{T}(\mathbf{u}, p) = 2\mathcal{T}(\mathbf{u}) - Ip$
$\mathcal{T}$	: Viscous Stress Tensor	$\mathcal{T}(\mathbf{u}) = \mu \mathbf{S}(\mathbf{u})$
$\mathbf{S}$	: Strain Tensor	$\mathbf{S}(\mathbf{u}) = \frac{1}{2}(\nabla \mathbf{u} + \nabla \mathbf{u}^\top)$
$V/V_h$	: Velocity function space	
$Q/Q_h$	: Pressure function space	
$X/X_h$	: Temperature function space	
$\mathbf{u}$	: Velocity Vector Field	
$T$	: Temperature Scalar Field	
$p$	: Pressure Scalar Field	

## Turbulence Modelization

$R$	:	Reynolds Stress Tensor
$k$	:	Kinetic Energy
$\mu_t$	:	Eddy Viscosity

## Reduced Basis Method

$\mathcal{D}$	:	Parameter Space	$\mathcal{D} \in \mathbb{R}^p, p \in \mathbb{N}^*$
$\Xi$	:	Super sampling for RBM offline part	
$X_N$	:	Reduced Basis Space of size $N$	
$S_N$	:	Sample of parameters on which $X_N$ is built	



## **ANR Project CHORUS**

This thesis has been funded by the French National Research Agency (ANR) and the project CHORUS (Common Horizon of Open Research in Uncertainty for Simulation) started in 2013. This project aimed to support the research activities and the software development in the domain of the uncertainty quantification. The work was split into three main topics: the academic research, the software development and the benchmarking on industrial test-cases. All the implementation was made under the open source license to benefit to a large community.

The quantification of uncertainty is the capacity to measure the effect of the uncertainty in numerical models and simulations. In the industry, this subject is particularly important during the early phases of design and certifications. Although the topic is well defined and mastered in the academic community, the application to industrial problems raises some scalability issues. The majority of the methods are based on many query techniques which imply to perform a considerable number of parametric simulations to collect and analyze the data. This kind of procedure leads to scalability issue when we are working on large dimension cases. To tackle that scalability issue there exist two main solutions: the exploitation of the growing HPC (High-Performance Computing) capabilities and the development of efficient and accurate model order reduction methods. Those two topics are highly interconnected and have been widely studied in the context of this ANR project.



The project gathered different industrial and academic research teams to work on these subjects. Our goal was to improve the robustness, the scalability and the reliability of different reduction modeling techniques used for uncertainty quantification:

- the low rank and tensor reduction
- the Multi fidelity kriging
- the PGD
- the Reduced Basis
- the Bayesian calibration

In parallel of the research efforts, industrials proposed different concrete applications. The idea was to confront the academic work to the resolution of complex test-cases.

## Airbus Use Case

Airbus Group was one of the leading company in the project. They proposed different subjects of research and test-cases. In particular, they were interested in the development of a reliable reduced model for aerothermal simulations. This research topic has a double objective. First, it makes conceivable the use of numerical simulations in the design and certification phases. The actual full order aerothermal models require huge computation capabilities. With that constraint, it is not practicable to perform a large number of simulations, necessary for the uncertainty quantification or the design optimization. In a second time, Airbus Group was interested in the interface between their environment conditioning system (ECS) and a reduced aerothermal model. This ECS is a numerical model, based on the resolution of multiple ordinary differential equations (ODE). Among other things, it control the pressure and the quality of the air in the aircraft. The integration of a reduced high fidelity 3D model in the ECS would highly increase its possibilities. However, it requires real-time embedded simulations which are impossible with a full order model but achievable with a reduced one.

Airbus proposed to work on an aerothermal flow problem in an airplane. The geometry is decomposed into two sub-domains: the cabin on the top with passengers and the bay on the bottom with various electric equipment. A cooling system controls the temperature in both areas, see figure 1

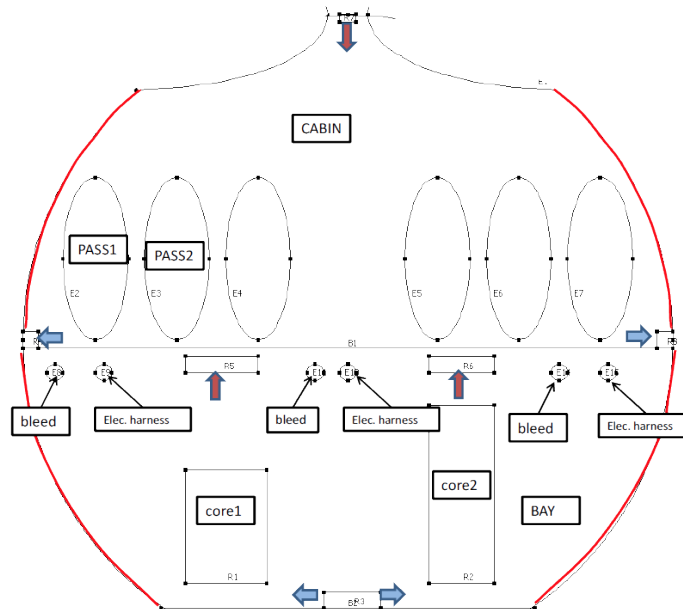


Figure 1 – Description of the Cabin Use Case Proposed by Airbus Group

The development of a reliable reduced order model for this test case was the starting point of this thesis.

## Modeling Aerothermal Flows

Before thinking about the reduction of the model, we had to implement a reliable *truth* approximation. That kind of simulations implies to model the fluid flow and the heat transfers in a complex coupled system. The temperature of the fluid is subject to two phenomena: the usual diffusion and the convection through the effect of the velocity of the fluid. These heat transfers can be modeled by a classic advection-diffusion equation. In addition, we have to determine the fluid behavior, governed by the laws of the fluid dynamics.

The Navier-Stokes equations are the key to computational fluid dynamics (CFD). This non-linear system provides a mathematical model for the velocity and pressure field of a Newtonian fluid. The original equations are designed for compressible fluids. In our case, we will only consider flow with reasonably low velocity, and we can use the incompressibility assumption to simplify the equations. The effect of the temperature of the fluid will be model using the Boussinesq approximation to express the density of the fluid as a function of its temperature. Even with a

simplified version of the Navier-Stokes equations, there exists no known analytical solutions to this problem. Thus the numerical resolution is the only actual way to perform CFD[79]. Many approaches exist for the numerical treatment of these equations: finite volume, finite element, finite difference, spectral element, boundary element,... All those methods are based on a discretization of the spatial domain and a resolution of the equations projected on the discrete space.

In the context of this thesis, we chose to work with the finite element method (FEM). This technique is particularly adapted for the resolution of algebraic equations on complex domains. It can be performed on non-structured discretizations and provide accurate approximations with robust convergence theories [17]. We chose this method for its reliability and its convenient formulation. We were also motivated in our choice by the use of the reduced basis method which was initially designed for the finite element formulation[80].

However, in our case, the physical parameters of the air imply almost negligible diffusion compared to the convection effects. Using the FEM for the resolution of that kind of convection dominated problems, usually produces non-physical oscillations. Thus, the finite element community developed various solutions to tackle this annoying instability issue. In this thesis, we chose to focus on the streamline diffusion methods [46, 7], which allows stabilizing the problem by adding numerical diffusivity in the direction of the convection field. These techniques were widely studied to find the most appropriate formulation for our model.

Those stabilization methods allow recovering a numerically stable solution. Nevertheless, when the Reynolds number is growing, these methods do not take care of the energy dissipated in the smallest scales. To capture these effects, the naive solution would be to use ridiculously thin discretization[57]. This solution is called Direct Numerical Simulation and is not achievable for industrial cases. The alternative is the use of turbulence models to compute an additional turbulent viscosity, modeling the loss of energy in the smallest scales. There exist two main families of turbulence models. The large eddy simulation (LES) is used for transient simulations. The Reynolds average Navier-Stokes (RANS) is adapted for stationary problems, and that is why decide to implement two RANS models: the Spalart-Allmaras and the  $k - \omega$  SST [41].

The cited methods allow performing fluid and aerothermal simulations, including for high Reynolds or turbulent flows. The next step is now to build a reduced version of this model.

## **Model Order Reduction**

The model order reduction methods are particularly adapted in the context of parametrized Partial Differential Equations (PDE). The physic and the mesh are always the same, but the parameters of the equations may vary (boundary

conditions, geometrical parametrization, source terms, ... ). It is interesting to develop a reduced model when it becomes necessary to solve a parametrized PDE multiple times, with different sets of parameters. The reduced models are usually built from an expensive learning procedure to gather information on the model. Then, it exploits the invariant structures of the problem to compute a new reduced order approximation.

We chose to focus our research on the Reduced Basis Method (RBM), one of the principal model order reduction methods, which presents some interesting advantages [80, 103, 81, 83, 88]:

- Rigorous a posteriori error estimators
- An efficient Offline/Online decomposition.

The idea of the reduced basis method is to build a small (relatively to the finite element dimension) dimension space. Then the resolution of the PDE in this low dimension space is very fast and produce accurate approximation as long as the reduced space is rich enough. There exist different algorithms to build this reduced space and the methods are usually problem-dependant. In the context of aerothermal simulations, the constructions of an efficient reduced model also require some advanced techniques to deal with non-linearity: the Empirical Interpolation Method (EIM)[3, 65, 23, 37] will be an essential tool for our reduced basis framework.

### **Feel++: Finite Element Embedded Library**

Feel++ is a finite element embedded library in C++. It offers a very convenient language designed for the resolution of PDE, relying on a Domain Specific Language (DSEL). The goal of Feel++ is to provide a language close to the mathematical one to solve complex PDE. The idea is to provide to scientists a framework in which they can express in a language close to the mathematics the strategy they propose for solving complex systems of PDE and generate a high performance code. To this goal, Feel++ uses the last standards of C++ and meta-programming (through the template and boost meta-programming library) to have a maximum of generality. The library also offers seamless parallel computing. This High Performance Computing capability is essential for the treatment of concrete industrial problems.



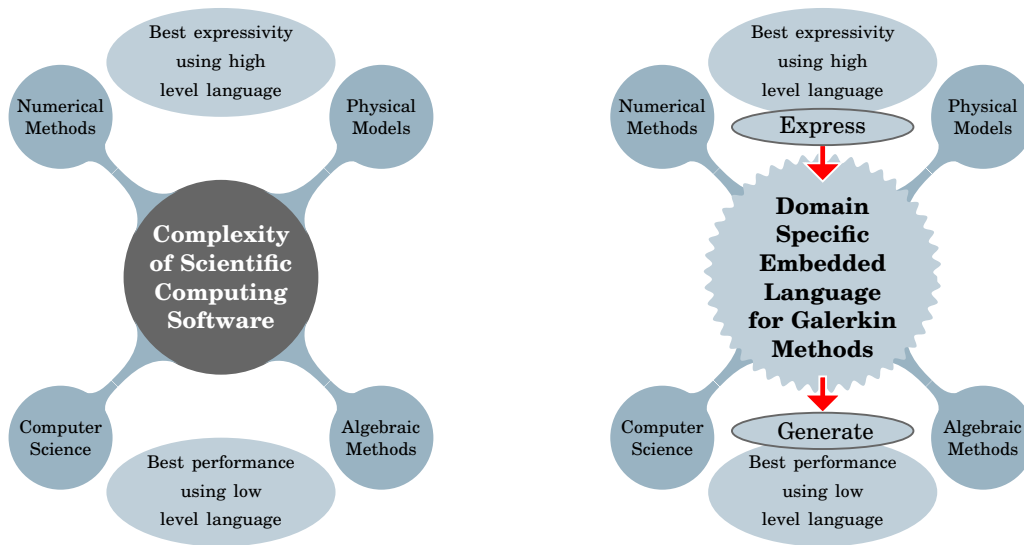


Figure 2 – The DSEL offered by Feel++ provides high level language to break the complexity of scientific computing software while keeping the performances of a low level language.

Feel++ is an open-source library developed under the LGPL and GPL License. It is used in numerous scientific projects funded by different institutions such as the French National Research Agency or the European Commission.

During this thesis, all the applications have been developed using Feel++. We also widely contributed to the reduced basis framework of the library.

## Plan

Just after this introduction, the chapter 1 introduces useful mathematical notions.

Then, the manuscript is organized into three parts.

The part I is dedicated to the aerothermal simulation using the finite element method.

In chapter 2 we introduce the governing equations, we present our discretization method and the numerical strategy implemented for the resolution of the non-linear system. The last section of the chapter is dedicated to the validation of our implementation with numerical results on different benchmarks.

In chapter 3, we detail the stabilization method studied during this thesis. We remind the different formulation of the SDM such as the Streamline Upwind/Petrov-Galerkin (SUPG) or the Galerkin Least Square (GLS). Then we discuss the design of the stabilization parameter which significantly impact the efficiency of the method. As a complement of the SDM, we also compare the different discontinuity capturing method, essential for high order stabilization. The efficiency of the methods introduced in this chapter will be discussed with various numerical applications.

In chapter 4, we propose a very brief review of the existing turbulent models,

and we focus on the RANS method. We detail two models in particular: the Spalart-Allmaras and  $k-\omega$  SST. The implementation of these models is validated by numerical results on the well-known NACA0012 benchmark proposed on the NASA turbulence modeling resource website.

The part II presents our work on the model order reduction. In chapter 5, we remind the principle of the RBM for coercive and saddle-point elliptic problems. We give an example of an application with a reduced model of Stokes flow in a parametrized geometry.

In chapter 6, we present the Empirical Interpolation Method (EIM) and a new implementation of this method developed during this thesis. This new development allows constructing an affine approximation for discrete operators. We also detail the interface of this version of EIM with the Simultaneous EIM and RB (SER) algorithm allowing a decrease of the computational cost. The implementation and the efficiency of the method are discussed with the resolution of a standard benchmark in the last section.

In chapter 7, we finally propose our solution for the reduction of an aerothermal model. We detail the treatment of the quadratic non-linearity, the reduction of the stabilization operators and the construction of the reduced spaces. The method is illustrated with two applications. The first one is a simulation of the heat transfers in the neighboring of an electronic component, refreshed by a cooling system. The second one is a simplified version of the use-case proposed by Airbus Group; we proposed a reduced model of the aerothermal flow into the airplane cabin.

The last part III is an overview of the implementation work realized during this thesis. In chapter 8, we present the design and the different features implemented in our aerothermal library. These tools and algorithms are now progressively integrated into the Feel++ toolboxes.

In chapter 9, we give an overview of the developments in Feel++ RB framework. The main contributions are the implementation of the reduced basis by block for multi-physics problem and the reduced basis for aerothermal problems.

In chapter 10, we finally detail the implementation of the EIM for discrete operators. This development has been entirely realized during this thesis and involved interesting methods of parallel and metaprogramming.

## Publications

in preparation | *Review and Implementation of Streamline Diffusion Methods on Anisotropic Meshes, Application to Aerothermal Simulations*  
with C. Prud'homme

in preparation | *Implementation of RANS models in Feel++ library*  
with C. Prud'homme, Y. Hoarau, V. Chabannes

in preparation | *Simultaneous EIM and RB Construction for Non-linear Operators, Application to Aerothermal Problems*  
with C. Prud'homme, V. Chabannes

## Oral Presentations

2018 | **MoRePaS**, Nantes, France  
*High Reynolds Aerothermal Simulations and Reduced Basis*

2017 | **Feel++ Users Days**, Strasbourg, France  
*Aerothermal Simulation and Model Order Reduction, Using the Open-Source Framework Feel++*

2018 | **ANR CHORUS Workshop**, Paris, France  
*Model Order Reduction for Multi-physics Problems, Using the Open-Source Framework Feel++*

2017 | **A3F**, Strasbourg, France  
*Aerothermal Simulation and Model Order Reduction, using the Open-Source Framework Feel++*

2015 | **SimRace**, Paris, France  
*Solving Strategy for Large Scale Aerothermal Simulation using the Open-Source Framework Feel++*

2015 | **Feel++ Users Days**, Strasbourg, France  
*Aerothermal Simulations : Towards Reduced Basis Applications*

2014 | **Research Group on Model Order Reduction**, Porquerolles, France

## Posters

2018 | **CANUM**, Cap d'Agde, France  
Poster: *High Reynolds Aerothermal Simulations and Reduced Basis*

2015 | **Journée Poster de l'école doctorale**, Strasbourg, France  
Poster: *La Méthode des Bases Réduites Appliquée à des Simulations d'Aérothermie*

---

**Contents**

---

<b>1</b>	<b>Function Spaces . . . . .</b>	<b>9</b>
1.1	Lebesgue Space . . . . .	10
1.2	Sobolev Space . . . . .	10
1.3	Dual Space . . . . .	11
<b>2</b>	<b>The Finite Element Method . . . . .</b>	<b>12</b>
2.1	Variational Formulation . . . . .	12
2.2	Finite Element Definition . . . . .	14
2.3	Galerkin Projection . . . . .	15

---

This first chapter is dedicated to the definition of some mathematical notions and notations, relevant to the context of this document.

## 1 Function Spaces

In the following sections, and more generally in this thesis report,  $\Omega$  is supposed to be a bounded domain of  $\mathbb{R}^d$ ,  $d = 1, 2, 3$ , with a Lipschitz boundary  $\partial\Omega = \Gamma$ .

## 1.1 Lebesgue Space

Lebesgue spaces,  $\mathcal{L}^p(\Omega)$ , are spaces of functions on  $\Omega$  for which the  $p$ -th power of its absolute value is Lebesgue integrable,

$$\mathcal{L}^p(\Omega) = \{v : \Omega \rightarrow \mathbb{R}, \int_{\Omega} |v|^p d\Omega < \infty\} \quad (1.1)$$

$\mathcal{L}^p(\Omega)$  is a Banach space with the usual  $p$ -norm  $\|\cdot\|_p$ ,

$$\|v\|_p = \left( \int_{\Omega} |v|^p d\Omega \right)^{1/p}. \quad (1.2)$$

We can extend this notion to the particular case of  $p = \infty$  using the essential supremum of the function  $v$  as  $\infty$ -norm,

$$\mathcal{L}^{\infty}(\Omega) = \{v : \Omega \rightarrow \mathbb{R}^n, \|v\|_{\infty} < \infty\}, \quad \|v\|_{\infty} = \operatorname{ess\,sup}_{x \in \Omega} |v(x)|, \quad (1.3)$$

where the essential supremum of a function is defined as the supremum of the function *almost everywhere*, i.e. except on a set of measure zero. This notion is an adaptation of the usual global supremum to measure theory.

In the following report, we essentially use the common  $\mathcal{L}^2(\Omega)$  space of square-integrable functions,

$$\mathcal{L}^2(\Omega) = \{v : \Omega \rightarrow \mathbb{R}^n, \|v\|_2 < \infty\}, \quad \|v\|_2 = \sqrt{\int_{\Omega} |v|^2 d\Omega}. \quad (1.4)$$

$\mathcal{L}^2$  is particularly interesting since he has an Hilbert structure, with its usual scalar product

$$(u, v)_0 = \int_{\Omega} u \cdot v d\Omega. \quad (1.5)$$

## 1.2 Sobolev Space

A Sobolev space  $W^{k,p}(\Omega)$  is a subspace of  $\mathcal{L}^p(\Omega)$  such that for any  $v \in W^{k,p}(\Omega)$  all the weak derivatives of  $v$ , up to the order  $k$ , are in  $\mathcal{L}^p(\Omega)$ ,

$$W^{k,p}(\Omega) = \{v \in \mathcal{L}^p(\Omega), D^{\alpha}v \in \mathcal{L}^p(\Omega), |\alpha| \leq k\}, \quad (1.6)$$

where  $D^{\alpha}v$  is the generalised weak derivative operator, defined for any multi-index  $\alpha = (\alpha_1, \dots, \alpha_d)$ , such that  $\alpha_i \geq 0$ ,  $1 \leq i \leq d$ , and for any  $x = (x_1, \dots, x_d) \in \Omega$  as

$$D^{\alpha} = \frac{\partial^{|\alpha|}}{\partial x_1^{\alpha_1} \dots \partial x_d^{\alpha_d}}, \quad |\alpha| = \sum_{i=1}^d \alpha_i. \quad (1.7)$$

Sobolev spaces have Banach structure with the norm

$$\|v\|_{k,p} = \left( \sum_{|\alpha| \leq k} \|D^\alpha v\|_p^p \right)^{1/p}, \quad 1 \leq p \leq +\infty \quad (1.8)$$

$$\|v\|_{k,\infty} = \max_{|\alpha| \leq k} \|D^\alpha v\|_\infty.$$

In the context of parametrized differential equations, the Sobolev spaces with  $p = 2$  are particularly interesting since they have an Hilbert structure, with the scalar product

$$(u, v)_k = \int_{\Omega} u \cdot v \, d\Omega + \sum_{0 < |\alpha| \leq k} \int_{\Omega} D^\alpha u \cdot D^\alpha v \, d\Omega = (u, v)_0 + \sum_{0 < |\alpha| \leq k} \int_{\Omega} D^\alpha u \cdot D^\alpha v \, d\Omega, \quad (1.9)$$

and the associated norm

$$\|v\|_k = \sqrt{(v, v)_k}. \quad (1.10)$$

In this case, we use the notation  $H^k(\Omega) = W^{k,2}(\Omega)$ ,

$$H^k(\Omega) = \{v \in \mathcal{L}^2(\Omega), D^\alpha v \in \mathcal{L}^2(\Omega), |\alpha| \leq k\}. \quad (1.11)$$

### 1.3 Dual Space

Let  $X$  be an Hilbert space with scalar product  $(\cdot, \cdot)_X$  and associated norm  $\|\cdot\|_X$ . The dual spaces  $X'$  of  $X$  is the set of continuous linear forms on  $X$ ,

$$X' = \{l : X \rightarrow \mathbb{R}, l \text{ continuous linear}\}. \quad (1.12)$$

For any  $l \in X'$ , we define its dual norm by

$$\|l\|_{X'} = \sup_{v \in X} \frac{|l(v)|}{\|v\|_X} \quad (1.13)$$

The Riesz representation theorem establishes that  $X$  and  $X'$  are isometrically isomorphic. A direct consequence of this theorem is the existence, for any  $l \in X'$ , of a Riesz representation  $\hat{l} \in X$  such that

$$(\hat{l}, v)_X = l(v), \quad \forall v \in X. \quad (1.14)$$

This Riesz representation is very convenient to evaluate the dual norm of a linear form  $l$ , using the isometrical relation between  $X$  and  $X'$ ,

$$\|\hat{l}\|_X = \|l\|_{X'} = \sup_{v \in X} \frac{|l(v)|}{\|v\|_X}. \quad (1.15)$$

This relation is widely used in the reduced basis methodology, see chapter 5.

## 2 The Finite Element Method

### 2.1 Variational Formulation

In a weak formulation, an equation is no longer required to hold absolutely as in its strong formulation. Instead, the weak formulation requires a solution in the sense of distributions and might not be defined on a set of measure zero. This formulation is crucial in the context of functional analysis and especially for treatment of PDEs.

In the present report, we assume that all studied PDE admit a weak formulation, written as find  $u \in X$  such that

$$a(u, v) = f(v), \quad \forall v \in X, \quad (1.16)$$

where  $X$  is an Hilbert space on  $\Omega$ ,  $a : X \times X \rightarrow \mathbb{R}$  is a bilinear form and  $f : X \rightarrow \mathbb{R}$  is a linear form.  $u$  is called the trial function and lives in the trial space  $X$  while  $v$  are the test functions living in the test space  $X$ .

We now remind two fundamental theorems that ensure the well-posedness of the examined problems. In the rest of this report, when we describe a problem and its weak formulation, we assume that it fulfills the conditions of one of these theorems - depending on the nature of the problem.

**Theorem 1 (Lax-Milgram).**

*We consider the problem (1.16) and we furthermore assume that*

(a) *the linear form  $l$  is continuous,*

$$\gamma_l = \sup_{v \in X} \frac{|l(v)|}{\|v\|_X} < \infty \quad (1.17)$$

(b) *the bilinear form  $a$  is continuous,*

$$\gamma_a = \sup_{u \in X} \sup_{v \in X} \frac{|a(u, v)|}{\|u\|_X \|v\|_X} < \infty \quad (1.18)$$

(c) *the bilinear form  $a$  is coercive,*

$$\alpha_a = \inf_{v \in X} \frac{a(v, v)}{\|v\|_X^2} > 0. \quad (1.19)$$

*Under these assumptions, there exists an unique solution  $u \in X$  to the weak formulation (1.16).*

*Proof.* See [60].

□

**Remark 1.** *The proof of this theorem comes from the calculus of variations; that is why the term of variational formulation is often used to refer to the weak formulation.*

**Theorem 2 (Brezzi).**

*In the context of fluid simulation we will consider saddle point problems expressed under the variational formulation: find  $(\mathbf{u}, p) \in V \times Q$  such that*

$$\begin{cases} a(\mathbf{u}, \mathbf{v}) + b(\mathbf{v}, p) = f(\mathbf{v}), & \forall \mathbf{v} \in V, \\ b(\mathbf{u}, q) = g(q), & \forall q \in Q, \end{cases} \quad (1.20)$$

*where  $V$  and  $Q$  are two Hilbert spaces with associated scalar products  $(\cdot, \cdot)_V$ ,  $(\cdot, \cdot)_Q$  and norms  $\|\cdot\|_V$ ,  $\|\cdot\|_Q$ .  $Z$  is the composite space  $Z = V \times Q$  with inner product  $(\cdot, \cdot)_Z = (\cdot, \cdot)_V + (\cdot, \cdot)_Q$  and norm  $\|\cdot\|_Z = \sqrt{(\cdot, \cdot)_Z}$ .  $a(\cdot, \cdot) : V \times V \rightarrow \mathbb{R}$ ,  $b(\cdot, \cdot) : V \times Q \rightarrow \mathbb{R}$  and  $f(\cdot) : V \rightarrow \mathbb{R}$ ,  $g(\cdot) : Q \rightarrow \mathbb{R}$  are bilinear and linear bounded forms respectively.*

*We furthermore assume that :*

(a)  *$a$  and  $b$  are continuous*

$$\gamma_a = \sup_{\mathbf{u} \in V} \sup_{\mathbf{v} \in V} \frac{|a(\mathbf{u}, \mathbf{v})|}{\|\mathbf{u}\|_V \|\mathbf{v}\|_V} < \infty \quad (1.21)$$

$$\gamma_b = \sup_{\mathbf{v} \in V} \sup_{q \in Q} \frac{|b(\mathbf{v}, q)|}{\|\mathbf{v}\|_V \|q\|_Q} < \infty, \quad (1.22)$$

(b)  *$a$  is coercive on  $V$*

$$\alpha_a = \inf_{\mathbf{v} \in V} \frac{a(\mathbf{v}, \mathbf{v})}{\|\mathbf{v}\|_V^2} > 0, \quad (1.23)$$

(c)  *$b$  satisfies the inf-sup condition (or Ladyzhenskaya–Babuška–Brezzi condition), expressed either using the Brezzi definition,*

$$\beta_{Br} = \inf_{q \in Q} \sup_{\mathbf{v} \in V} \frac{b(\mathbf{v}, q)}{\|\mathbf{v}\|_V \|q\|_Q} > 0, \quad (1.24)$$

*or the Babuška definition,*

$$\beta_{Ba} = \inf_{(\mathbf{u}, p) \in Z} \sup_{(\mathbf{v}, q) \in Z} \frac{a(\mathbf{u}, \mathbf{v}) + b(\mathbf{v}, p) + b(\mathbf{u}, q)}{\|(\mathbf{u}, p)\|_Z \|(\mathbf{v}, q)\|_Z} > 0. \quad (1.25)$$

*Under these assumptions, there exists a unique couple of solutions  $(\mathbf{u}, p) \in Z$  to the problem (1.20).*

*Proof.* See [5]. □



## 2.2 Finite Element Definition

Using the formalism introduced in [17], a finite element is defined by a tuple  $(K, P, \Sigma)$ .  $K$  is a compact, connected and non-empty geometrical element with a Lipschitz boundary - typically a simplex or a hypercube.  $P$  is a finite dimensional function space on  $K$ .  $\Sigma$  is a set of  $N_l$  linear forms  $\{\sigma_i : P \rightarrow \mathbb{R}\}_{i=1}^{N_l}$  such that

$$\begin{aligned} P &\rightarrow \mathbb{R}^{N_l} \\ p &\rightarrow (\sigma_1(p), \dots, \sigma_{N_l}(p)) \end{aligned} \quad (1.26)$$

is bijective. Those linear forms are called the degrees of freedom (DoF) of the finite element.

We find many kinds of finite element in the literature: Lagrange, Hermite, Raviart-Thomas, Crouzeix-Raviart, Nedelec,... Lagrange finite element is probably the most common type and will be used in this work. For more details on the finite elements available in the Feel++ library, we refer the reader to [15, 24].

### Lagrange Finite Element

We denote by  $\mathbb{P}^k(K)$  the set of polynomial functions on  $K$  of maximum degree  $k$ . We introduce the interpolation points  $\{t_i\}_{i=1}^{N_l}$  defined on the element  $K$ . The linear forms in  $\Sigma$  are then defined as

$$\begin{aligned} \sigma_i : \mathbb{P}^k(K) &\rightarrow \mathbb{R} \\ p &\rightarrow p(t_i) \end{aligned} \quad (1.27)$$

A Lagrange finite element of order  $k$  is defined by the tuple  $(K, \mathbb{P}^k(K), \Sigma)$  and the local basis functions  $\{\phi_i\}_{i=1}^{N_l}$  are set such that

$$\sigma_i(\phi_j) = \delta_{ij}, \quad \forall 1 \leq i, j \leq N_l, \quad (1.28)$$

where  $\delta_{ij}$  is the Kronecker symbol.

### Lagrange Finite Element Space

We introduce now a conforming partition of the domain  $\Omega$

$$\mathcal{T}_h = \{K_m\}_{m=1}^{N_e} \quad (1.29)$$

where all the elements  $K_i$  are of same nature - simplex or hypercubes. For any Hilbert space  $X(\Omega)$  on  $\Omega$ , we can define its discrete representation on the mesh  $\mathcal{T}_h$ , as the continuous Lagrange finite element space of order  $k$

$$X_h(\Omega) \equiv \mathbb{L}_{ch}^k(\Omega) = \{v \in X(\Omega) \cap C^0(\Omega), v|_{K_m} \in \mathbb{P}^k(K_m)\}. \quad (1.30)$$

The subscript  $c$  stands for the continuous property which is actually not mandatory, depending on the projection method. Hereafter we will only consider the Galerkin projection with continuous Lagrange space  $\mathbb{L}_{ch}^k$ . The subscript  $h$  stands for the characteristic length of the mesh  $\mathcal{T}_h$ . We usually denote by  $\mathcal{N}$  the dimension of the FE space  $X_h(\Omega)$  and use the convenient notation  $X_h \equiv X_h(\Omega)$  when it is not ambiguous.

The FE space  $X_h$  is an Hilbert space with scalar product  $(\cdot, \cdot)_X$  and associated norm  $\|\cdot\|_X$  inherited from  $X$ .

**Remark 2.** *It was convenient in the previous definition to make the assumption  $\Omega_h = \Omega$ , with  $\Omega_h = \cup_{K \in \mathcal{T}_h} K$ . In practice, we only have an approximation  $\Omega_h \approx \Omega$  and the FE element space  $X_h$  is not defined on  $\Omega$  but only on the approximated domain  $\Omega_h$ .*

### 2.3 Galerkin Projection

We use the Galerkin method as internal approximation method. In order to numerically solve the problem (1.16), we introduce a FE subspace  $X_h \subset X$ . The exact solution  $u$  of (1.16) will be approximated by the discrete solution  $u_h \in X_h$  such that

$$a(u_h, v_h) = f(v_h), \quad \forall v_h \in X_h. \quad (1.31)$$

Since the solution  $u_h \in X_h = \text{span}\{\phi_i, 1 \leq i \leq \mathcal{N}\}$ , we have,

$$u_h = \sum_{j=1}^{\mathcal{N}} w_h^j \phi_j. \quad (1.32)$$

The variational formulation can then be rewritten as find  $(u_h^1, \dots, u_h^{\mathcal{N}})$  such that

$$\sum_{i=1}^{\mathcal{N}} u_h^j a(\phi_j, \phi_i) = f(\phi_i), \quad \forall 1 \leq i \leq \mathcal{N}. \quad (1.33)$$

And finally, introducing the matrix  $\mathbf{A}_h$ ,  $(\mathbf{A}_h)_{ij} = a(\phi_j, \phi_i)$  and the vector  $\mathbf{F}_h$ ,  $(\mathbf{F}_h)_i = f(\phi_i)$ , the Galerkin projection  $u_h$  of the continuous solution  $u$  is the solution of the linear system

$$\mathbf{A}_h u_h = \mathbf{F}_h \quad (1.34)$$

#### Well-Posedness of the Discrete Problem

Assuming that the problem (1.16) is well-posed in the sense of the Lax-Milgram theorem 1, we have the following relations:

$$\gamma \geq \gamma_h = \sup_{v_h \in X_h} \frac{|l(v_h)|}{\|v_h\|_X}, \quad (1.35)$$

$$\gamma_a \geq \gamma_{ah} = \sup_{u_h \in X_h} \sup_{v_h \in X_h} \frac{|a(u_h, v_h)|}{\|u_h\|_X \|v_h\|_X}, \quad (1.36)$$

$$\alpha_a \leq \alpha_{ah} = \inf_{v_h \in X_h} \frac{a(v_h, v_h)}{\|v_h\|_X^2}. \quad (1.37)$$

The properties of continuity and coercivity are then inherited from the continuous problem, and it guarantees the well-posedness of discrete problem (1.34), [60].

In the case of saddle point problem, the property of continuity is also inherited from the continuous problem. However, the inf-sup condition is not automatically satisfied for any FE space  $V_h \times Q_h$ . A very convenient solution to satisfy this necessary condition is to choose specific finite element space. For instance, in the context of this thesis, we use Taylor-Hood spaces of degree  $k$ ,

$$\text{TH}_{ch}^k = [\mathbb{L}_{ch}^{k+1}]^d \times \mathbb{L}_{ch}^k. \quad (1.38)$$

This discretization ensures the inf-sup condition to be satisfied [5], and then the well-posedness of the discrete problem is ensured by the Brezzi theorem 2.

## Convergence of the Method

The finite element methodology provides bounds on the error between the exact solution  $u$  and the FE approximation  $u_h$ , [17]. This error depends on the following parameters:

- the polynomial order  $k$  of the FE space,
- the nature of the considered norm: usually the  $H^m$ -norm  $\|\cdot\|_m$ , as defined in section 1.2. The case  $m = 0$  stands for the  $\mathcal{L}^2$ -norm,
- the characteristic length  $h$  of the mesh  $\mathcal{T}_h$ ,
- the regularity of the solution.

Considering these elements, we have the existence of a constant  $C \in \mathbb{R}$ ,  $C > 0$  such that

$$\|u - u_h\|_m \leq Ch^{k+1-m}. \quad (1.39)$$

This bound ensures the convergence of the FE solution  $u_h$  when the characteristic length of the mesh decrease,

$$u_h \xrightarrow{h \rightarrow 0} u. \quad (1.40)$$

**Part I**

**Aerothermal Simulations**



---

Modeling of an Aerothermal Problem

---

**Contents**


---

<b>1</b>	<b>The Navier Stokes Equations and the Boussinesq Approximation . . . . .</b>	<b>20</b>
<b>2</b>	<b>Finite Element Approach . . . . .</b>	<b>23</b>
2.1	Variational Formulation . . . . .	23
2.2	Discretization . . . . .	24
<b>3</b>	<b>Solution Strategy . . . . .</b>	<b>26</b>
3.1	Non-Linear Iterative Methods . . . . .	28
3.2	Pseudo Transient Continuation . . . . .	30
<b>4</b>	<b>Validation . . . . .</b>	<b>32</b>
4.1	Validation of the Fluid Solver . . . . .	33
4.2	Validation of the Slip Boundary Condition Implementation .	37
4.3	Validation of the Aerothermal Solver . . . . .	41

---

In this first chapter we describe the full resolution process, from the continuous level to the computation of a discrete solution. We start with the formulation of the governing equations arising from different physical laws and approximations. In the second section, we detail the discretization of the continuous problem, using the FEM. And finally, we present our resolution strategy for the non-linear system, in section 3. The last section of this chapter is dedicated to numerical results and especially to the validation of our implementation.

# 1 The Navier Stokes Equations and the Boussinesq Approximation

In the context of this thesis, the fluid is always supposed Newtonian and incompressible. Besides, we use the Boussinesq approximation for modeling the buoyancy force. With these hypotheses, it only remains three quantities of interest in the aerothermal system: the velocity ( $\mathbf{u}$ ), the pressure ( $p$ ) and the temperature ( $T$ ).

In this first section, we briefly remind the laws of conservation involved in the establishment of the Navier-Stokes equations [79]. These equations mainly derive from the principles of continuity of mass, momentum, and energy. Those physical rules combined with the Boussinesq approximation lead to a coupled system between the three quantities,  $\mathbf{u}$ ,  $p$ , and  $T$ .

## Derivation of continuity equation

This equation arises from the fundamental law of Newton mechanics that states the conservation of mass ( $M$ ) in an arbitrary control volume ( $V_m$ ) which leads to

$$0 = \frac{dM}{dt} = \frac{d}{dt} \int_{V_m} \rho dV = \int_{V_m} \frac{\partial \rho}{\partial t} dV + \int_{S_m} \rho \mathbf{u} \cdot \mathbf{n} dS = \int_{V_m} \left( \frac{\partial \rho}{\partial t} + \nabla \cdot (\rho \mathbf{u}) \right) dV, \quad (2.1)$$

where  $S_m$  is the boundary surface of the volume  $V_m$ ,  $\mathbf{n}$  is the outward normal vector to the surface and  $\rho$  is the density of the fluid. Since the volume  $V_m$  is arbitrary chosen, the integrand has to be zero and leads to the continuity equation,

$$\frac{\partial \rho}{\partial t} + \nabla \cdot (\rho \mathbf{u}) = 0. \quad (2.2)$$

In the particular case of incompressible fluids, the density ( $\rho$ ) is considered constant, this leads to the simplified continuity equation,

$$\boxed{\nabla \cdot \mathbf{u} = 0.} \quad (2.3)$$

## Derivation of Momentum Equation

One way to obtain this second equation is the application of the second Newton's law to a fluid volume element. It leads to the first form of the momentum equation,

$$\rho \left( \frac{\partial \mathbf{u}}{\partial t} + \mathbf{u} \cdot \nabla \mathbf{u} \right) = \mathbf{s}, \quad (2.4)$$

where  $\mathbf{s}$  is the momentum source. This term is composed of two kinds of forces: the surface forces act only on the surface of a control volume whereas body forces act on the whole control volume. Thus, by splitting  $\mathbf{s}$ , we obtain the Cauchy equation

$$\rho \left( \frac{\partial \mathbf{u}}{\partial t} + \mathbf{u} \cdot \nabla \mathbf{u} \right) = \nabla \cdot \mathbf{T} + \mathbf{f} \quad (2.5)$$

where  $\mathbf{f}$  represents the sum of different body forces applied on the fluid and  $\mathbf{T}$  is the Cauchy stress tensor, used to express the surface forces acting on the fluid. This equation can again be rewritten by introducing the viscous stress tensor  $\mathcal{T}$  and the pressure field  $p$ ,

$$\rho\left(\frac{\partial \mathbf{u}}{\partial t} + \mathbf{u} \cdot \nabla \mathbf{u}\right) = -\nabla p + \nabla \cdot \mathcal{T} + \mathbf{f}. \quad (2.6)$$

A fluid is defined as Newtonian [4] when the following assumptions are verified:

- The viscous stress tensor is a linear function of the strain rate tensor,  $\mathcal{T} = 2\mu\mathbf{S}$ , where  $\mathbf{S}(\mathbf{u}) = \frac{1}{2}(\nabla \mathbf{u} + \nabla \mathbf{u}^\top)$ ,
- The fluid is isotropic,
- $\nabla \cdot \mathcal{T} = \mathbf{0}$  at rest.

And by adding the incompressible hypothesis (2.3), we obtain

$$\boxed{\rho\left(\frac{\partial \mathbf{u}}{\partial t} + \mathbf{u} \cdot \nabla \mathbf{u}\right) = -\nabla p + 2\nabla \cdot (\mu\mathbf{S}(\mathbf{u})) + \mathbf{f}.} \quad (2.7)$$

Henceforth, the only body force considered is the gravity:  $\mathbf{f} = \rho\mathbf{g}$ , where  $\mathbf{g}$  is the gravitational field.

### The Boussinesq Approximation

The Boussinesq approximation states that density differences are sufficiently small to be neglected, except where they appear in terms multiplied by  $\mathbf{g}$ . So we can assume that  $\rho = \rho_0$  except in the buoyancy force which is the only force acting on the fluid (in the considered model).

The Boussinesq approximation also links the density variation to the temperature field with the relation

$$\rho - \rho_0 = -\rho_0\beta(T - T_0), \quad (2.8)$$

where  $\rho_0$  and  $T_0$  are reference density and temperature and  $\beta$  is the coefficient of thermal expansion. We can then rewrite the buoyancy force

$$\mathbf{f} = \rho\mathbf{g} = \rho_0\mathbf{g} - \rho_0\beta(T - T_0)\mathbf{g}. \quad (2.9)$$

This yields to a new momentum equation

$$\rho\left(\frac{\partial \mathbf{u}}{\partial t} + \mathbf{u} \cdot \nabla \mathbf{u}\right) = -\nabla p + 2\nabla \cdot (\mu\mathbf{S}(\mathbf{u})) + \rho_0\mathbf{g} - \rho_0\beta(T - T_0)\mathbf{g}. \quad (2.10)$$



To avoid the evaluation of the term  $\rho_0 \mathbf{g}$ , we use a pressure shift and we define the new pressure  $p^* = p + \rho_0 \mathbf{g} h$ , where  $h$  is the elevation. From here, we will just ignore the \* superscript and denote by  $p$  the shifted pressure.

The Boussinesq approximation also gives a simplified energy equation [98], by considering the different hypothesis exposed and the fact that the heat capacity is constant,

$$\boxed{\frac{\partial T}{\partial t} + \mathbf{u} \cdot \nabla T = \nabla \cdot (\kappa \nabla T) + \frac{J}{\rho C_p}} \quad (2.11)$$

where  $\kappa$  is the thermal diffusivity,  $C_p$  is the heat capacity of the fluid and  $J$  is the rate per unit volume of internal heat production. This last quantity is always supposed to be zero in the present work.

### Strong Formulation

We obtain a non-linear system of three coupled equations

$$\begin{cases} \rho \left( \frac{\partial \mathbf{u}}{\partial t} + \mathbf{u} \cdot \nabla \mathbf{u} \right) + \nabla p - 2 \nabla \cdot (\mu \mathbf{S}(\mathbf{u})) = -\rho \beta (T - T_0) \mathbf{g}, \\ \nabla \cdot \mathbf{u} = 0, \\ \frac{\partial T}{\partial t} + \mathbf{u} \cdot \nabla T - \nabla \cdot (\kappa \nabla T) = 0. \end{cases} \quad (2.12)$$

This system is completed with boundary conditions and initial values for the different fields of interest. These closure conditions will be detailed for each studied model and test case but we can already enumerate some usual boundary conditions and notations.

We consider the problem (2.12) on a domain  $\Omega \in \mathbb{R}^d$ ,  $d = 2, 3$ , with a sufficiently smooth boundary,  $\partial \Omega = \Gamma$ . This boundary is partitioned (in two different ways) to settle the different boundary conditions

$$\Gamma = \Gamma_D^F \cup \Gamma_N^F \cup \Gamma_S^F = \Gamma_D^T \cup \Gamma_N^T \cup \Gamma_R^T, \quad (2.13)$$

where the  $F$  and  $T$  superscript stand for fluid or temperature and the  $D$ ,  $N$ ,  $S$  and  $R$  subscript stand for Dirichlet, Neumann, Slip and Robin conditions. On these subsets of  $\Gamma$ , the boundary conditions can be described as

$$\begin{aligned} \mathbf{u} &= \mathbf{u}_D(\mathbf{x}, t), & \text{on } \Gamma_D^F, & \text{Dirichlet condition on Fluid} \\ \mathbf{T}(\mathbf{u}, p) \mathbf{n} &= \mathbf{0}, & \text{on } \Gamma_N^F, & \text{Neumann condition on Fluid} \\ \mathbf{u} \cdot \mathbf{n} &= 0, \quad (\mathbf{T}(\mathbf{u}, p) \mathbf{n})_\tau = \mathbf{u}_S & \text{on } \Gamma_S^F, & \text{Slip condition on Fluid} \\ T &= T_D(\mathbf{x}, t) & \text{on } \Gamma_D^T, & \text{Dirichlet condition on Temperature} \\ \kappa \nabla T \cdot \mathbf{n} &= \phi_\Gamma(\mathbf{x}, t) & \text{on } \Gamma_N^T, & \text{Neumann condition on Temperature} \\ \kappa \nabla T \cdot \mathbf{n} &= \kappa_\Gamma (T_R(\mathbf{x}, t) - T) & \text{on } \Gamma_R^T, & \text{Robin condition on Temperature} \end{aligned} \quad (2.14)$$

where  $\mathbf{u}_D$ ,  $\mathbf{u}_S$ ,  $T_D$ ,  $F_\Gamma$  and  $\kappa_\Gamma$  are at least  $\mathcal{L}^2$  functions on  $\Gamma$ .

## 2 Finite Element Approach

We now detail the application of the finite element method (FEM) to the aerothermal system (2.12). We are mainly interested in the steady state of the model (when it exists) and then the stationary system reads

$$\begin{cases} \rho \mathbf{u} \cdot \nabla \mathbf{u} + \nabla p - 2\nabla \cdot (\mu \mathbf{S}(\mathbf{u})) = -\rho\beta(T - T_0)\mathbf{g}, \\ \nabla \cdot \mathbf{u} = 0, \\ \mathbf{u} \cdot \nabla T - \nabla \cdot (\kappa \nabla T) = 0, \\ +\text{Boundary Conditions 2.14.} \end{cases} \quad (2.15)$$

### 2.1 Variational Formulation

We start by introducing the Hilbert spaces  $V$ ,  $Q$  and  $X$  with the scalar products  $(\cdot, \cdot)_V$ ,  $(\cdot, \cdot)_Q$ ,  $(\cdot, \cdot)_X$  and associated norms  $\|\cdot\|_V$ ,  $\|\cdot\|_Q$ ,  $\|\cdot\|_X$ , respectively. We also define the composite space  $Z = V \times Q$  with the scalar product  $(\cdot, \cdot)_Z = (\cdot, \cdot)_V + (\cdot, \cdot)_Q$  and associated norm  $\|\cdot\|_Z = \sqrt{\|\cdot\|_V^2 + \|\cdot\|_Q^2}$ .

We apply the scalar product to each row of the system (2.15) with its respective test function  $\mathbf{v} \in V$ ,  $q \in Q$  or  $S \in X$ . Then, by using standard integration by part formulas, we obtain the variational formulation, where the weak solution  $(\mathbf{u}, p, T) \in V \times Q \times X$  is such that

$$\begin{aligned} \int_{\Omega} \rho(\mathbf{u} \cdot \nabla \mathbf{u}) \cdot \mathbf{v} \, d\Omega + \int_{\Omega} 2\mu \mathbf{S}(\mathbf{u}) : \mathbf{S}(\mathbf{v}) \, d\Omega - \int_{\Gamma} \mathbf{T}(\mathbf{u}, p) \mathbf{n} \cdot \mathbf{v} \, d\Gamma \\ - \int_{\Omega} p \nabla \cdot \mathbf{v} \, d\Omega - \alpha \int_{\Omega} (\nabla \cdot \mathbf{u}) q \, d\Omega = - \int_{\Omega} \rho\beta(T - T_0)\mathbf{g} \cdot \mathbf{v} \, d\Omega, \\ \int_{\Omega} (\mathbf{u} \cdot \nabla T) S \, d\Omega + \int_{\Omega} \kappa \nabla T \cdot \nabla S \, d\Omega - \int_{\Gamma} \kappa \nabla T \mathbf{n} \cdot S \, d\Gamma = 0, \end{aligned} \quad \forall (\mathbf{v}, q, S) \in V \times Q \times X, \quad (2.16)$$

where  $\alpha = \pm 1$  determines when we use symmetric or asymmetric formulation.

The terms  $\int_{\Gamma} \mathbf{T}(\mathbf{u}, p) \mathbf{n} \cdot \mathbf{v} \, d\Gamma$  and  $\int_{\Gamma} \kappa \nabla T \cdot \mathbf{n} S \, d\Gamma$  allow to apply the boundary conditions 2.14. The test spaces  $V$  and  $X$  are supposed to verify  $\mathbf{v}|_{\Gamma_D^F} = 0, \forall \mathbf{v} \in V$  and  $T|_{\Gamma_D^T} = 0, \forall T \in X$ . Thus the terms  $\int_{\Gamma} \mathbf{T}(\mathbf{u}, p) \mathbf{n} \cdot \mathbf{v} \, d\Gamma$  and  $\int_{\Gamma} \kappa \nabla T \cdot \mathbf{n} S \, d\Gamma$  are zero on  $\Gamma_D^F$  and  $\Gamma_D^T$ , respectively, and the Dirichlet boundary conditions will be imposed strongly as described later. The Neumann boundary conditions on fluid also set the term  $\int_{\Gamma} \mathbf{T}(\mathbf{u}, p) \mathbf{n} \cdot \mathbf{v} \, d\Gamma$  to zero since the Cauchy stress tensor is supposed to be null. The treatment of the slip boundary condition was not trivial and will be detailed in a further paragraph. It remains the Robin and Neumann boundary conditions on the temperature which leads to

$$\int_{\Gamma} \kappa \nabla T \cdot \mathbf{n} S \, d\Gamma = \int_{\Gamma_N^T} \phi_{\Gamma} S \, d\Gamma + \int_{\Gamma_R^T} \kappa_{\Gamma} (T_R - T) S \, d\Gamma. \quad (2.17)$$

## 2.2 Discretization

In order to solve the problem, (2.16) we introduce the discrete spaces

$$V_h \times Q_h = \mathbb{T}\mathbb{H}_{ch}^k \subset V \times Q, \quad X_h = \mathbb{L}_{ch}^k \subset X, \quad (2.18)$$

see section 2 for a description of these FE spaces. The resolution on these discrete spaces leads to the computation of the FE approximation  $(\mathbf{u}_h, p_h, T_h) \in V_h \times Q_h \times X_h$  such that

$$\begin{aligned} c(\mathbf{u}_h, \mathbf{v}_h; \mathbf{u}_h) + d_u(\mathbf{u}_h, \mathbf{v}_h) + \alpha b(\mathbf{v}_h, p_h) + b(\mathbf{u}_h, q_h) &= f(\mathbf{v}_h; T_h), \\ a(T_h, S_h; \mathbf{u}_h) + d_T(T_h, S_h) &= g(S), \\ \forall (\mathbf{v}_h, q_h, S_h) &\in V_h \times Q_h \times X_h, \end{aligned} \quad (2.19)$$

where,  $\forall (\mathbf{u}, \mathbf{v}, \mathbf{w}) \in V^3, \forall q \in Q, \forall (T, S) \in X^2$ ,

$$\begin{aligned} c(\mathbf{u}, \mathbf{v}; \mathbf{w}) &= \int_{\Omega} \rho(\mathbf{w} \cdot \nabla \mathbf{u}) \cdot \mathbf{v} \, d\Omega, & d_u(\mathbf{u}, \mathbf{v}) &= \int_{\Omega} 2\mu \mathbf{S}(\mathbf{u}) : \mathbf{S}(\mathbf{v}) \, d\Omega, \\ b(\mathbf{u}, q) &= - \int_{\Omega} (\nabla \cdot \mathbf{u}) q \, d\Omega, \\ a(T, S; \mathbf{u}) &= \int_{\Omega} (\mathbf{u} \cdot \nabla T) S \, d\Omega, & d_T(T, S) &= \int_{\Omega} \kappa \nabla T \cdot \nabla S \, d\Omega + \int_{\Gamma_R^T} \kappa_{\Gamma} T S \, d\Gamma \\ f(\mathbf{v}; T) &= \int_{\Omega} \rho \beta (T - T_0) \mathbf{g} \cdot \mathbf{v} \, d\Omega, & g(S) &= \int_{\Gamma_N^T} \phi_{\Gamma} S \, d\Gamma + \int_{\Gamma_R^T} \kappa_{\Gamma} T_R S \, d\Gamma \end{aligned} \quad (2.20)$$

By introducing the basis  $\mathcal{B}_u = \{\zeta_i\}_{i=1}^{\mathcal{N}_u}$ ,  $\mathcal{B}_p = \{\eta_i\}_{i=1}^{\mathcal{N}_p}$  and  $\mathcal{B}_T = \{\xi_i\}_{i=1}^{\mathcal{N}_T}$  of the spaces  $V_h, Q_h$  and  $X_h$  respectively, we can rewrite system (2.20) in the matrix formulation: find  $(\mathbf{u}_h, p_h, T_h) \in V_h \times Q_h \times X_h$  such that

$$\begin{bmatrix} \mathbf{D}_u + \mathbf{C}(\mathbf{u}_h) & \mathbf{B}^T & 0 \\ \alpha \mathbf{B} & 0 & 0 \\ 0 & 0 & \mathbf{A}(\mathbf{u}_h) + \mathbf{D}_T \end{bmatrix} \begin{bmatrix} \mathbf{u}_h \\ p_h \\ T_h \end{bmatrix} = \begin{bmatrix} \mathbf{F}(T_h) \\ 0 \\ G \end{bmatrix}, \quad (2.21)$$

with, for any  $1 \leq i, j \leq \mathcal{N}_u, 1 \leq m \leq \mathcal{N}_p, 1 \leq k, l \leq \mathcal{N}_T$ ,

$$\begin{aligned} (\mathbf{D}_u)_{ij} &= d_u(\zeta_j, \zeta_i), & (\mathbf{C}(\mathbf{u}_h))_{ij} &= c(\zeta_j, \zeta_i; \mathbf{u}_h), & (\mathbf{F}(T_h))_i &= f(\zeta_i; T_h), \\ (\mathbf{B})_{mj} &= b(\zeta_j, \eta_m) \\ (\mathbf{D}_T)_{kl} &= d_T(\xi_l, \xi_k), & (\mathbf{A}(\mathbf{u}_h))_{kl} &= a(\xi_l, \xi_k; \mathbf{u}_h), & (\mathbf{G})_k &= g(\xi_k). \end{aligned} \quad (2.22)$$

### Dirichlet Boundary Conditions

As explained in the previous paragraph, we usually strongly impose the Dirichlet boundary conditions. We modify the algebraic structures after the assembly of the system (2.21) to force the value on the Dirichlet boundaries.

We introduce the application  $\varphi_u : [1 : \mathcal{N}_u] \rightarrow \Omega$  such that for any  $i \in [1 : \mathcal{N}_u]$ ,  $\mathbf{x} = \varphi_u(i)$  is the point in  $\Omega$  corresponding to the  $i$ -th DoF. We can now define  $\mathcal{I}_{\Gamma_D} = \{i \in [1 : \mathcal{N}_u], \varphi_u(i) \in \Gamma_D\}$ . The entries of the matrix  $\mathbf{M} = \mathbf{D}_u + \mathbf{C}(\mathbf{u}_h)$  and the entries of the vector  $\mathbf{F}(T_h)$  are then modified such that

$$\begin{aligned} (\mathbf{M})_{ij} &= 0, & \forall i, j \in \mathcal{I}_{\Gamma_D}, i \neq j \\ (\mathbf{M})_{ii} &= 1, & \forall i \in \mathcal{I}_{\Gamma_D}, \\ (\mathbf{F})_i &= \mathbf{u}_\Gamma(\varphi_u(i)), & \forall i \in \mathcal{I}_{\Gamma_D}. \end{aligned} \quad (2.23)$$

We apply the same process in order to impose the value  $T_\Gamma$  of temperature field on the boundary  $\Gamma_D^T$ .

Nevertheless, we also implemented a weak formulation to compare the results of the two methods. The weak formulation does not set any constraint on the function spaces, and use Nitsche's penalization methods to impose the Dirichlet conditions through the variational formulation. When imposing the Dirichlet conditions weakly, we add the following terms to the variational formulation:

$$- \int_{\Gamma_D^F} \mathbf{T}(\mathbf{u}_h, p_h) \mathbf{n} \cdot \mathbf{v}_h + (\mathbf{u}_h - \mathbf{u}_D) \cdot \mathbf{T}(\mathbf{v}_h, q_h) \mathbf{n} + \gamma_h^F (\mathbf{u}_h - \mathbf{u}_D) \cdot \mathbf{v}_h \, d\Gamma = 0, \quad (2.24)$$

$$- \int_{\Gamma_D^T} \kappa \nabla T_h \cdot \mathbf{n} S_h + \kappa (T_h - T_D) \nabla S_h \cdot \mathbf{n} + \gamma_h^T (T_h - T_D) S_h \, d\Gamma = 0, \quad (2.25)$$

where  $\gamma_h^T > 0, \gamma_h^F > 0$  are suitable constants defined by elements.

## Slip Boundary Conditions

The slip boundary conditions are mainly required for the wall functions of the turbulence models, see chapter 4. We can strongly impose this kind of condition when the concerned boundary is parallel to the axis. In this particular case, we can edit the algebraic structures as we do for the Dirichlet conditions. However, in a general context - on curved boundaries - the application of slip conditions is not trivial.

Kuzmin *et al.* proposed in [59] an algebraic version but this technique requires an additional linear resolution and it was not appropriate to our solving strategy. As an alternative solution, we initially implemented a penalization formulation method following the design presented in [100]. This technique is close to the Nitsche's approach and is very convenient since it only adds extra terms in the variational formulation. Unfortunately, the convergence rate of this method is very low, and we preferred looking for alternative solutions. We nevertheless detail this formulation, and we present our convergence study for Stokes problem in the appendix 2.

After this unfruitful effort, we found two options. Dione proposes the first one in his thesis work [27]. In this review of the penalization method, he suggested some modifications to boost the convergence rate. Those adjustments mainly aim to

improve the quality of the geometric approximation either with a continuation of the normal-to-the-face vector field or with high-order elements on the curved boundaries. This approach was interesting, but we finally chose to implement the alternative method proposed in [100] using Lagrange multipliers. The formulation is a bit more constraining as it requires introducing a new FE space for the Lagrange multiplier, but it significantly improves the accuracy of the results.

For any point of the boundary  $\partial\Omega$ , we note  $\mathbf{n}$  the outgoing unit normal vector to  $\partial\Omega$ . We also introduce, for any vector  $\mathbf{a}$ , its normal and tangential components respectively denoted by  $(\mathbf{a})_n$  and  $(\mathbf{a})_\tau$ .

$$\begin{aligned} (\mathbf{a})_n &= (\mathbf{n} \cdot \mathbf{a})\mathbf{n} \\ (\mathbf{a})_\tau &= \mathbf{a} - (\mathbf{a})_n \end{aligned} \quad (2.26)$$

We define the Lagrange multiplier space  $\Lambda_h$ , which is a  $\mathbb{P}^0$  discontinuous space on the boundary  $\Gamma_S$ . The fluid part of the variational formulation (2.19) is completed into: find  $(\mathbf{u}_h, p_h, \xi_h) \in V_h \times Q_h \times \Lambda_h$  such that

$$\mathcal{B}_{\gamma,h}((\mathbf{u}_h, p_h, \xi_h), (\mathbf{v}_h, q_h, \lambda_h)) = \mathcal{L}_{\gamma,h}(\mathbf{v}_h, q_h, \lambda_h), \quad \forall (\mathbf{v}_h, q_h, \lambda_h) \in V_h \times Q_h \times \Lambda_h, \quad (2.27)$$

where

$$\begin{aligned} \mathcal{B}_{\gamma,h}((\mathbf{u}_h, p_h, \xi_h), (\mathbf{v}_h, q_h, \lambda_h)) &= \int_{\Omega} (\mathbf{u}_h \cdot \nabla \mathbf{u}_h) \cdot \mathbf{v}_h \, d\Omega + \int_{\Omega} 2\mu \mathbf{S}(\mathbf{u}_h) : \mathbf{S}(\mathbf{v}_h) \, d\Omega \\ &\quad - \int_{\Omega} p_h \nabla \cdot \mathbf{v}_h \, d\Omega - \alpha \int_{\Omega} (\nabla \cdot \mathbf{u}_h) q_h \, d\Omega \\ &\quad + \int_{\Gamma} \xi_h \mathbf{v}_h \cdot \mathbf{n} \, d\Gamma + \int_{\Gamma} \mathbf{u}_h \cdot \mathbf{n} \lambda_h \, d\Gamma \\ &\quad - \gamma \int_{\Gamma} h_f (\xi_h + \mathbf{nT}(\mathbf{u}_h, p_h)\mathbf{n}) (\lambda_h + \epsilon \mathbf{nT}(\mathbf{v}_h, q_h)\mathbf{n}) \, d\Gamma, \end{aligned} \quad (2.28)$$

and

$$\mathcal{L}_{\gamma,h}(\mathbf{v}, q, \lambda) = \int_{\Omega} \mathbf{f} \cdot \mathbf{v} \, d\Omega + \int_{\Gamma} \mathbf{g}_\tau \cdot (\mathbf{v})_\tau \, d\Gamma \quad (2.29)$$

where  $\gamma$  is a positive constant and  $\epsilon = 1, -1$  and  $h_f$  is the characteristic length of the current face. The values of  $\gamma$  and  $\epsilon$  will be discussed in the further numerical tests in the section 4.2.

### 3 Solution Strategy

The Galerkin projection on the FE spaces defined in the previous section requires the resolution of the non-linear system (2.21). The actual resolution of such systems might be non-trivial for different reasons:

- The numerical instability of the model: this issue is prevalent when using the FEM to solve convection dominated problems. We usually overcome this problem by adding stabilization terms in the discrete formulation. The chapter 3 is dedicated to such techniques.
- The choice of the resolution scheme: there are many possibilities to treat the non-linearity of the system. We present and discuss some of these options in the present section.
- The preconditioning of the system: this point might not be an issue when the size of the problem remains reasonable (typically in 2D). In that case, a direct resolution is usually achievable using, for instance, an LU decomposition. However, in 3D we need proper preconditioning strategies to succeed solving large linear systems. This issue is directly related to the choice of the resolution scheme.

During this thesis work, we mainly considered two iterative methods for the resolution of non-linear problems: the Picard algorithm (or fixed point algorithm) and the Newton algorithm. The Newton method theoretically presents an attractive quadratic convergence rate. However, the algorithm is very sensitive to the initial guess, and the convergence radius decrease very fast for high Reynolds number. This point rapidly became an issue, and we had to find new solutions. We implemented the Picard algorithm as an alternative. This method has slow convergence rate compared to the Newton algorithm however, it is usually more robust and less sensitive to the initial guess. The Picard algorithm offered a new option but was still unsuccessful to solve aerothermal problems for high Reynolds numbers.

At this point, we decided to implement a transient version of our model. Our initial motivation was to understand the evolution of the flow when the Reynolds number was growing. But we finally considered the transient resolution as a way to reach the stationary state. We tested different time integration algorithms. The use of the standard backward differentiation formulas, present in Feel++, produced a considerable number of iterations to reach the stationary state. Then we started looking at adaptive time-step, and we implemented the algorithm proposed by Kay *et al.* in [54, 38]. We give implementation details and a review of the method in the appendix 1. Unfortunately, as illustrated in our numerical results, the technique was not entirely stable. We finally dropped the idea of using transient resolution, and we start studying the pseudo time-step methodology.

In the current section, we briefly remind the mentioned iterative methods and the different resolution schemes we envisaged. Then we will detail the pseudo-transient continuation and discuss its advantages. We dedicate the last part to the different preconditioning strategies.

### 3.1 Non-Linear Iterative Methods

#### Picard Algorithm

We consider the generic non linear discrete problem: find  $w \in \mathbb{R}^{\mathcal{N}}$  such that

$$M(w)w = F(w), \quad (2.30)$$

where  $M : \mathbb{R}^{\mathcal{N}} \rightarrow \mathbb{R}^{\mathcal{N} \times \mathcal{N}}$  and  $F : \mathbb{R}^{\mathcal{N}} \rightarrow \mathbb{R}^{\mathcal{N}}$ . This problem can be solved using Picard iterations. From an initial guess  $w^0$ , we build a sequence  $(w^n)$ , hopefully convergent such that  $w^n \rightarrow w$ . For any  $n \geq 0$ , assuming that  $w^n$  is known,  $w^{n+1}$  is computed by solving the linear system

$$M(w^n)w^{n+1} = R(w^n). \quad (2.31)$$

The Picard algorithm presents a wider radius of convergence than the Newton algorithm. However, this method has potentially very slow convergence rate and is very sensitive to local minima.

#### Newton Algorithm

We consider the generic non-linear discrete system: find  $w \in \mathbb{R}^{\mathcal{N}}$  such that

$$\mathcal{F}(w) = 0 \quad (2.32)$$

where  $\mathcal{F} : \mathbb{R}^{\mathcal{N}} \rightarrow \mathbb{R}^{\mathcal{N}}$  is supposed to be differentiable. Here again we iteratively build a sequence  $(w^n)$  such that  $w^n \rightarrow w$ , from an initial guess  $w^0$ . For any  $n \geq 0$ , assuming that  $w^n$  is known, we set  $w^{n+1} = w^n + \delta^n$ , with  $\delta^n$  solution of the linear system

$$\mathcal{F}'(w^n)\delta^n = -\mathcal{F}(w^n), \quad (2.33)$$

where  $\mathcal{F}'$  is the Jacobian matrix of  $\mathcal{F}$ . This method has a theoretical quadratic convergence, but is very sensitive to the choice of the initial guess  $w^0$ . Local minima of the function  $\mathcal{F}$  might also be an issue.

#### Resolutions Schemes

These two algorithms work with the same idea: splitting the non-linear resolution in multiple linear resolutions. This technique might be very costly if the algorithm is not able to quickly converge. In both cases, the choice of an appropriate initial guess is crucial since it affects both the capability to converge and the number of required iterations.

Once we dispose of these iterative solvers, we have to decide if we want to treat the system (2.21) in a monolithic way or by splitting the problem. The monolithic

resolution consists in the direct application of an iterative solver to the non-linear system

$$\begin{bmatrix} \mathbf{D}_u + \mathbf{C}(\mathbf{u}_h) & \mathbf{B}^\top & 0 \\ \alpha \mathbf{B} & 0 & 0 \\ 0 & 0 & \mathbf{A}(\mathbf{u}_h) + \mathbf{D}_T \end{bmatrix} \begin{bmatrix} \mathbf{u}_h \\ p_h \\ T_h \end{bmatrix} = \begin{bmatrix} \mathbf{F}(T_h) \\ 0 \\ \mathbf{G} \end{bmatrix}. \quad (2.34)$$

The decoupled resolution consists in fixed point iterations between the resolution of the non-linear system (using our favorite iterative method)

$$\begin{bmatrix} \mathbf{D}_u + \mathbf{C}(\mathbf{u}_h) & \mathbf{B}^\top \\ \alpha \mathbf{B} & 0 \end{bmatrix} \begin{bmatrix} \mathbf{u}_h \\ p_h \end{bmatrix} = \begin{bmatrix} \mathbf{F}(T_h) \\ 0 \end{bmatrix} \quad (2.35)$$

and the linear one

$$[\mathbf{A}(\mathbf{u}_h) + \mathbf{D}_T] [T_h] = [\mathbf{G}], \quad (2.36)$$

The different options were implemented and tested. Each configuration has its advantages and will be chosen following the needs. For instance, the decoupled version is particularly efficient for transient resolution. In this particular case, we can use an Oseen linearization for the fluid system to run quick simulations. On the other hand, this decoupled scheme is hardly practicable for steady resolution since we have to manage the non-linearity between  $\mathbf{u}$  and  $p$  and then between the fluid and  $T$ . When we finally decided to only work with the steady state directly, we abandoned the decoupled resolution.

## Continuation

The remaining issue was the choice of the initial guess. Our framework was still unable to reach the solution for more dynamic problems (smaller viscosity and diffusivity). To improve the convergence of our solver, we envisaged two solutions: a physical continuation and a transient continuation. We finally implemented a mixed solution in which both kind of continuation can be used.

The physical continuation iteratively solves the system with evolving parameters. Each new resolution uses the previous solution as an initial guess. In practice, the continuation is made on the viscosity  $\mu$  and the thermal diffusivity  $\kappa$ . For a continuation of  $N$  steps, we will successively solve

$$\mathcal{F}(\mathbf{u}^k, p^k, T^k; \mu^k, \kappa^k) = 0 \quad (2.37)$$

where  $\mu^k = \varepsilon(k)\mu$  and  $\kappa^k = \varepsilon(k)\kappa$ ,  $\varepsilon(N) = 1$ . There are many possibilities for the design of the continuation function  $\varepsilon$ . We chose a log-equidistributed partitioning of the viscosity range with

$$\varepsilon(k) = \frac{\mu^{k/N}}{\mu}. \quad (2.38)$$

This technique is very efficient and particularly stable for high Reynolds simulations. However, we usually need a well tuned continuation to guaranty the convergence.



To remedy this issue, we implemented the pseudo-transient continuation which remarkably improves the convergence of the Newton algorithm. This technique allows to remove the physical continuation or, at least, to only consider few steps.

### 3.2 Pseudo Transient Continuation

The pseudo-transient continuation [55, 19] is a modification of the Newton algorithm to improve the convergence capabilities of the method. The continuation is made by imposing local transient states to the system, to progressively reach the steady state.

We consider a Differential Algebraic Equation

$$\mathcal{F}(\boldsymbol{w}) = 0 \quad (2.39)$$

resulting of a FEM discretization. In a generic formulation, we suppose  $\boldsymbol{w} \in \mathcal{N}$  and  $\mathcal{F} : \mathbb{R}^{\mathcal{N}} \rightarrow \mathbb{R}^{\mathcal{N}}$  differentiable. The non-linear system  $\mathcal{F}(\boldsymbol{w}) = 0$  has to be solved using an iterative method of successive linear resolutions. For such techniques, the choice of the initial guess is crucial. For instance, for the Navier-Stokes system, for a low Reynolds number, a common initialization is to start by resolving the associated stokes problem. For larger Reynolds numbers, the convergence radius of iterative methods becomes smaller, and we have to choose a better initial guess.

The pseudo-transient continuation ( $\Psi$ tc) method is particularly adapted for problems deriving from a transient dynamic. From the generic stationary equation (2.39), we can define the transient problem associated

$$\boldsymbol{w}' + \mathcal{F}(\boldsymbol{w}) = 0, \boldsymbol{w}(0) = \boldsymbol{w}_0, \quad (2.40)$$

where  $\boldsymbol{w}' = \frac{\partial \boldsymbol{w}}{\partial t}$ . Then, it is natural to seek a root of  $\mathcal{F}$  as the limit of  $\boldsymbol{w}(t)$

$$\boldsymbol{w} = \lim_{t \rightarrow \infty} \boldsymbol{w}(t). \quad (2.41)$$

When we are not interested in the transient states, we can reach this limit faster by considering non-physic intermediary states of the problem. This idea is the main ingredient of the pseudo-transient continuation.

In [55], Kelley *et al.* propose a continuation algorithm that integrates iteratively

$$\boldsymbol{V} \boldsymbol{w}' = \mathcal{F}(\boldsymbol{w}), \quad \boldsymbol{w}(0) = \boldsymbol{w}_0 \quad (2.42)$$

by a variable time step  $\delta_n$ . This time step is supposed to grow as the residual  $\mathcal{F}(\boldsymbol{w}_n)$  approaches zero. Here  $\boldsymbol{V}$  is a nonsingular matrix used to improve the scaling of the problem and usually build to balance the local CFL number through the domain.

The convergence results for  $\Psi$ tc in [55] are extended in [19] for systems in which it is more convenient to not parabolize all the equations. For instance, in

: Newton Algorithm	: $\Psi$ tc Algorithm
1 Choose $\varepsilon$ ;	1 Choose $\varepsilon$ ;
2 Set $\mathbf{w} = \mathbf{w}_0$ ;	2 Set $\mathbf{w} = \mathbf{w}_0$ and $\delta = \delta_0$ ;
3 <b>repeat</b>	3 <b>repeat</b>
4     Solve $\mathcal{F}'(\mathbf{w})s = -\mathcal{F}(\mathbf{w})$ ;	4     Solve $(\frac{1}{\delta}D + \mathcal{F}'(\mathbf{w}))s = -\mathcal{F}(\mathbf{w})$ ;
5     Set $\mathbf{w} = \mathbf{w} + s$ ;	5     Set $\mathbf{w} = \mathbf{w} + s$ ;
6     Evaluate $\mathcal{F}(\mathbf{w})$ ;	6     Evaluate $\mathcal{F}(\mathbf{w})$ ;
7 <b>until</b> $\ \mathcal{F}(\mathbf{w})\  < \varepsilon$ ;	7     Update $\delta$ ;
	8 <b>until</b> $\ \mathcal{F}(\mathbf{w})\  < \varepsilon$ ;

Figure 2.1 – Comparison of  $\Psi$ tc and Newton algorithms

the incompressible Navier-Stokes problem, it is entirely natural only to evolve the velocity through pseudo-time stepping. It results in the equations,

$$D\mathbf{w}' = -\mathcal{F}(\mathbf{w}), \quad (2.43)$$

with

$$D = \begin{bmatrix} \mathbf{V} & \\ & 0 \end{bmatrix} \quad (2.44)$$

with  $\mathbf{V}$  a non-singular matrix. The design of this matrix will be discussed in a further paragraph.

The procedure of  $\Psi$ tc is very close to the usual Newton method, but we add the matrix  $\frac{1}{\delta_n}D$  to the Jacobian matrix,  $\mathcal{F}'$ , of the function  $\mathcal{F}$ . The pseudo transient iteration then reads

$$\mathbf{w}_{n+1} = \mathbf{w}_n - (1/\delta_n D + \mathcal{F}'(\mathbf{w}_n))^{-1} \mathcal{F}(\mathbf{w}_n). \quad (2.45)$$

As a comparison we remind the Newton algorithm and present the  $\Psi$ tc algorithm in figure 2.1.

### Choice of parameter $\delta_n$

An updating procedure for the pseudo time step  $\delta$  is proposed in [55] using the Switched Evolution Relaxation. This method increases the time step in inverse proportion of the residual reduction.

$$\delta_{n+1} = \delta_n \frac{\|\mathcal{F}(\mathbf{w}_{n-1})\|}{\|\mathcal{F}(\mathbf{w}_n)\|} = \delta_0 \frac{\|\mathcal{F}(\mathbf{w}_0)\|}{\|\mathcal{F}(\mathbf{w}_n)\|} \quad (2.46)$$

In order to control the evolution of this pseudo time step, a bounding function  $\phi$  is often added

$$\delta_{n+1} = \phi\left(\delta_n \frac{\|\mathcal{F}(\mathbf{w}_{n-1})\|}{\|\mathcal{F}(\mathbf{w}_n)\|}\right), \quad (2.47)$$

with

$$\phi(\xi) = \min(\xi, \delta_{\max}), \quad (2.48)$$

where  $\delta_{\max}$  is a custom threshold parameter. When this bounding function is used, the system switches into a steady state resolution after some iterations. The initial value  $\delta_0$  is another custom parameter and has to be chosen depending of the dynamic of the problem.

Another version of the Switched Evolution Relaxation proposes not to consider the norm of the residual  $\|\mathcal{F}(\mathbf{w})\|$  but the norm of the increment  $\|\mathbf{w}_n - \mathbf{w}_{n-1}\|$ . However, this method is more prone to convergence failures and requires a uniformly well conditioned Jacobian. The update of the time step then reads

$$\delta_{n+1} = \phi(\delta_n \|\mathbf{w}_n - \mathbf{w}_{n-1}\|). \quad (2.49)$$

This alternative was not studied in this thesis work since the first version gives satisfying results.

## The Scaling Matrix

The construction of the scaling matrix  $\mathbf{V}$  is made to balance the local CFL number,

$$C_K = \Delta t \sum_{i=1}^d \frac{u_i}{h_i} \quad (2.50)$$

where  $\Delta t$  is the time step,  $u_i$  is the component of the velocity field  $h_i$  is the characteristic length of the cell in the  $i$ -eme direction of the element  $K$ . Hence it naturally appears to choose the scaling factor per-element

$$\omega_K = \frac{\|\mathbf{u}\|_2}{h} \quad (2.51)$$

where  $\|\mathbf{u}\|_2$  is the usual norm-2 of the convection field on the center of the element and  $h$  is the characteristic length of the element. The evaluation of  $h$  is discussed in chapter 3.

## 4 Validation

In this section, we present the results obtained with our aerothermal framework on different test-cases. The following problems are usually fundamental and aim to validate our implementation, see chapter 8. By default the configuration is the following:

- $\mathbb{TH}_{ch}^1$  FE space for the fluid part,

- $\mathbb{L}_{ch}^2$  FE space for the temperature,
- domain meshed with linear simplex,
- unstructured homogeneous mesh of characteristic size  $h$ ,
- monolithic resolution with Newton algorithm and eventually  $\Psi$ tc continuation,
- LU direct solver in 2D.

This default configuration may be edited for certain test-cases, it will be specified where appropriate.

The common quantities of interest are the following errors (usually normalized):

- $\mathcal{L}^2$ -error on the velocity:  $\|\mathbf{u} - \mathbf{u}_h\|_0$ ,
- $H^1$ -error on the velocity:  $\|\mathbf{u} - \mathbf{u}_h\|_1$ ,
- $\mathcal{L}^2$ -error on the pressure (up to an additive constant):  $\|p - p_h\|_0$ ,
- $\mathcal{L}^2$ -error on the temperature:  $\|T - T_h\|_0$ .

## 4.1 Validation of the Fluid Solver

### Fluid Convergence Study, Kovaszny Test Case

We run a convergence study on the well known case proposed by Kovaszny in [58]. The geometry of the problem is a 2D grid,  $\Omega = [-0.5, 2] \times [-0.5, 1.5]$ . The analytic solution is given by

$$\begin{aligned} u &= (1 - e^{\lambda x} \cos(2\pi y), \frac{\lambda}{2\pi} e^{\lambda x} \sin(2\pi y)) \\ p &= p_0 - e^{\lambda x} \end{aligned} \quad (2.52)$$

where  $\lambda = \frac{\text{Re}}{2} - \sqrt{\frac{\text{Re}^2}{4} + 4\pi^2}$ , see figure 2.2. For the present study we fixed  $\text{Re} = \frac{\rho}{\mu u} = 40$ . We impose the analytic solution as strong Dirichlet condition on the boundaries of the domain, and we treat the non-linearity with Newton algorithm. We measured the three errors :  $\|\mathbf{u}_h - \mathbf{u}_e\|_0$ ,  $\|\mathbf{u}_h - \mathbf{u}_e\|_1$  and  $\|p_h - p_e\|_0$ . The evolution of these quantities with respect to the characteristic size  $h$  is presented in figure 2.3. The study is performed with or without stabilization method (SUPG<sub>1d</sub>, see chapter 3 ), for different polynomial orders :  $\text{TH}_h^1$ ,  $\text{TH}_h^2$  and  $\text{TH}_h^3$ . We retrieve the expected convergence rates for  $\mathbf{u}$  and  $p$ . We also notice that the use of stabilization method does not influence the results. It proves the consistency of the method: the stabilization terms vanish when the characteristic length of the mesh goes to zero. Those results are identical with other configuration for the stabilization method.

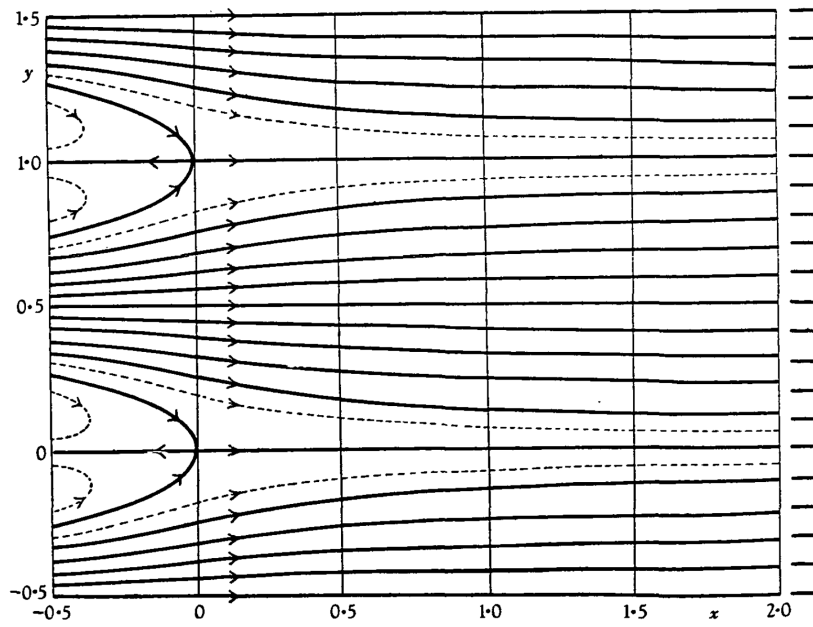


Figure 2.2 – Streamline of the analytic solution, capture from [58]

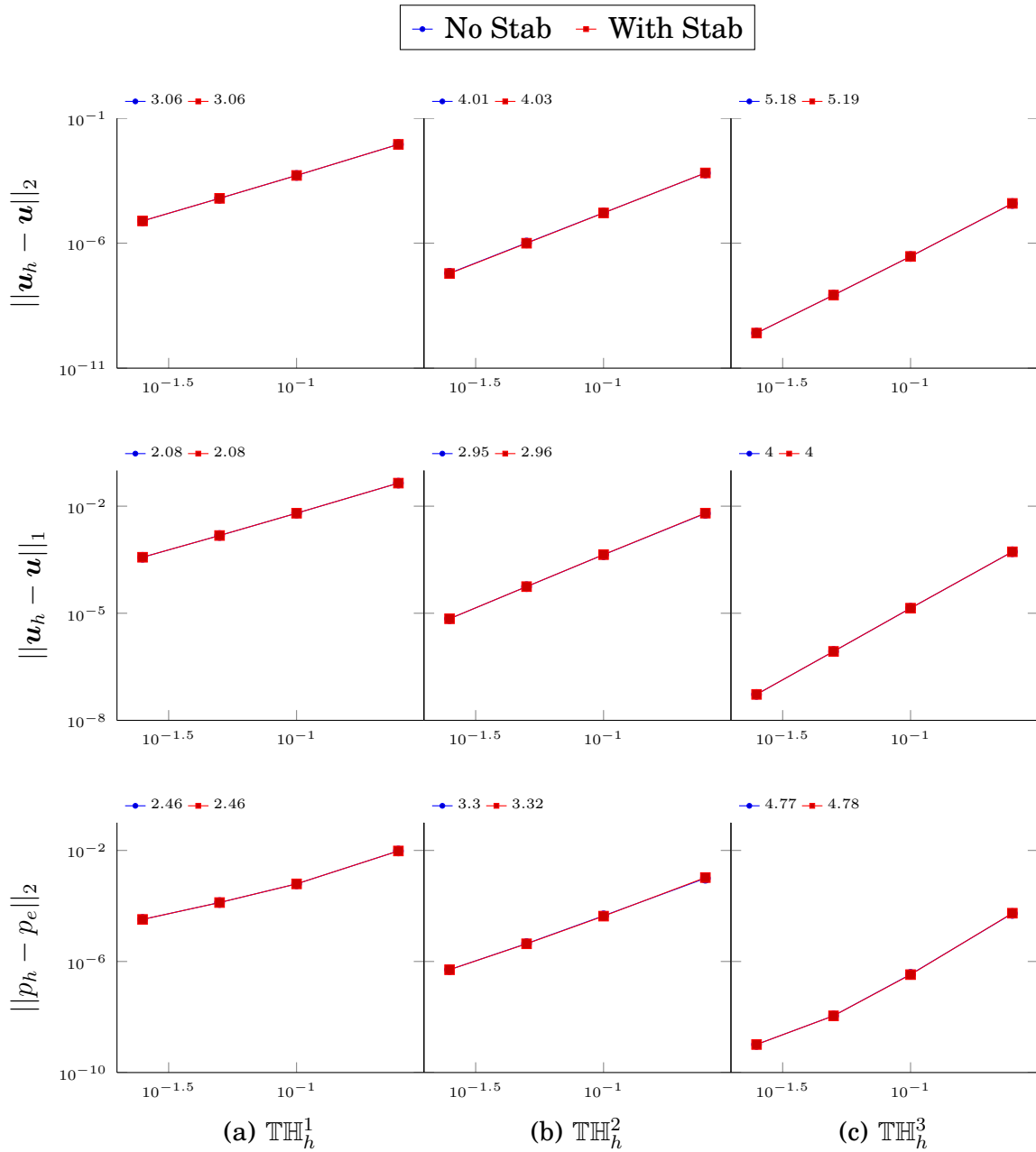


Figure 2.3 – Kovasznay test-case. Convergence study for different polynomial orders, with or without stabilization method. Log/Log scale. Stab and No-Stab plots are perfectly stacked

**Turek 2D Benchmark: Flow Around a Cylinder, Re=20**

This benchmark is proposed on <http://www.featflow.de>. The geometry of the model is a pipe without a circular cylinder of radius  $r = 0.05$ , as described in figure 2.4. No-slip boundary conditions are imposed on upper and lower walls and

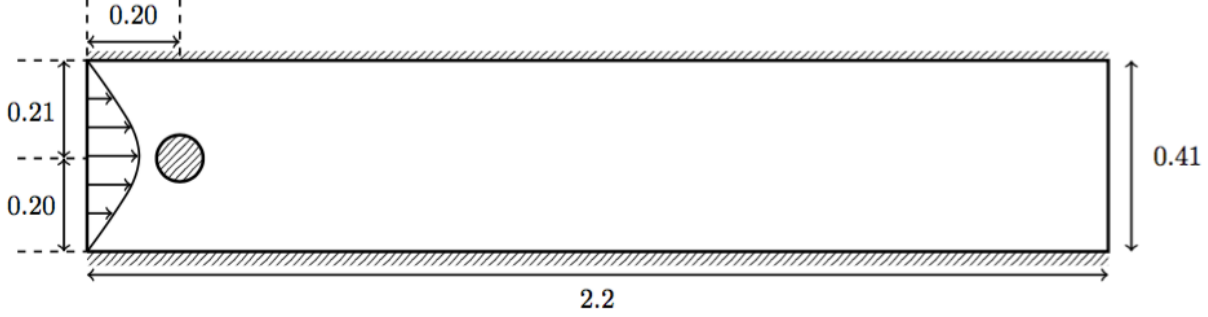


Figure 2.4 – Geometry of the Problem

on the boundary of the cylinder. The left edge is the input of the pipe where a Poiseuille profile is prescribed,

$$u(0, y) = \left( \frac{4U_y(0.41 - y)}{0.41^2}, 0 \right) \quad (2.53)$$

with the maximum amplitude of the velocity  $U_y = 0.3$ . The right edge is the output of the pipe and with a do-nothing condition

$$\bar{\sigma} = 0. \quad (2.54)$$

The parameters of the fluid are set to  $\rho = 1.0$  and  $\mu = 0.001$ . With this configuration, the Reynolds number obtained is  $Re = 20$  and the system turns into a stationary state.

In this benchmark, we are interested in three quantities. The lift and drag coefficients are respectively defined by

$$C_D = \frac{2}{\bar{U}^2 L} F_D, \quad C_L = \frac{2}{\bar{U}^2 L} F_L \quad (2.55)$$

with  $\bar{U} = 0.2$  the mean velocity at the input,  $L = 0.1$  the characteristic length of the flow configuration and  $(F_D, F_L)$  such that

$$(F_D, F_L) = \int_{\Gamma_C} \bar{\sigma} n d\Gamma \quad (2.56)$$

where  $\Gamma_C$  is the boundary of the cylinder and  $n$  is the outer normal vector of the circle. The last quantity of interest is the pressure difference  $\Delta_p = p(a_1) - p(a_2)$  between the points  $a_1 = (0.15, 0.2)$  and  $a_2 = (0.25, 0.2)$ . We will compare our results with the reference values:  $C_D = 5.57953523384$ ,  $C_L = 0.010618948146$  and  $\Delta_p = 0.11752016697$ .

We tested the framework with or without stabilization method (SUPG1<sub>d</sub>, see chapter 3). The results obtained with our fluid solver are coherent with the reference values, see table 2.5.

	No Stab		With Stab	
$C_D$	5.5769	[0.047]	5.5771	[0.043]
$C_L$	0.010629	[0.100]	0.010630	[0.110]
$\Delta_p$	0.11751	[0.008]	0.11751	[0.007]

Figure 2.5 – Turek benchmark: flow around a cylinder,  $Re = 20$ . Numerical results for the drag coefficient, lift coefficient and pressure difference. In bracket: relative error in %, compared with the reference values.

## 4.2 Validation of the Slip Boundary Condition Implementation

We present, in this subsection, a convergence study for two problems with slip conditions: one with a straight boundary and the second one with a curved boundary. The principle will be identical for both test-cases: we solve the Navier-Stokes equations on a closed domain, where a manufactured function is imposed on the boundary, using slip conditions. We can then compare the numerical solutions with the exact solutions and we evaluate the quantities:  $\|\mathbf{u}_h - \mathbf{u}_e\|_0$ ,  $\|\mathbf{u}_h - \mathbf{u}_e\|_1$  and  $\|p_h - p_e\|_0$ . For each problem, we display the evolution of these errors as a function of the characteristic size of the mesh. To determine the best configuration, we performed simulations for different values of the constants  $\alpha$ ,  $\gamma$  and  $\epsilon$ . We remind the role of these constants:

- $\alpha = \pm 1$  determine the treatment of the pressure/velocity term in the Navier-Stokes weak formulation (symmetric or skew-symmetric), see (2.16),
- $\gamma$  is a positive penalty constant used in the stabilized formulation for slip boundary conditions, see (2.28)
- $\epsilon = \pm 1$  determine the form of the stabilization terms in this formulation for slip boundary conditions: symmetric or skew-symmetric, see (2.28).

We also performed one simulation with strongly-imposed Dirichlet conditions to compare the convergence rates.

### Slip Conditions on a Square

The domain is a square  $\Omega = [-1, 1] \times [-1, 1]$ . The viscosity is supposed to be constant  $\mu = 1$ . The manufactured solution is

$$\begin{aligned} \mathbf{u} &= (2y(1 - x^2), -2x(1 - y^2)), \\ p &= x^2 + y^2. \end{aligned} \tag{2.57}$$

With this solution we can evaluate the slip condition on the boundary

$$\begin{aligned} \mathbf{u} \cdot \mathbf{n} &= 0, \\ (\mathbf{T}(\mathbf{u}, p)\mathbf{n})_\tau &= (2y(1 - x^2), -2x(1 - y^2)). \end{aligned} \tag{2.58}$$



The source term is chosen to fit with this solution. The discrete solution is computed for different values of  $h = 0.2, 0.1, 0.05, 0.025$ . The method only converges when  $\epsilon = \alpha$ , so we presented results for  $(\alpha, \epsilon) = (1, 1)$  and  $(\alpha, \epsilon) = (-1, -1)$ . We also tested different values for  $\gamma = 10^{-4}, 10^{-2}, 1, 10^2, 10^4$ . The convergence study is summarized on figure 2.6.

### Slip Conditions on a Ring

The second test-case is made on a ring, defined by  $1 \leq x^2 + y^2 \leq 4$ . The exact solution

$$\begin{aligned} \mathbf{u} &= (-y\sqrt{x^2 + y^2}, x\sqrt{x^2 + y^2}), \\ p &= x^2 + y^2, \end{aligned} \quad (2.59)$$

is imposed with a Dirichlet condition on the internal boundary ( $x^2 + y^2 = 1$ ) and with the slip condition

$$\begin{aligned} \mathbf{u} \cdot \mathbf{n} &= 0, \\ (\mathbf{T}(\mathbf{u}, p)\mathbf{n})_\tau &= (-2y, 2x), \end{aligned} \quad (2.60)$$

on the external boundary ( $x^2 + y^2 = 4$ ). Again the source term  $f$  is designed to obtain the exact solution in the domain. The results for this test-case can be seen on figure 2.7.

Considering the results of the figures 2.6 and 2.7, we observe that:

- as announced in [100], the version with  $(\alpha, \epsilon) = (-1, -1)$  is more robust, and seems to be unconditionally stable,
- the convergence is globally better for small values of  $\gamma$ ,
- we almost retrieve the theoretical convergence rate, except for the  $\mathcal{L}^2$  error on  $\mathbf{u}$ , on the ring,
- with a suited configuration, the method is as efficient for curved or straight boundary.

From those observations, we can choose our favorite setup for later simulations with slip boundary conditions:  $\gamma = 10^{-4}$ ,  $\alpha = -1$  and  $\epsilon = -1$ . These values will be used by default. For that reason, the parameter  $\alpha$  will not appear in the variational formulation anymore: we always use the anti-symmetric formulation.

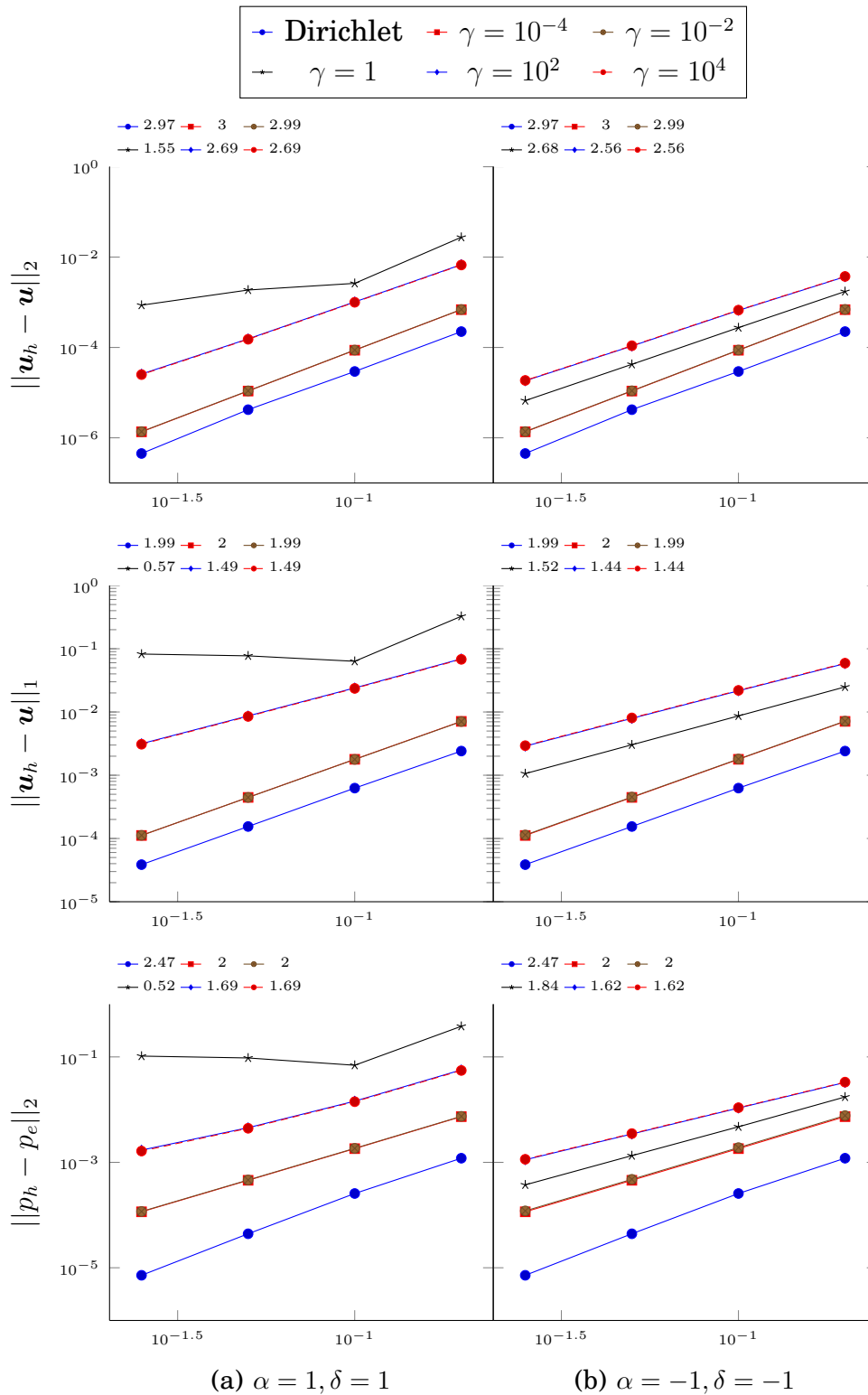


Figure 2.6 – Slip conditions applied on a squared geometry. Convergence study of the errors between exact and FE solutions as function of the characteristic length  $h$ . Log/Log scale. Linear regression and slope displayed for different values of  $\alpha, \gamma$  and  $\delta$ .

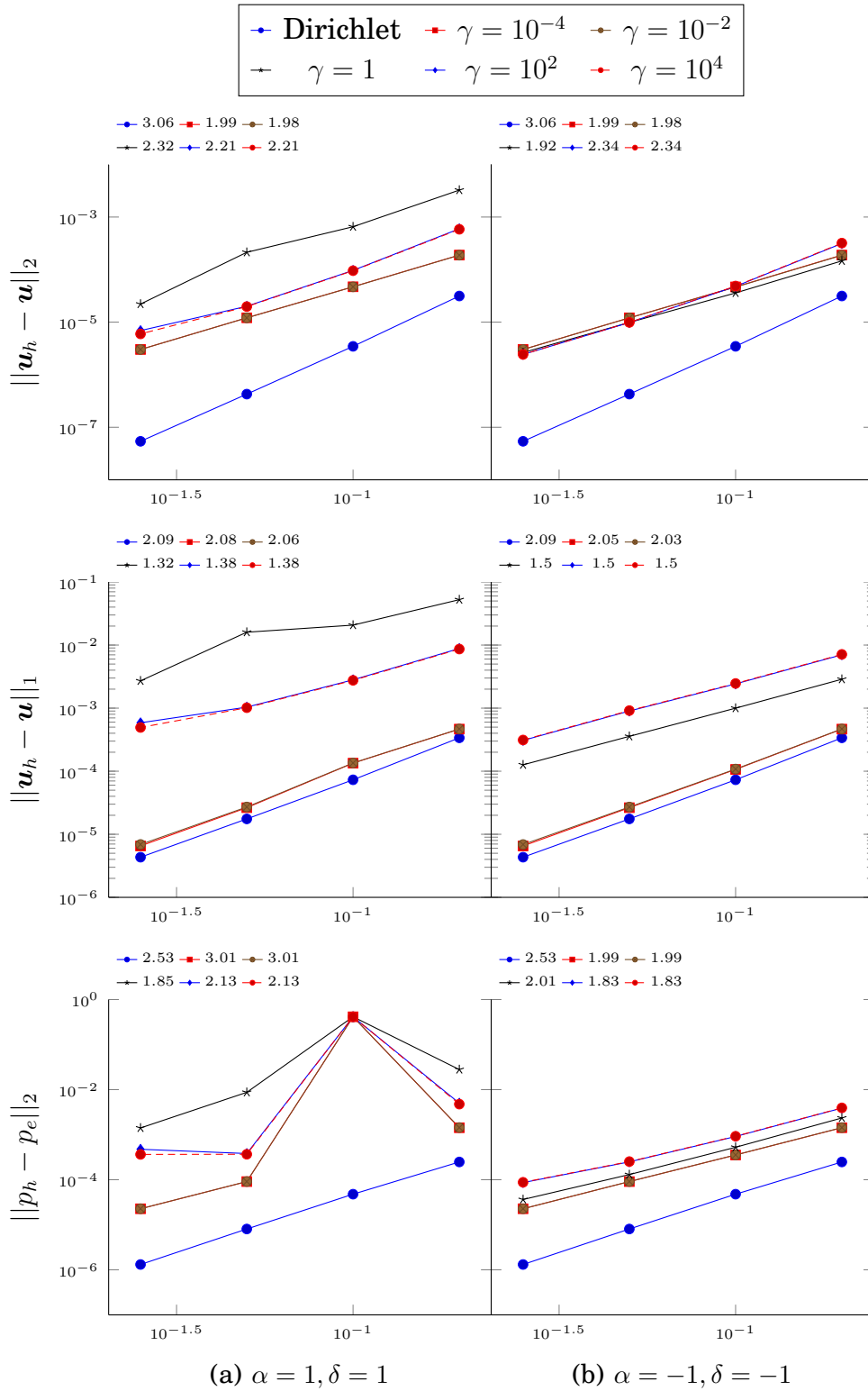


Figure 2.7 – Slip conditions applied on a ring geometry. Convergence study of the errors between exact and FE solutions as function of the characteristic length  $h$ . Log/Log scale. Linear regression and slope displayed for different values of  $\alpha$ ,  $\gamma$  and  $\delta$ .

### 4.3 Validation of the Aerothermal Solver

#### The Buoyancy Driven Cavity

We test here our aerothermal framework on the closed insulated heated cavity benchmark. The domain is a square  $\Omega = [0, 1] \times [0, 1]$ , the left and right walls have fixed temperature, respectively to  $T = 1$  and  $T = 0$ . The bottom and top walls are supposed to be adiabatic. We also impose a no-slip boundary condition,  $\mathbf{u} = (0, 0)$ , on the whole boundary. This setup is summarized in figure 2.8. We solved the aerothermal problem with different parameters. For this benchmark, the aerothermal system is usually written in a dimensionless formalism and the regime can be characterized using the Rayleigh number,

$$\text{Ra} = \frac{\mathbf{g}\beta}{\nu\kappa}\Theta L^3, \quad (2.61)$$

where  $\nu$ ,  $\mathbf{g}$ ,  $\beta$  and  $\kappa$  are the physical parameters defined in previously,  $\Theta$  is the characteristic temperature difference and  $L$  is the characteristic length of the problem.

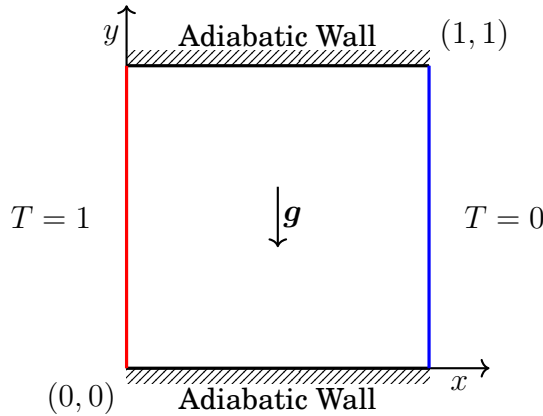


Figure 2.8 – Buoyancy Driven Cavity Setup

We performed simulations for  $10^3 \leq \text{Ra} \leq 10^8$ , and we measured the following quantities:

- the maximum vertical velocity at the mid-height (value and location), table 2.1,
- the maximum horizontal velocity at the mid-width (value and location), table 2.2,
- the maximum Nusselt number on the hot wall (value and location), table 2.3,
- the averaged Nusselt number on the hot wall, table 2.4.

Ra	Present	Ref.[106]	Ref.[25]	Ref.[67]	Ref.[69]
$10^3$	3.697 (0.178)	3.686 (0.188)	3.679 (0.179)	3.73 (0.182)	3.696 (0.179)
$10^4$	19.63 (0.118)	19.79 (0.120)	19.51 (0.120)	19.90 (0.124)	19.61 (0.119)
$10^5$	68.59 (0.068)	70.63 (0.072)	68.22 (0.066)	70.0 (0.068)	68.69 (0.067)
$10^6$	221.1 (0.038)	227.1 (0.040)	216.8 (0.039)	228.0 (0.039)	220.8 (0.038)
$10^7$	697.6 (0.022)	714.5 (0.022)	-	698.0 (0.024)	703.3 (0.022)
$10^8$	2206 (0.12)	2259 (0.012)	-	-	2223 (0.013)

Table 2.1 – Buoyancy Driven Cavity. Maximum vertical velocity ( $u_y$ ) at the mid-height ( $y = 0.5$ ). Value ( $x$  position) for different Rayleigh numbers

Ra	Present	Ref.[106]	Ref.[25]	Ref.[67]	Ref.[69]
$10^3$	3.649 (0.812)	3.489 (0.813)	3.634 (0.813)	3.680 (0.817)	3.649 (0.813)
$10^4$	16.18 (0.822)	16.12 (0.815)	16.20 (0.823)	16.10 (0.817)	16.18 (0.824)
$10^5$	34.73 (0.852)	33.39 (0.835)	34.81 (0.855)	34.00 (0.857)	34.77 (0.854)
$10^6$	64.82 (0.848)	65.40 (0.860)	65.33 (0.851)	65.40 (0.875)	64.69 (0.846)
$10^7$	148.6 (0.878)	143.6 (0.922)	-	139.7 (0.919)	145.3 (0.885)
$10^8$	321 (0.927)	296.71 (0.930)	-	-	283.7 (0.946)

Table 2.2 – Buoyancy Driven Cavity. Maximum horizontal velocity ( $u_x$ ) at the mid-width ( $x = 0.5$ ). Value ( $y$  position) for different Rayleigh numbers

Our results are compared with reference values from [106], [25], [67] and [69]. We can reasonably conclude in the accuracy of our aerothermal framework. The measured values perfectly fit with the references, even for high Rayleigh numbers ( $Ra \geq 10^7$ ) when the resolution is definitely not trivial.

Ra	Present	Ref.[106]	Ref.[25]	Ref.[67]	Ref.[69]
$10^3$	1.506 (0.088)	1.501 (0.08)	1.50 (0.092)	1.47 (0.109)	1.506 (0.090)
$10^4$	3.531 (0.142)	3.579 (0.13)	3.53 (0.143)	3.47 (0.125)	3.530 (0.143)
$10^5$	7.721 (0.082)	7.945 (0.08)	7.71 (0.08)	7.71 (0.08)	7.708 (0.084)
$10^6$	17.56 (0.038)	17.86 (0.03)	17.92 (0.038)	17.46 (0.039)	17.53 (0.038)
$10^7$	39.92 (0.018)	38.6 (0.015)	-	30.46 (0.024)	41.02 (0.039)
$10^8$	93.63 (0.008)	91.16 (0.010)	-	-	91.21 (0.067)

Table 2.3 – Buoyancy Driven Cavity. Maximum Nusselt number on the hot wall ( $x = 0$ ). Value ( $y$  position) for different Rayleigh numbers.

Ra	Present	Ref.[106]	Ref.[25]	Ref.[67]
$10^3$	1.118	1.117	1.12	1.074
$10^4$	2.245	2.254	2.243	2.084
$10^5$	4.522	4.598	4.52	4.3
$10^6$	8.835	8.976	8.8	8.743
$10^7$	16.63	16.66	-	13.99
$10^8$	31.12	31.49	-	-

Table 2.4 – Buoyancy Driven Cavity. Averaged Nusselt number on the hot wall ( $x = 0$ ). Value for different Rayleigh numbers.

## Conclusion

In this chapter, we detailed the different techniques implemented for the resolution of an aerothermal problem:

- formulation: incompressible Navier-Stokes equations coupled with an energy equation in the Boussinesq approximation,
- discretization: finite element method
- non linear resolution:
  - iterative method: Newton algorithm
  - continuation method on physical parameters
  - pseudo-transient continuation

In this framework, we developed an efficient tool for the resolution of aerothermal problems. We also proposed techniques to reach high Reynolds numerical solutions. And eventually, the numerical results confirmed the adequacy of our implementation.

However, in this chapter, we were only working on small academic problems, and we were free to refine the meshes when necessary. This refining solution is hardly practicable for industrial problems, and we have to consider new approaches. To avoid the Direct Numerical Simulation (DNS), we will introduce turbulence models in the chapter 4. But in a first instance, we have to deal with the instability issue of the Galerkin method. This delicate point is detailed in the next chapter.



---

Stabilization Methods

---

**Contents**


---

<b>1</b>	<b>Streamline Diffusion Methods . . . . .</b>	<b>47</b>
1.1	Stabilized Formulation . . . . .	48
1.2	The Navier-Stokes Equations . . . . .	48
1.3	Convergence of the Method . . . . .	49
<b>2</b>	<b>Stabilization Parameter . . . . .</b>	<b>50</b>
2.1	Extension to Advection Diffusion Reaction Problems . . . . .	52
2.2	Choice of the characteristic length . . . . .	53
<b>3</b>	<b>Shock Capturing Method . . . . .</b>	<b>55</b>
3.1	Isotropic Artificial Diffusion . . . . .	55
3.2	Artificial Diffusion Added Orthogonally to Streamlines . . . . .	56
<b>4</b>	<b>Validation . . . . .</b>	<b>57</b>
4.1	Comparison of the Stabilization Methods . . . . .	57
4.2	Shock Capturing Methods . . . . .	64
4.3	Stabilized Navier-Stokes Equations . . . . .	66

---

The instability of the Galerkin method is a well-known issue of advection/convection dominated problems. When we started working on aerothermal simulations, the physical parameters were very smooth and not convection dominated. In these conditions, the stability of the solution was not an issue. Our straightforward recipe



was to refine the mesh when oscillations appeared. However, when we started looking at industrial configurations - with realistic physical parameters - it was not possible to refine meshes that much. We had to introduce new ingredients.

We naturally studied turbulence models and started working on the implementation. Nevertheless, it appeared very soon that the turbulence model would not be enough. We needed a specific feature to handle the numerical instability. We considered two kinds of method: the Continuous Interior Penalty (CIP) [11], and the Streamline Diffusion Methods (SDM) [46]. After some tests with the CIP method, we were not fully satisfied: the procedure was either not stable enough or too diffusive. We finally chose to focus our efforts on the SDM which are, by the way, the most common techniques to deal with turbulent flows.

The SDM were first introduced by Hughes and Brooks [46, 7], with the Streamline Upwind/Petrov-Galerkin (SUPG) method for advection dominated flows. Then, this technique has been extended into two new stable Petrov-Galerkin formulations: the Galerkin Least Square (GLS), Hughes *et al.* [43] and the Douglas-Wang (or SGS) for Stokes equations, Douglas and Wang [29]. In the following section, we propose a review of these different approaches for the advection-diffusion-reaction equation and the Navier-Stokes system.

Our implementation of those methods has been enriched progressively to face different issues:

- First SUPG implementation: efficient on homogeneous isotropic meshes. The design of the stabilization parameter was very basic and did not take into account the polynomial order of the finite element space.
- We implemented a new design for the stabilization parameter, and we introduced the weighting coefficient  $\lambda_K$ . The stabilization on anisotropic meshes (boundary layers) with very stretched cells remains inefficient.
- We expected some improvement with the implementation of the GLS and SGS formulations. The GLS helps in some configuration, but the SGS does not provide any significant upgrade. For that reason, we only propose numerical results for GLS and SUPG in this chapter.
- A significant enhancement came with the implementation of a new approach for the evaluation of the characteristic length. The new design is adapted to triangular cells and takes into account the direction of the advection field. Thus, the stabilization parameter became more efficient in the boundary layers, and we finally obtained stable solutions. However, the SDM still produces some under/overshoots in regions of sharp discontinuities. These small oscillations usually generate new kind of instability within the turbulence transport equations.
- The under/overshoots issue have been covered with the implementation of shock capturing (or discontinuity capturing) methods. With this new feature,

we can finally envisage the modeling of turbulent flows.

We give the details of those developments in the further sections. In the last section, we also propose a comparison of the different configurations with some relevant test cases. Hopefully, at the end of the chapter we will be able to determine the optimal stabilization technique for each kind of problems.

## 1 Streamline Diffusion Methods

Let  $X$  an Hilbert space on  $\Omega$ , with scalar product  $(\cdot, \cdot)_X$  and associated norm  $\|\cdot\|_X$ . We introduce a stationary advection-diffusion problem, where  $u \in X$  is the solution of

$$\mathbf{a} \cdot \nabla u - \nabla \cdot (\kappa \nabla u) + \sigma u = f, \quad (3.1)$$

where  $\mathbf{a} \in [\mathcal{L}^2(\Omega)]^d$  is the advection flow satisfying  $\nabla \cdot \mathbf{a} = 0$  on  $\Omega$ ,  $\kappa \in H^1(\Omega)$ ,  $\kappa(\mathbf{x}) > 0$  is the diffusivity coefficient,  $\sigma > 0 \in \mathcal{L}^2(\Omega)$  is a reaction coefficient and  $f \in \mathcal{L}^2(\Omega)$  is source function.

To solve this problem, we introduce a suitable FE space  $X_h \subset X$  and we write the variational formulation of equation (3.1), as find  $u_h \in X_h$  such that

$$\int_{\Omega} \mathbf{a} \cdot \nabla u_h v_h \, d\Omega + \int_{\Omega} \kappa \nabla u_h \cdot \nabla v_h \, d\Omega + \int_{\Omega} \sigma u_h v_h \, d\Omega = \int_{\Omega} f v_h \, d\Omega, \quad \forall v_h \in X_h. \quad (3.2)$$

This discrete problem is called advection (or convection) dominated when the Peclet number is larger than 1,

$$\text{Pe} = \frac{\|\mathbf{a}\|_{\infty} L}{2\kappa} > 1, \quad (3.3)$$

with  $L$  the characteristic length of the domain. This definition might be inappropriate when the nature of the flow  $\mathbf{a}$  brutally changes in the domain or when the mesh is not homogeneous. For that reason we usually prefer the local Peclet number

$$\text{Pe}_h = \frac{|\mathbf{a}|_p h}{2\kappa}, \quad (3.4)$$

where  $|\cdot|_p$  denotes the  $p$ -norm and  $h$  is the characteristic length of the current element. Using the local Peclet number we can define a problem as advection dominated as soon as its max on the domain is larger than 1,

$$\max_{\mathbf{x} \in \Omega} \text{Pe}_h(\mathbf{x}) > 1. \quad (3.5)$$

In the context of advection-dominated flow, the error bound on the Galerkin projection  $u_h$  does not guarantee a good approximation anymore. In practice the system becomes numerically unstable and non-physical oscillations appear. A trivial solution would be to refine the mesh to reduce the Peclet number, but this approach might be very costly and not achievable in practice. The alternative solution is to add a stabilization operator to take care of these numerical oscillations.

## 1.1 Stabilized Formulation

We first introduce the stabilized formulation, using the decomposition of the operators, as proposed in [43]. We denote by  $\mathcal{L}$  the advection diffusion reaction operator, defined as

$$\mathcal{L}u = \mathbf{a} \cdot \nabla u - \nabla \cdot (\kappa \nabla u) + \sigma u, \forall u \in X \quad (3.6)$$

and we also define its symmetric and skew-symmetric contributions

$$\mathcal{L}_S u = -\nabla \cdot (\kappa \nabla u) + \sigma u, \quad \mathcal{L}_{SS} = \mathbf{a} \cdot \nabla u, \quad \forall u \in X \quad (3.7)$$

respectively.

The stabilized formulation of equation (3.2) can be written as, find  $u_h \in X_h$  such that

$$a(u_h, v_h) + a_K(u_h, v_h) = l(v_h) + l_K(v_h), \quad \forall v_h \in X_h, \quad (3.8)$$

where for any  $(u, v) \in X^2$

$$a(u, v) = \int_{\Omega} \mathbf{a} \cdot \nabla uv \, d\Omega + \int_{\Omega} \kappa \nabla u \cdot \nabla v \, d\Omega + \int_{\Omega} \sigma uv \, d\Omega, \quad l(v) = \int_{\Omega} f v \, d\Omega, \quad (3.9)$$

and the stabilization terms  $a_K$  and  $l_K$  are defined, for any  $(u_h, v_h) \in X_h$ , such that

$$a_K(u_h, v_h) = \sum_{K \in \mathcal{T}_h} \int_K \tau_h \mathcal{L}u_h (\mathcal{L}_{SS}v_h + \gamma \mathcal{L}_S v_h) \, d\Omega, \quad (3.10)$$

$$l_K(v_h) = \sum_{K \in \mathcal{T}_h} \int_K \tau_h f (\mathcal{L}_{SS}v_h + \gamma \mathcal{L}_S v_h) \, d\Omega. \quad (3.11)$$

The computation of the stabilization parameter  $\tau_h$  is detailed in the following section. The constant  $\gamma$  determines the nature of the stabilization method:

- $\gamma = 0$  : SUPG,
- $\gamma = 1$  : GLS,
- $\gamma = -1$  : DW.

**Remark 3.** *The stabilization terms (3.10) are consistent with respect to the strong formulation of the problems. They vanish if we plug the exact solution in the equations.*

## 1.2 The Navier-Stokes Equations

The stationary Navier-Stokes system

$$\begin{cases} \rho \mathbf{u} \nabla \mathbf{u} - \nabla \cdot (\mu \nabla \mathbf{u}) + \nabla p = \mathbf{g} \\ \nabla \cdot \mathbf{u} = 0 \end{cases} \quad (3.12)$$

can be seen as a particular case of the generic advection diffusion reaction problem presented in former section by identifying  $f = g - \nabla p$ ,  $\mathbf{a} = \rho \mathbf{u}$ ,  $\kappa = \mu$  and  $\sigma = 0$ . We can then extend the stabilization operator described previously and write the stabilized variational formulation for NS system, as find  $(\mathbf{u}_h, p_h) \in V_h \times Q_h$  such that

$$\begin{aligned} & \int_{\Omega} \rho \mathbf{u}_h \nabla \mathbf{u}_h \cdot \mathbf{v}_h \, d\Omega + \int_{\Omega} \mu \nabla \mathbf{u}_h : \nabla \mathbf{v}_h \, d\Omega \\ & - \int_{\Omega} p_h \nabla \cdot \mathbf{v}_h \, d\Omega + \int_{\Omega} (\nabla \cdot \mathbf{u}_h) q_h \, d\Omega + a_K^{NS}(\mathbf{u}_h, \mathbf{v}_h; p_h) = \int_{\Omega} g \mathbf{v}_h \, d\Omega + l_K^{NS}(\mathbf{u}_h, \mathbf{v}_h) \end{aligned} \quad (3.13)$$

where

$$a_K^{NS}(\mathbf{u}_h, \mathbf{v}_h; p_h) = \sum_{K \in \mathcal{T}_h} \int_K \tau_h (\rho \mathbf{u}_h \nabla \mathbf{u}_h - \nabla \cdot (\mu \nabla \mathbf{u}_h) + \nabla p_h) \cdot (\rho \mathbf{u}_h \nabla \mathbf{v}_h - \gamma \nabla \cdot (\mu \nabla \mathbf{v}_h)) \, d\Omega \quad (3.14)$$

$$+ \sum_{K \in \mathcal{T}_h} \int_K \delta_h \nabla \cdot \mathbf{u}_h \nabla \cdot \mathbf{v}_h \, d\Omega \quad (3.15)$$

$$l_K^{NS}(\mathbf{u}_h, \mathbf{v}_h) = \sum_{K \in \mathcal{T}_h} \int_K \tau_h g \cdot (\rho \mathbf{u}_h \nabla \mathbf{v}_h - \gamma \nabla \cdot (\mu \nabla \mathbf{v}_h)) \, d\Omega \quad (3.16)$$

with  $\gamma$  and  $\tau_h$  defined in previous paragraph. We choose here to keep the term  $l_K^{NS}$  on the right-hand side since the value of the convection field in that term is usually taken at the previous linear iteration, see section 3 for more details on the non-linear resolution.

**Remark 4.** *We add a stabilization term on the divergence. This term is controlled by a new parameter  $\delta_h = 2|\mathbf{u}_h|_p^2 \tau_h$ , where  $|\cdot|_p$  is the  $p$ -norm as used in the definition of  $\tau_h$ , see next section. The divergence stabilization term is not mandatory, but we observed a clear improvement of the results when using it. In further sections, we suppose that this term is present in the stabilized formulation by default.*

### 1.3 Convergence of the Method

We will not detail the error analysis of the SDM in this manuscript. The convergence of this formulation has been studied, for instance, in [43, 44, 53]. We remind here the global result: for sufficiently smooth  $u$ , there exists a constant  $C_u$  such that

$$\|u_h - u\|^2 \leq C_u h^{2l} \quad (3.17)$$

where

$$\begin{aligned} 2l &= 2k + 1, & \alpha \text{ large} \\ 2l &= 2k, & \alpha \text{ small} \end{aligned} \quad (3.18)$$

with  $k$  the polynomial order of the finite element space and  $\alpha$  is the weight local Peclet number introduced in (3.26), see section 2.

We refer the reader to the cited references for a rigorous formulation of the hypothesis and a proof of this result.

## 2 Stabilization Parameter

The present stabilization methods are particularly sensitive to the choice of the parameter  $\tau_h$ . Many different definitions of this parameter can be found in the literature. We present in this section a non-exhaustive list of suitable parameters determined from different mathematical error analysis [45, 43, 61, 62]. We first detail the design of the parameter for advection-diffusion problem and then an extension of this definition for problems with reaction term. In the last subsection, we present different approaches for the delicate computation of the characteristic length of an element.

The required ingredient for the stabilization parameter is the evaluation of the local Peclet number,  $Pe_h$ , presented in (3.4). For the Navier Stokes equation, the Peclet number is replaced by the Reynolds number but definition and roles are identical,

$$Re_h = \frac{|\mathbf{u}|_p h}{2\nu}. \quad (3.19)$$

These quantities mainly depend on the characteristic size  $h$  and have real impacts on the efficiency of the method. For now, we do not specify the form of the length  $h$ . We will focus on this delicate detail in the subsection 2.2.

The value of the local Peclet/Reynolds number will determine the type of problem we are dealing with. Small Peclet/Reynolds numbers correspond to diffusive problems and are usually stable. Large Peclet/Reynolds numbers define convection dominated cases and require adapted stabilization. For the rest of the section, we use Peclet number and associated notations  $(\mathbf{a}, \kappa)$  in the formulas. The computation of stabilization parameter is equivalent for Navier Stokes equations, by replacing the corresponding quantities  $Pe_h \sim Re_h$ ,  $\mathbf{a} \sim \mathbf{u}$  and  $\kappa \sim \nu$ .

The stabilization parameter  $\tau_h$  has to satisfy the asymptotic behaviors

$$\tau_h = O\left(\frac{h}{|\mathbf{a}|}\right), \quad Pe_h \text{ large}, \quad \tau_h = O\left(\frac{h^2}{\kappa}\right), \quad Pe_h \text{ small} \quad (3.20)$$

In [8], Hughes *et al.* determine the optimal stabilization parameter, for 1-D problem with linear FE, as

$$\tau_h^0 = (\coth(Pe_h) - \frac{1}{Pe_h}). \quad (3.21)$$

This formula has been generalized for multi-dimensional problems but requires some modifications for higher FE polynomial order. An extension of this formula

was proposed in [45, 43], with the introduction of a function  $\xi(\alpha)$  to locally switch between diffusive and convective problems. We can then define the stabilization parameter as

$$\tau_h^1 = \frac{h}{2|\mathbf{a}|_2} \xi(\alpha). \quad (3.22)$$

where  $\alpha$  is a weighted version of  $Pe_h$  and  $\xi(\alpha)$  has the properties described in figure 3.1, with  $m \leq 2C$  and  $C$  is a constant satisfying the inverse estimate inequality

$$C \sum_{K \in \mathcal{T}_h} h^2 \|\Delta v\|_K^2 \leq \|\nabla v\|^2, \quad v \in X_h \quad (3.23)$$

where  $\|\cdot\|_K$  is the  $\mathcal{L}^2$ -norm in the cell  $K$ . The practical evaluation of the constant

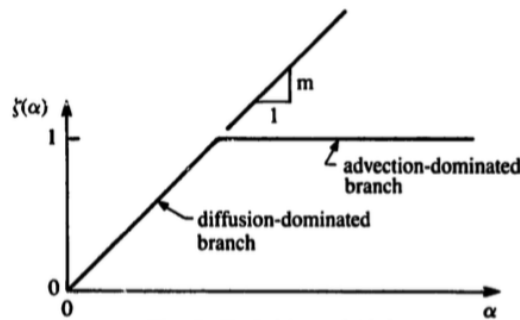


Fig. 2. Definition of  $\xi(\alpha)$ .

Figure 3.1 – Definition of  $\xi(\alpha)$ , from [43]

$C$  is explained in [48] : some basic values of  $C$  are given for particular shapes and polynomial degrees, for instance  $C = 1/24$  for biquadratic rectangular elements. A general method is also proposed in order to evaluate  $C$  in any cases using the solution of the local eigen value problem

$$\lambda_K = \max_{v \in (\mathbb{P}^k(K) \setminus \mathbb{R})} \frac{\|\Delta v\|_K^2}{\|\nabla v\|_K^2} \quad (3.24)$$

we can then define

$$C = \frac{1}{\lambda_K h^2}. \quad (3.25)$$

With this definition,  $C$  is a constant independent of the shape of the cell as long as we use the definition of  $h$  given in [48], this definition is detailed in the subsection 2.2. We actually figured out that  $C$  has to be a constant on the whole mesh. This property is particularly important on anisotropic meshes. To recover this constant value we always use the  $h_h$  definition (3.45) to compute  $C$ . Then, we can use any definition for the characteristic length to compute the local Peclet number.

Hughes *et al.* proposed different formulas for  $\xi(\alpha)$  first in [45] and [43] and finally in [61, 62]. This last version is more convenient to use and became the most common formulation, we recall here the main ingredients:

$$\xi(\alpha) = \begin{cases} \alpha, & 0 \leq \alpha \leq 1, \\ 1, & \alpha \geq 1, \end{cases} \quad (3.26)$$

$$\alpha = \frac{m|a|_p h}{2\kappa}, \quad (3.27)$$

$$|a|_p = \begin{cases} \left( \sum_{i=1}^d |a_i|^p \right)^{1/p}, & 1 \leq p < \infty, \\ \max_{i=1, \dots, d} |a_i|, & p = \infty, \end{cases} \quad (3.28)$$

$$m = \min\left\{\frac{1}{3}, 2C\right\}. \quad (3.29)$$

with  $C$  defined as in (3.23) and then  $\tau$  is defined as in (3.22).

**Remark 5.** *The interest of the weighted Peclet/Reynolds number  $\alpha$  is to take into account the effect of the polynomial FE degree. Indeed, polynomial degree  $k \geq 2$  has non-negligible effects on local Peclet/Reynolds number since it highly improves the quality of the discrete approximation.*

We will compare this definition of  $\tau_h^1$  with a formula proposed by Burda *et al.* in [10] using the approximation  $h \approx \frac{1}{\sqrt{\lambda_K}}$  which also takes into account the influence of the polynomial order,

$$\tau_h^2 = \frac{\xi(\alpha)}{\sqrt{\lambda_K} |a|_p}, \quad (3.30)$$

$$\xi(\alpha) = \begin{cases} \alpha, & 0 \leq \alpha \leq 1, \\ 1, & \alpha \geq 1, \end{cases} \quad (3.31)$$

$$\alpha = \frac{|a|_p}{2\kappa\sqrt{\lambda_K}}, \quad (3.32)$$

$$(3.33)$$

where  $\lambda_K$  is defined as in (3.24) and  $|\cdot|_p$  is the usual  $p$ -norm defined in (3.26).

## 2.1 Extension to Advection Diffusion Reaction Problems

The parameter  $\tau_h^1$  defined in (3.22) is very efficient for advection dominated flows, but he does not take care of the potentially high reaction coefficient  $\sigma$ . It is not an an issue for the resolution of the aerothermal model presented in chapter I but we met some difficulties with the stabilization of the turbulence models equations, see chapter 4.

An extension of  $\tau_h^1$  for problem with reaction term is proposed in [33], with a new stabilization parameter  $\tau_h^3$  defined by

$$\tau_h^3 = \frac{h^2}{\sigma h^2 \xi(\alpha_\sigma) + \frac{4\kappa}{m} \xi(\alpha_\kappa)}, \quad (3.34)$$

$$\alpha_\sigma = \frac{2\kappa}{m\sigma h^2}, \quad (3.35)$$

$$\alpha_\kappa = \frac{m|\mathbf{a}|_p h}{2\kappa}, \quad (3.36)$$

$$\xi(\alpha) = \begin{cases} 1, & 0 \leq \alpha \leq 1, \\ \alpha, & \alpha \geq 1, \end{cases} \quad (3.37)$$

$$|\mathbf{a}|_p = \left( \sum_{i=1}^d |a_i|^p \right)^{1/p}, \quad 1 \leq p < \infty. \quad (3.38)$$

We did minor changes in the constants of  $\tau_h^3$  and  $\alpha_\kappa$  in order to recover exactly the definition of  $\tau_h^1$  when  $\sigma = 0$ .

## 2.2 Choice of the characteristic length

The definition of the characteristic size of a cell is a key ingredient in the computation of a stabilization parameter. This choice is particularly crucial for anisotropic meshes with important aspect ratios. For isotropic meshes, most of the following solutions will produce acceptable results.

We want to compare these different options for highly anisotropic meshes, as used in the turbulent boundary layers. The same kind of study has been made by Mittal, [73] for the stabilization of equal-order-interpolation velocity-pressure for Navier-Stokes equations.

**$h_m$  definition :**  $h = h_m$  is the length of the smallest edge of the cell.

**$h_d$  definition :** We use the projection of the current convection field  $\mathbf{a}$  on the directions of the cell. This definition was introduced by Hughes in [42] and is particularly convenient for anisotropic meshes. We first define the local convection field  $\hat{\mathbf{a}}$  on the reference element,  $\hat{K}$ :

$$\hat{\mathbf{a}} = \mathbf{J}^{-1} \mathbf{a} \quad (3.39)$$

where  $\mathbf{J}$  is the Jacobian of the geometric mapping from the reference element to the current cell. We can then define

$$h = h_d = \frac{|\mathbf{a}|_p \hat{h}}{|\hat{\mathbf{a}}|_p}, \quad (3.40)$$

where  $\hat{h}$  is the characteristic length of the reference element. This definition works particularly well on hypercubes, but we noticed some regularity issue when we applied this method on triangles. We finally figured out that these irregularities came from the shape of the reference cell. In Feel++, the reference simplex (in 2D) is a squared isosceles triangle. This shape is not symmetric by rotation and then the advection field  $\hat{\mathbf{a}}$  would depend on the rotation between the current cell and the reference element. To overcome this problem, we decided to refer  $\hat{\mathbf{a}}$  to a new reference triangle  $\tilde{K}$  which is equilateral. An illustration of this procedure



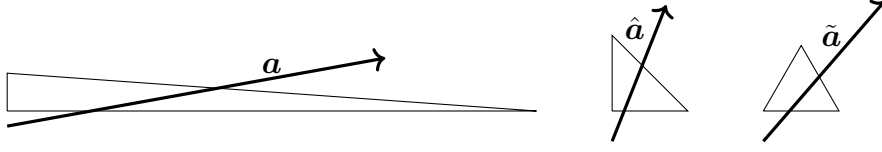


Figure 3.2 – Geometric transformation back on a reference element  $\tilde{K}$ .  
Computation of the local field  $\tilde{a}$ , used for the characteristic size  $h_d$ .

is proposed in figure 3.2. In our case the geometric mapping from  $\hat{K}$  to  $\tilde{K}$  can be written as

$$\tilde{X} = M\hat{X} \quad (3.41)$$

with

$$M = \begin{bmatrix} 1 & 1/2 \\ 0 & \sqrt{3}/2 \end{bmatrix}. \quad (3.42)$$

Then we can use the matrix  $M$  to report the field  $\hat{a}$  on  $\tilde{K}$ .

$$\tilde{a} = M\hat{a} = MJ^{-1}a \quad (3.43)$$

and finally the characteristic length is computed as

$$h = h_d = \frac{|\mathbf{a}|_p}{|\tilde{\mathbf{a}}|_p} \tilde{h}. \quad (3.44)$$

$h_h$  **definition** : Harari and Hughes proposed this definition in [48] and provided an interesting measure for anisotropic meshes. The idea is to compute the size of the cell using the area of the element and the average vertex-to-centroid distance. The formulation for triangle is

$$h_h = \frac{4S}{\sqrt{3 \sum_{a=1}^3 |x_i - x_c|^2}} \quad (3.45)$$

where  $S$  is the surface of the cell,  $x_i$  are the vertex of the triangle and  $x_c$  the centroid of the triangle. We have a similar definition for quadrangle elements

$$h_h = h_x h_y \sqrt{\frac{2}{h_x^2 + h_y^2}} \quad (3.46)$$

where  $h_x$  and  $h_y$  are the characteristic length of the quadrangle in both direction.

$h_e$  **definition** : an almost optimal definition for anisotropic elements would be to compute the inscribed ellipse for each element and then to take as  $h = h_e$  the radius of this ellipse in the direction of the advection field  $a$ . Even if this approach is interesting in theory, it appears impracticable to implement and we did not take time to study this version.

### 3 Shock Capturing Method

The SDM are particularly efficient to stabilize convection-dominated problems, however these methods also generate significant overshoots and undershoots in regions of high gradient. An example of this non-physical oscillations can be seen in figure 3.3

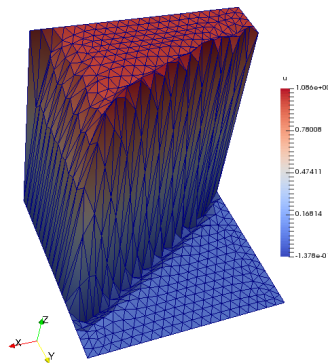


Figure 3.3 – Non-physical oscillations resulting from supg stabilization on the advection skew mesh test case 4.2

Different methods have been developed in the last decades to eliminate or at least diminish these oscillations. We particularly noted three families of methods. The non-linear upwinding techniques, introduced by Mizukami and Hughes [74], is very effective to remove the oscillations but is known to introduce too much numerical diffusion and then a loss of accuracy. Another solution is to add artificial diffusion. This diffusion can be either isotropic [45], orthogonal to the streamlines [52] or based on an edge stabilization [13]. In their implementation of a  $k - \varepsilon$  model, Kuzmin et al. [59] proposed an alternative algebraic solution, by applying an iterative flux correction. Note that all these methods are non-linear.

In this work, we chose to focus on the artificial diffusion methods. These methods are convenient to implement once the SUPG/GLS framework is in place. Also, they succeed in being very effective with no excessive artificial diffusivity. We implemented different versions proposed in the literature, and we give some details in this section.

#### 3.1 Isotropic Artificial Diffusion

Hughes *et al.*[45] first proposed to add a new term in the SUPG stabilization

$$(R_h(\mathbf{u}_h), \tau_h^{sc} \mathbf{a}^{\parallel} \cdot \nabla \mathbf{v}_h) \quad (3.47)$$

where  $R_h(\mathbf{u}_h)$  is the discrete residual,  $\mathbf{a}^\parallel$  is the projection of the convection field  $\mathbf{a}$  in the direction of  $\nabla \mathbf{u}_h$  and  $\tau_h^{sc}$  is a non negative stabilization parameter. Many definitions and improvements were developed after this first version, we remind hereafter the most common versions.

Galeão and do Carmo [28] proposed to replace the projected convection field  $\mathbf{a}^\parallel$  by the new vector field  $\mathbf{z}_h = R_h(\mathbf{u}_h) \frac{\nabla \mathbf{u}_h}{\|\nabla \mathbf{u}_h\|_2^2}$ . The new additional term then reads

$$(R_h(\mathbf{u}_h), \tau_h^{sc} \mathbf{z}_h \cdot \nabla \mathbf{v}_h), \quad (3.48)$$

with the following definition for  $\tau_h^{sc}$

$$\tau_h^{sc} = \tau_h \max \left( 0, \frac{\|\mathbf{a}\|_2}{\|\mathbf{z}_h\|_2} - 1 \right), \quad (3.49)$$

where  $\tau_h$  is the SUPG/GLS stabilization parameter as defined in section 2.

Almeida *et al.*[1] also use the formulation (3.48) but they introduce a new quantity  $\zeta_h$  in their definition of  $\tau_h^{sc}$  to reduce the artificial diffusion in regions where the numerical solution is smooth enough.

$$\tau_h^{sc} = \tau_h(\mathbf{a}) \max \left( 0, \frac{\|\mathbf{a}\|_2}{\|\mathbf{z}_h\|_2} - \zeta_h \right), \quad \text{with} \quad \zeta_h = \max \left( 1, \frac{\mathbf{a} \cdot \nabla \mathbf{u}_h}{R_h(\mathbf{u}_h)} \right). \quad (3.50)$$

### 3.2 Artificial Diffusion Added Orthogonally to Streamlines

This method, introduced by Johnson *et al.*[52] is very similar to the isotropic diffusion, but it is only applied in the direction orthogonal to Streamlines adding the term

$$(\tau_h^{sc} \bar{D} \nabla \mathbf{u}_h, \nabla \mathbf{v}_h) \quad (3.51)$$

where  $\bar{D} = \mathbb{I} - \mathbf{a} \otimes \frac{\mathbf{a}}{\|\mathbf{a}\|_2^2}$  is the projection orthogonally to vector  $\mathbf{a}$  and  $\tau_h^{sc}$  is defined as  $\tau_h^{sc} = \max \left( 0, \|\mathbf{a}\|_2 h^{3/2} - \kappa \right)$ .

Codina [18] gives a new definition for  $\tau_h^{sc}$ , based on the validation of the discrete maximum principle for several problems,

$$\tau_h^{sc} = \frac{1}{2} \max \left( 0, C - \frac{2\kappa}{\|\mathbf{a}\|_2 h} \right) h \frac{R_h(\mathbf{u}_h)}{\nabla \mathbf{u}_h}, \quad (3.52)$$

where  $C$  is a suitable positive constant (proposed constant are  $C = 0.6$  for linear elements and  $C = 0.35$  for quadratic elements).

Knobloch *et al.* proposed in [51] to replace  $|\mathbf{a}^\parallel|$  by the quantity  $Q_K(\mathbf{u}_h) = \frac{\|R_h(\mathbf{u}_h)\|_2}{\|\nabla \mathbf{u}_h\|_2}$  which leads to the formulation

$$\tau_h^{sc} = \frac{1}{2} \max \left( 0, C - \frac{2\kappa}{Q_K(\mathbf{u}_h) h} \right) h Q_K(\mathbf{u}_h). \quad (3.53)$$

In the same reference, Knobloch *et al.* also proposed a modified version of the shock capturing terms presented by Burman and Ern in [12] with

$$\tau_h^{sc} = \tau_h(\mathbf{a}) \frac{\|\mathbf{a}\|_2^2 \|R_h(\mathbf{u}_h)\|_2}{\|\mathbf{a}\|_2 \|\nabla \mathbf{u}_h\|_2 + \|R_h(\mathbf{u}_h)\|_2} \quad (3.54)$$

Considering all these methods and the results of section 4, we propose to use the shock capturing term (3.51) but with the new parameter

$$\tau_h^{sc} = \tau_h(\mathbf{a}) \max(0, \|\mathbf{a}\|_2^2 Q_K(\mathbf{u}_h) - Q_K^2(\mathbf{u}_h)). \quad (3.55)$$

We had no time to perform a proper convergence study for this parameter, but the numerical results presented in further sections are clearly encouraging.

## 4 Validation

We have implemented and compared the previously introduced methods on different well-known test cases. We proposed in this section a validation and a comparison of the different configurations. We tested the different stabilization methods : SUPG and GLS, with parameters 1 (3.26), 2 (3.30) and possibly 3 (3.34). We use the subscript  $m, d$  or  $h$  for the different  $h$  formulations proposed in section 2.2. We compare those methods with standard Galerkin method (SGM) when it makes sense. All the simulations are made on triangular meshes with order 2 FE.

### 4.1 Comparison of the Stabilization Methods

#### Advection Skew the Mesh

We propose here a qualitative observation of the results computed with different configurations, on the well known advection skew the mesh example, [8]. The advection field is constant  $\mathbf{a} = (\cos(\alpha), \sin(\alpha))$ , the diffusivity coefficient is  $\kappa = 10^{-8}$ . There is no reaction and no source term. The test case consists in the propagation of a discontinuity in quasi pure advection setup. The domain is a square  $\Omega = [0, 1] \times [0, 1]$  with the boundary conditions:  $u(x, y) = 1$  on  $\Gamma_1$  and  $u(x, y) = 0$  on  $\Gamma_0$ , with  $\Gamma_1 = \{(x, y), x = 0, 0 \leq y \leq 0.2\} \cup \{(x, y), y = 0\}$  and  $\Gamma_0 = \partial\Omega \setminus \Gamma_1$ . This setup is recapped on figure 3.4.

We compare the two methods (SUPG and GLS) for different values of  $\alpha$ ,  $\alpha = \text{atan}(0.5)$  on figure 3.5,  $\alpha = \text{atan}(1)$  on figure 3.6 and  $\alpha = \text{atan}(2)$  on figure 3.7. Those quantitative results allow some observations:

- All the configurations succeed in stabilizing the model where the SGM does not converge. We note the apparition of non-physical undershoot and overshoot near the discontinuities.

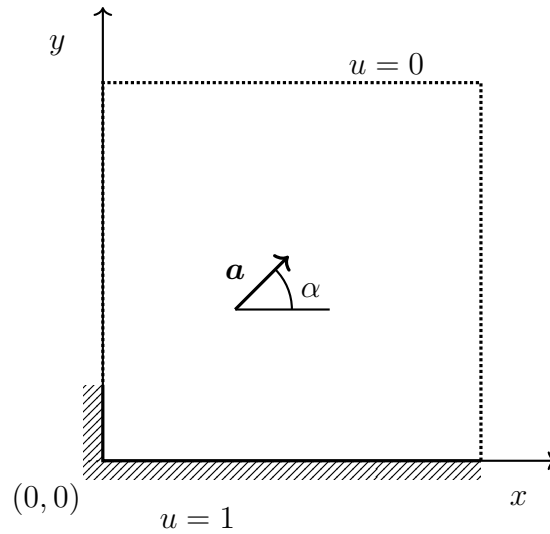


Figure 3.4 – Advection Skew the Mesh Setup

- The two methods (SUPG and GLS) give almost identical results and the stabilization parameter seems to have more influence than the method.
- The parameters  $\tau_h^2$  is more sensitive to discontinuities and boundaries.
- For  $\tau_h^1$ , the choice of the characteristic length has only slight influence.

### Thermal Boundary Layer

In this test case, we want to compare the different stabilization parameters for advection-diffusion-reaction problems and especially their accuracy on non-homogeneous anisotropic meshes.

We consider a rectangular domain  $\Omega = [0, 1] \times [0, 0.5]$ , the boundary conditions are :

$$\begin{aligned}
 u &= 0, & 1 \leq x \leq 1, y = 0, \\
 u &= 1, & x = 0, 0 \leq y \leq 0.5, \\
 u &= 1, & 0 \leq x \leq 1, y = 0.5 \\
 u &= 2y, & x = 1, 0 \leq y \leq 0.5
 \end{aligned} \tag{3.56}$$

On this geometry, we do consider two configurations, see figure 3.9. The first one with a coarse non-structured mesh, with a diffusion  $\kappa = 10^{-5}$ . The second case is set on an non-homogeneous anisotropic structured mesh, with  $\kappa = 10^{-8}$ . This second mesh presents high mesh ratio on the bottom boundary (up to  $10^5$ ). The stabilization on this kind of stretched elements is particularly complicated and requires a suitable definition of the characteristic length  $h$ . For both configurations, the advection field is defined as  $\mathbf{a} = (2y, 0)$ . The interest of using the advection-diffusion-reaction parameter  $\tau_h^3$  will be discussed with different values for the reaction parameter  $\sigma$ .

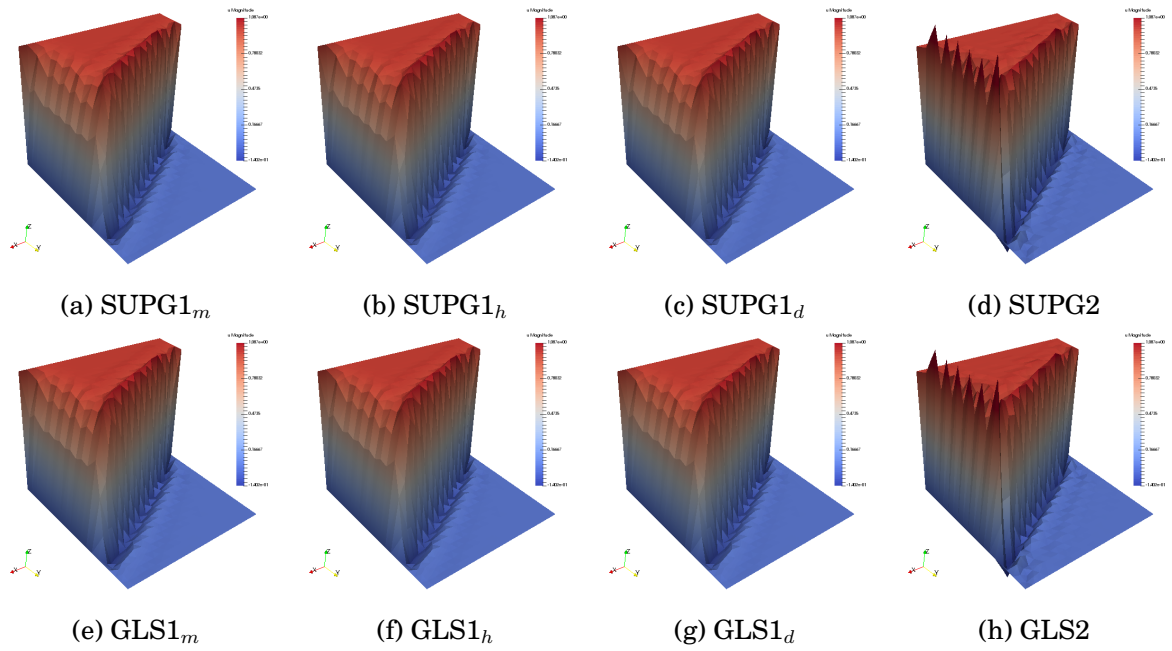


Figure 3.5 – Advection Skew the Mesh profiles for  $\alpha = \text{atan}(0.5)$ . Comparison of the two methods with different parameters and characteristic length

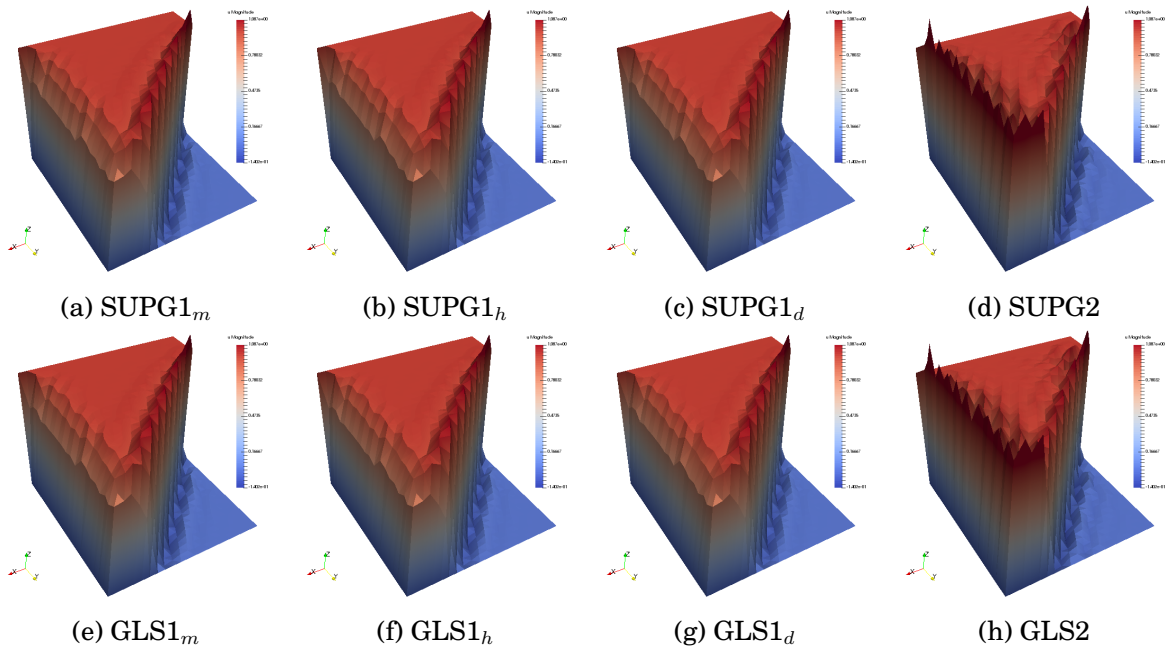


Figure 3.6 – Advection Skew the Mesh profiles for  $\alpha = \text{atan}(1)$ . Comparison of the two methods with different parameters and characteristic lengths

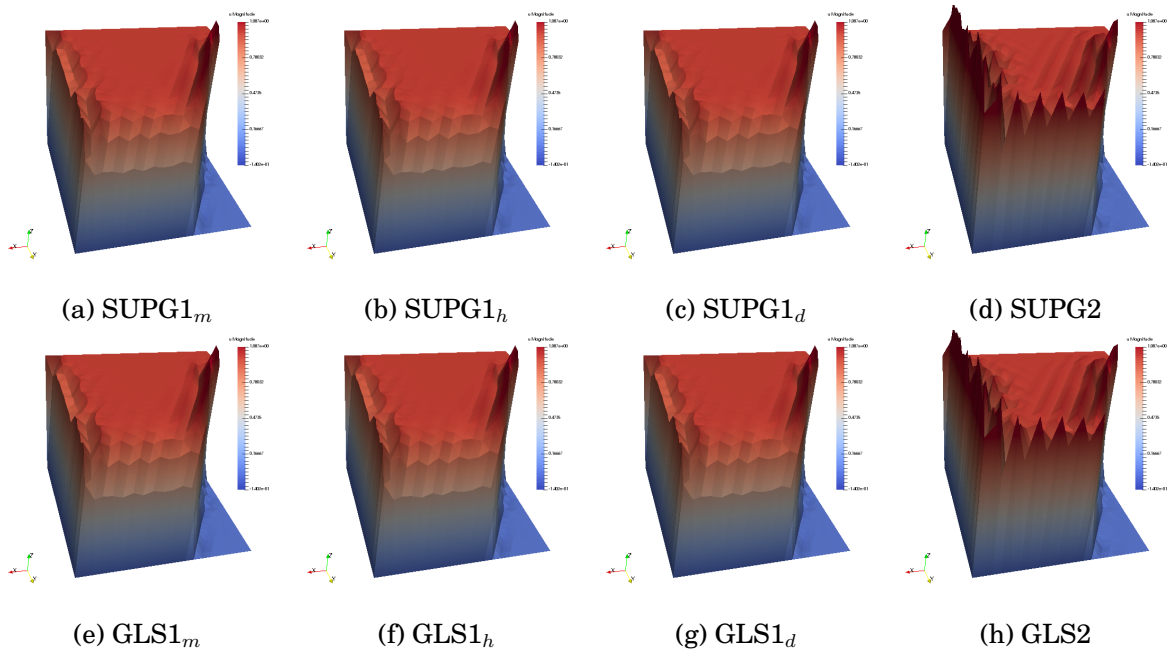


Figure 3.7 – Advection Skew the Mesh profiles for  $\alpha = \text{atan}(2)$ . Comparison of the two methods with different parameters and characteristic lengths

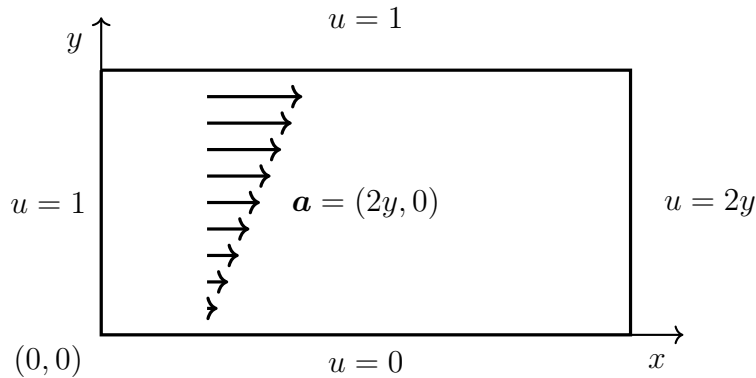


Figure 3.8 – Thermal Boundary Layer Setup

We first compare the results obtained with different methods, parameters, and characteristic lengths. This study is made on the unstructured homogeneous configuration. The setup presents sharp discontinuities and is not trivial to stabilize. The profiles are presented in figure 3.10. We also propose three profile cuts ( $x = 0.2$ ,  $x = 0.5$  and  $x = 0.8$ ) to accurately compare the methods, see figure 3.11. We globally retrieve the same conclusions we made in the previous test case: all configurations stabilize, with more oscillations for  $\tau_h^2$  parameter. From those first test cases, we can reasonably conclude that the choice of the characteristic length will not influence the stabilization method, at least on isotropic meshes. Also, the underlying method, SUPG or GLS, has very little influence on the solution. For those reasons, we will usually choose SUPG formulation (which is a bit less expensive to assemble) and the  $\tau_h^1$  parameter with  $h_m$  definition for advection-diffusion on isotropic meshes.

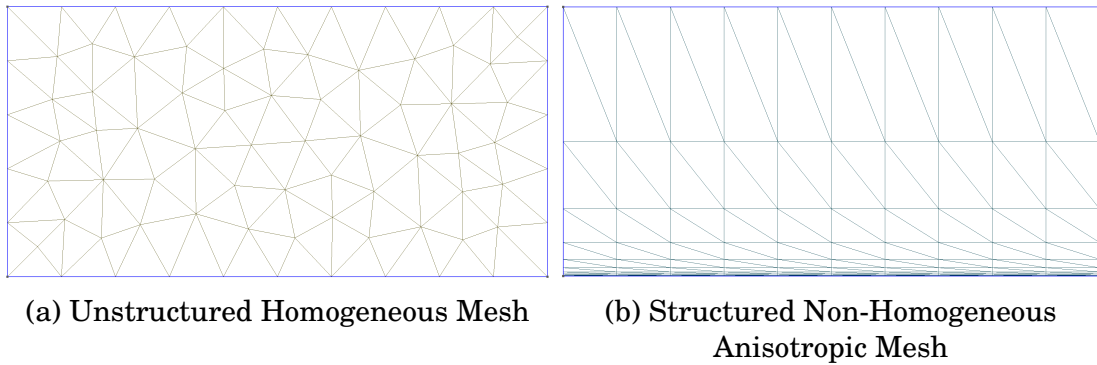


Figure 3.9 – Meshes Used for the two Different Configurations of the Thermal Layer Test Case

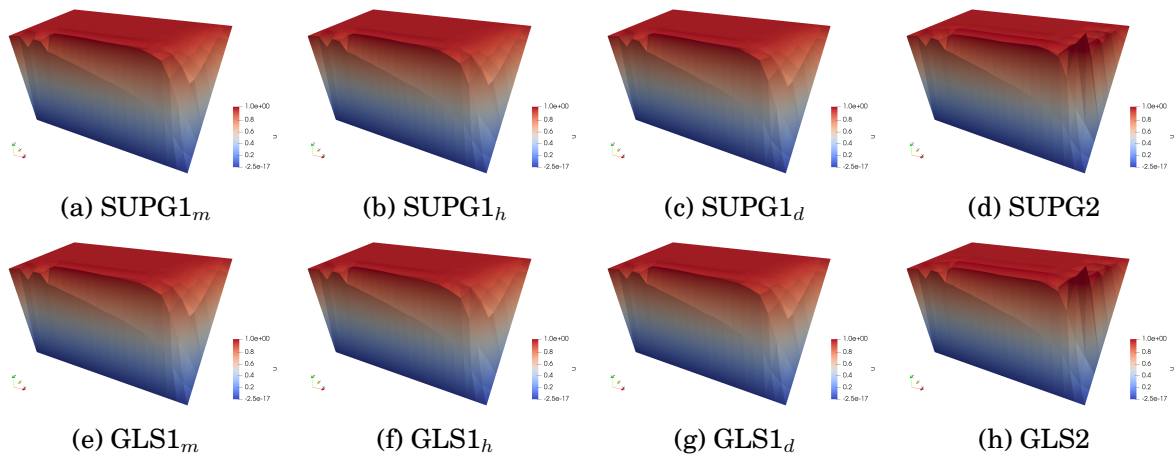


Figure 3.10 – Thermal Layer on Coarse Unstructured Homogeneous Mesh.  $\sigma = 0$ ,  $\kappa = 10^{-5}$ . Comparison of the two methods with different parameters and characteristic length

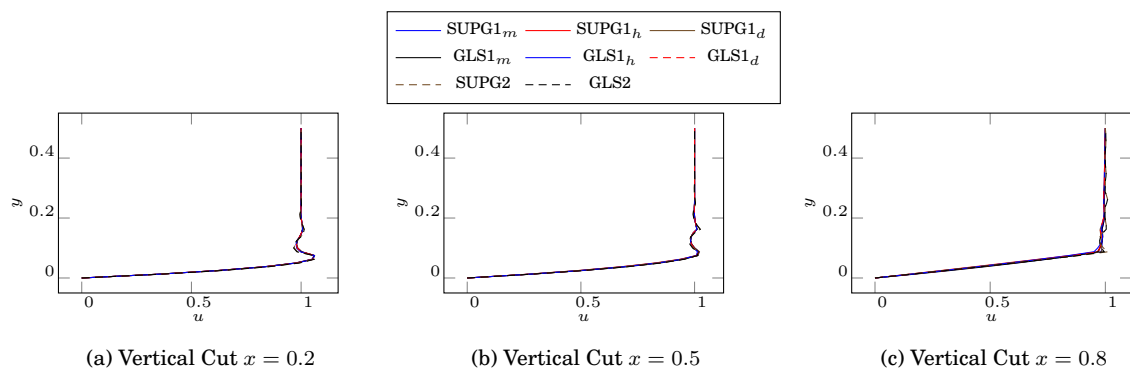


Figure 3.11 – Thermal Layer on Coarse Unstructured Homogeneous Mesh.  $\sigma = 0$ ,  $\kappa = 10^{-5}$ . Comparison of the two methods with different parameters and characteristic length



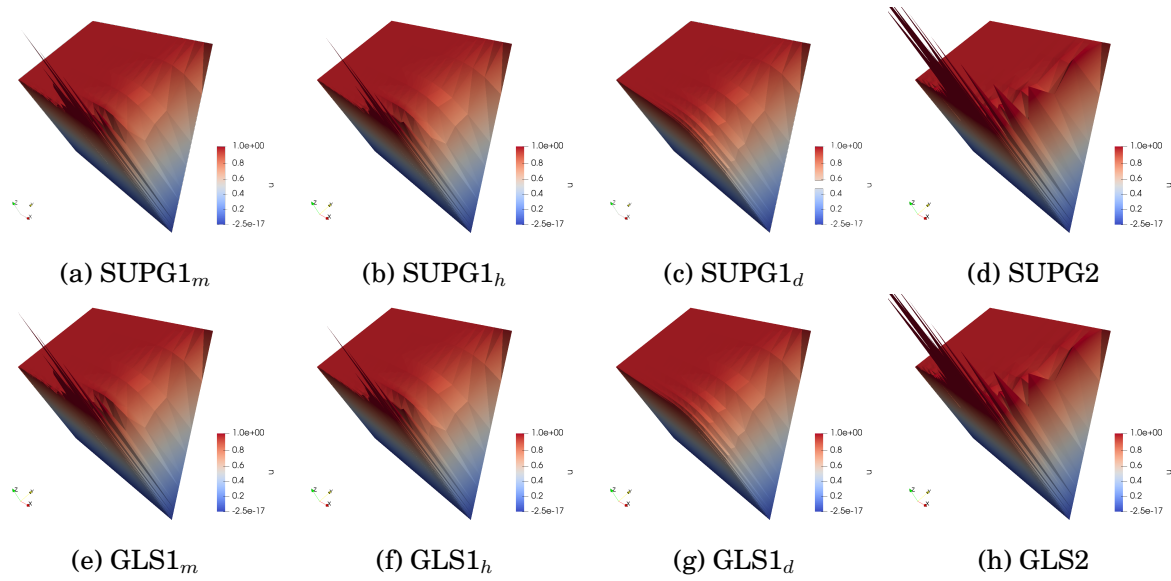


Figure 3.12 – Thermal Layer on Structured Anisotropic Mesh.  $\sigma = 0$ ,  $\kappa = 10^{-8}$ . Comparison of the two methods with different parameters and characteristic length

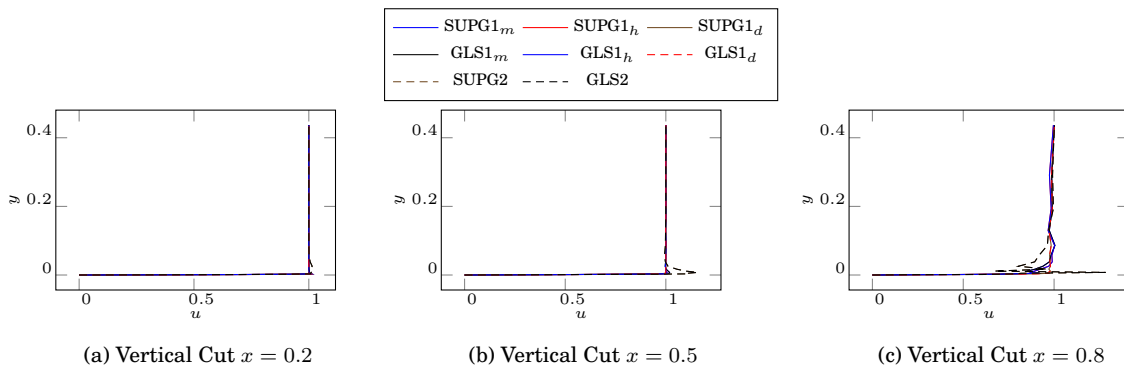


Figure 3.13 – Thermal Layer on Structured Anisotropic Mesh.  $\sigma = 0$ ,  $\kappa = 10^{-8}$ . Comparison of the two methods with different parameters and characteristic length

The second configuration proposed on this geometry is setup to illustrate the importance of the characteristic length definition for anisotropic meshes. On figures 3.12 and 3.13, we notice boundary oscillation for  $\tau_h^2$  and  $\tau_h^1$  with  $h_m$  and  $h_h$  versions. The  $h_d$  definition is the most accurate for anisotropic mesh since it takes into account the direction and the magnitude of the advection field. We will prefer this definition for problems with boundary layers. In this case again, the choice of the method between SUPG and GLS does not significantly influence the results.

As a last configuration on the thermal layer problem, we now add reaction term,  $\sigma = 1$ . This last study aims to discuss the interest of the parameter  $\tau_h^3$  which takes into account the value of the reaction coefficient  $\sigma$ . Since the parameter  $\tau_h^2$  is significantly less accurate than  $\tau_h^1$  we chose to only perform comparisons with this latter parameter. In both cases, we set  $h = h_d$  which is the most robust definition. The comparison is made on the unstructured homogeneous mesh with  $\kappa = 10^{-5}$ , figures 3.14 and 3.16, and also on the structured one with  $\kappa = 10^{-8}$ , figures 3.15

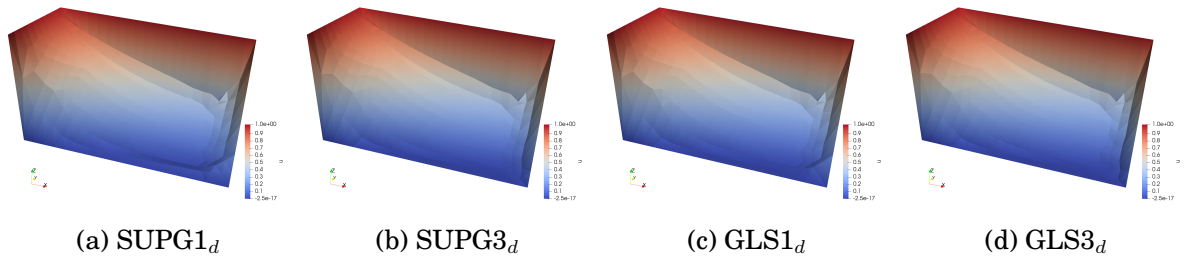


Figure 3.14 – Thermal Layer on Coarse Unstructured Homogeneous Mesh.  $\sigma = 1, \kappa = 10^{-5}$ . Comparison of the two methods with different parameters and characteristic length

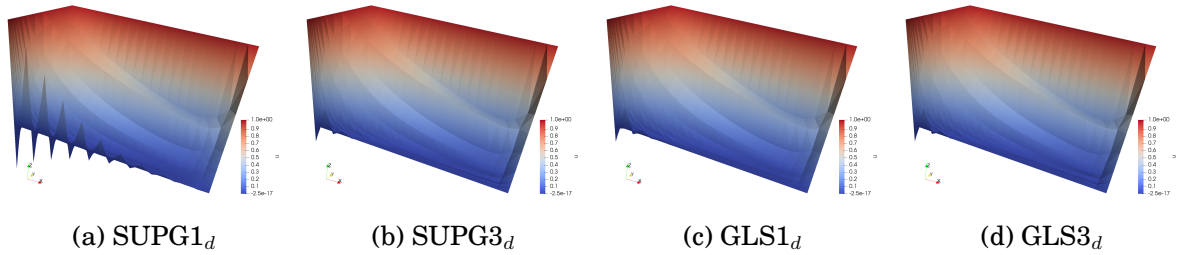


Figure 3.15 – Thermal Layer on Structured Anisotropic Mesh.  $\sigma = 1, \kappa = 10^{-8}$ . Comparison of the two methods with different parameters and characteristic length

and 3.17. On the homogeneous mesh, both parameters and both methods have comparable results. Nevertheless, we can denote a slight smoothing of the solution with parameter  $\tau_h^3$ . On the structured mesh, the activation of the reaction generates non-physical oscillations on the boundaries. Those oscillations are significantly bigger in the SUPG1<sub>d</sub> configuration. Using either GLS instead of SUPG or  $\tau_h^3$  instead of  $\tau_h^1$  reduces this numerical issue. However, we never succeeded in completely removing this annoying side effect. After many efforts, the best results we obtained are those presented for GLS3<sub>d</sub>. Note that the use of shock capturing methods does not improve the results.

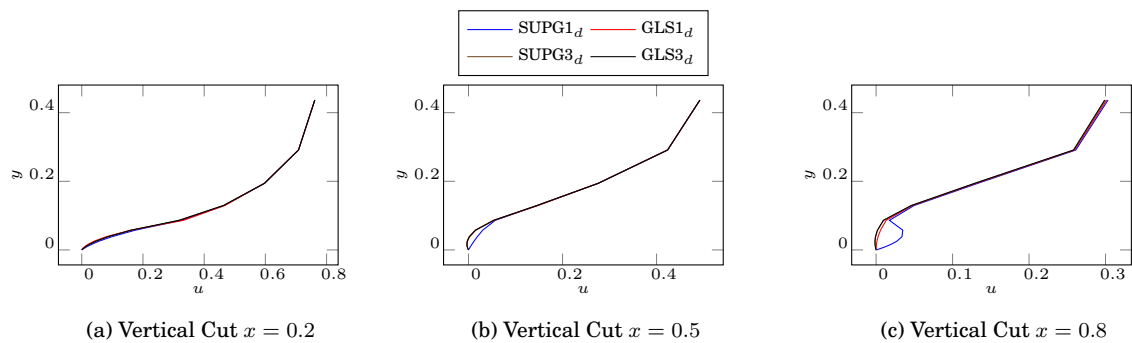


Figure 3.16 – Thermal Layer on Coarse Unstructured Homogeneous Mesh.  $\sigma = 1, \kappa = 10^{-5}$ . Comparison of the two methods with different parameters and characteristic length

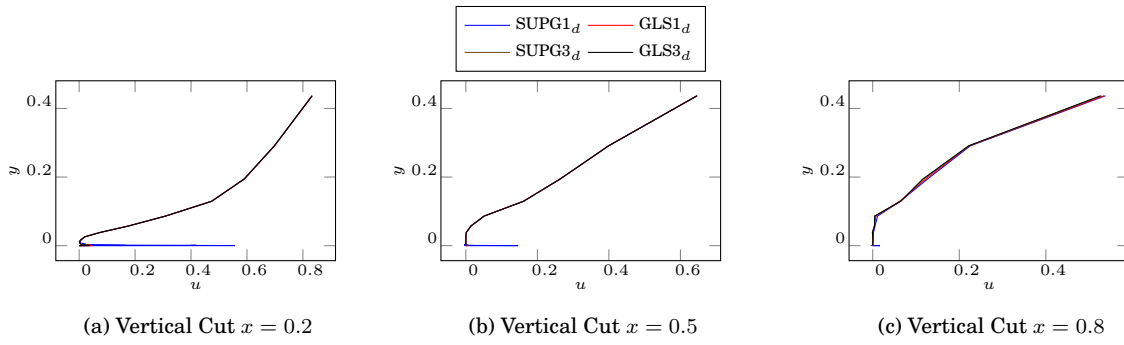


Figure 3.17 – Thermal Layer on Structured Anisotropic Mesh.  $\sigma = 1$ ,  $\kappa = 10^{-8}$ . Comparison of the two methods with different parameters and characteristic length

## 4.2 Shock Capturing Methods

The different methods tested in the present subsection are: ID-GC, from Galeão and Carmo (3.48) and (3.49), ID-A from Almeida *et al.* (3.50), OD-CK from Codina, modified by Knobloch *et al.* (3.53) and OD-K proposed by Knobloch *et al.* (3.54). We also present the results of our HOME shock capturing parameter, (3.55). According to the previous results, the choice of method and characteristic length definition does not have real influence on the solution, for this test case. All simulations are then performed with the method denoted by SUPG1<sub>d</sub>.

### Advection Skew the Mesh

We are considering again the advection skew the mesh test-case presented in section 4.1. The aim of the shock capturing methods is the reduction of the undershoots and overshoots produced by streamline diffusion methods, as identified in figure 3.3.

Starting with a qualitative observation of the profiles on figure 3.18 we notice reduction of the under/overshoots. To have a better idea, we propose a cut of these profiles at  $x = 0.5$ , see figure 3.19. All the methods are smoothing the profile and the under/overshoot entirely disappears except for OD-CK where the oscillations are only reduced.

### Advection in a Rotating Flow Field

We observed in the latter test case, the efficiency of the different methods. Nevertheless, those methods are reducing the non-physical oscillations by adding diffusivity. We now have to ensure that the considered methods are not too diffusive. For that, we compare their results on the advection in a rotating field test case.

The problem is defined on a unit square ( $-0.5 \leq x, y \leq 0.5$ ), see figure 3.20 with

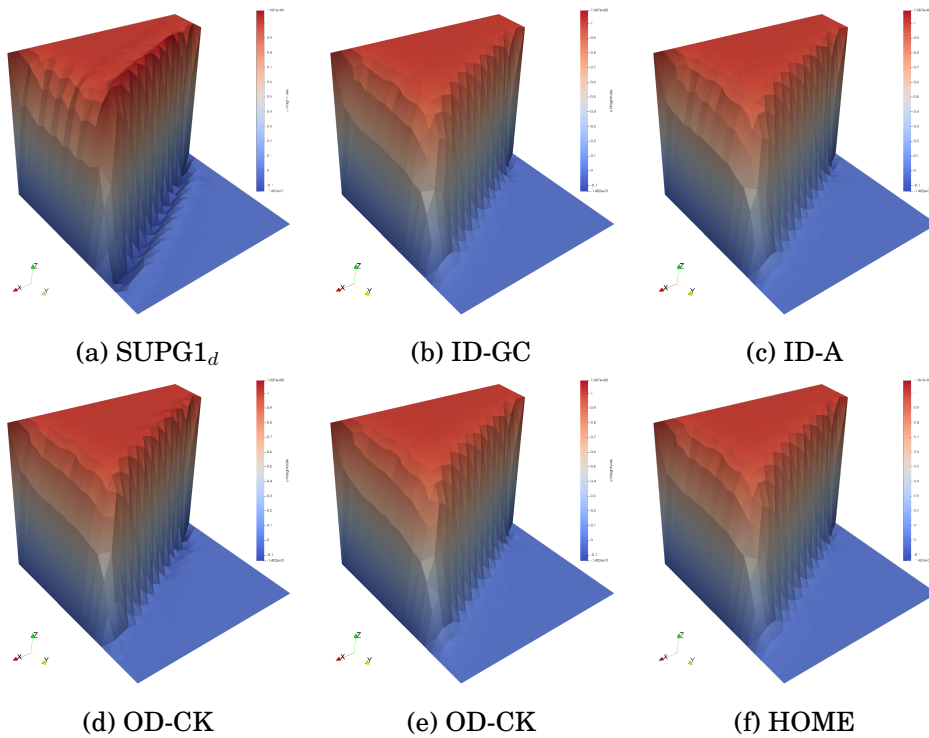


Figure 3.18 – Advection Skew the Mesh profiles for  $\alpha = \text{atan}(0.5)$ . Comparison of the different shock capturing methods, using SUPG $1_d$  as stabilization method.

a velocity flow given by  $a = (-y, x)$ . The boundary conditions are given by  $u = 0$  on all the external boundaries and we add an internal boundary condition with

$$u = \frac{1}{2}(\cos(4\pi y + \pi) + 1) \quad (3.57)$$

on the segment  $[OA]$ , defined by  $x = 0$ ,  $-0.5 \leq y \leq 0.5$ . The diffusivity is set to  $\kappa = 10^{-8}$ .

With this setup, the problem is almost purely advective, and so we should not see any loss of energy. If the numerical scheme is too diffusive, a discontinuity will appear along the segment  $[OA]$ . It is not supposed to happen with streamline diffusion methods which only add diffusivity along the streamline. However, for shock capturing method, we do see a discontinuity, see figure 3.21. As expected, the SUPG used solo does not cause any loss of energy. For the shock capturing methods, the results are very different. OD-K was one of the most effective methods when considering results on the previous test case, but it produces here a massive loss of energy. The discontinuity is less pronounced for ID-A and ID-GC and even less for OD-CK, but this latter method was not very convincing on the previous test case. Finally, our HOME method has the best results in term of energy loss, with a very small discontinuity. Those observations are very encouraging, and we will use this method in further numerical experiences.

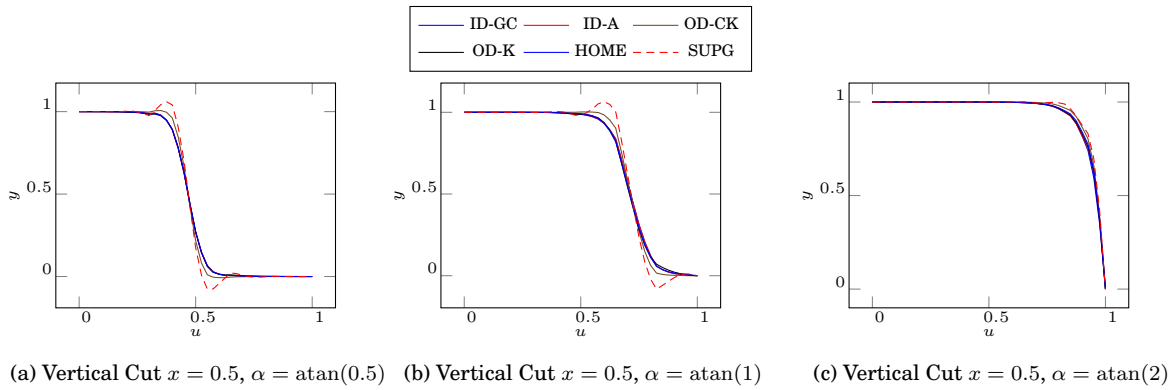


Figure 3.19 – Advection Skew the Mesh profiles for different values of  $\alpha$ . Comparison of the different shock capturing methods, using SUPG $_{1_d}$  as stabilization method.

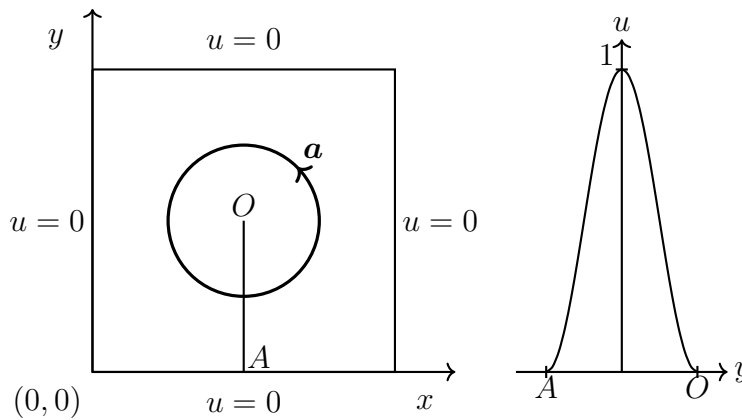


Figure 3.20 – Advection in a rotating flow field setup

### 4.3 Stabilized Navier-Stokes Equations

The previous preliminary results, on advection diffusion reaction problems, allows to efficiently chose the stabilization method, depending on the nature of the problem. In this section, we want to validate those results for the Navier-Stokes system.

#### The Driven Cavity

The steady NS equations are solved on a unit square  $\Omega = [0, 1] \times [0, 1]$ . A constant velocity field  $u = (1, 0)$  is imposed on the top boundary  $y = 1$ ,  $0 \leq x \leq 1$ . Homogeneous Dirichlet condition  $u = (0, 0)$  is imposed on other boundaries, see figure 3.22. We run simulation with four different configurations:  $\text{Re} = 400/h = 0.125$ ,  $\text{Re} = 1000/h = 0.125$ ,  $\text{Re} = 5000/h = 0.0625$  and  $\text{Re} = 10000/h = 0.03$ . All those simulations are performed on a non-structured mesh. The discretization is a standard Taylor-Hood FE space  $\mathbb{T}\mathbb{H}_h^1$ . We performed all the stabilized simulation with the stabilization term on divergence: without this term, the results are not

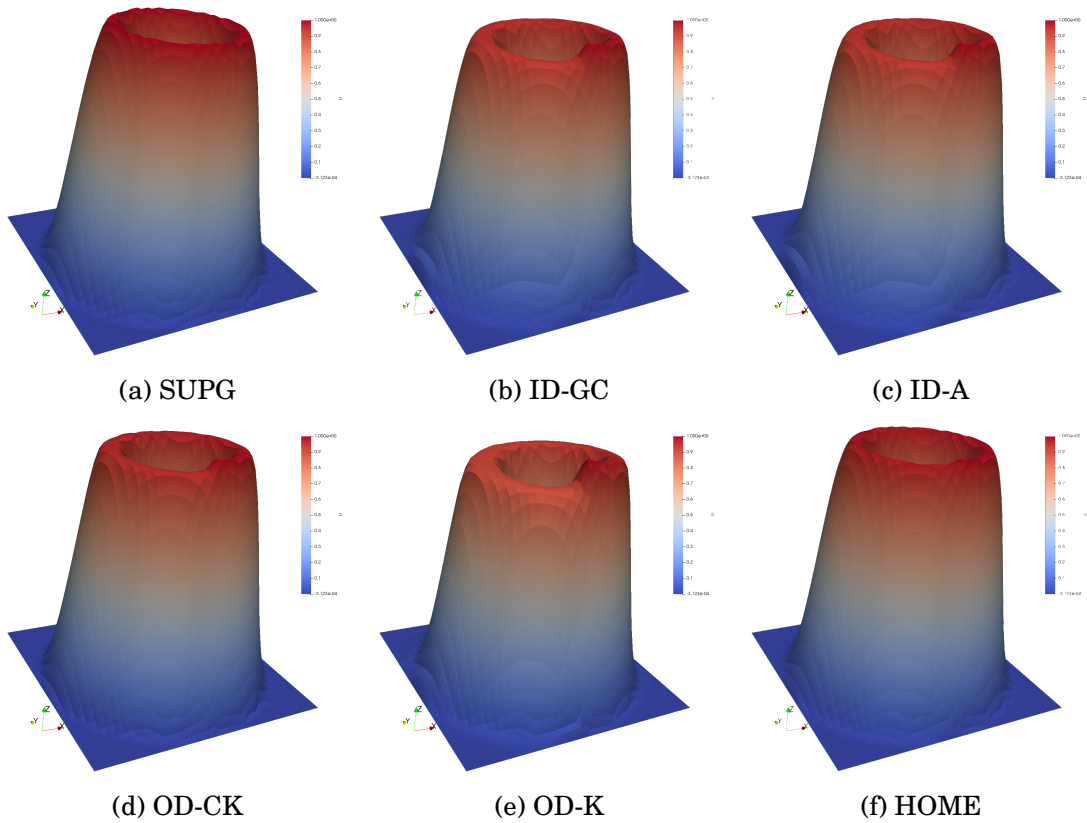


Figure 3.21 – Advection in a Rotating Flow Field.  
Comparison of the different shock capturing methods, using SUPG<sub>1d</sub> as stabilization method.

conclusive for high Reynolds numbers. Our data are compared to reference values from [77]. In this unpublished study, the author made an accurate compilation of different results from many references, [108, 99, 20, 40, 14, 92, 93, 31] and also proposed his own results.

As for the advection-diffusion problem, we will compare the two methods SUPG and GLS, with parameters  $\tau_h^1$  and  $\tau_h^2$ . As we figured out in the previous study, the influence of the characteristic length definition is negligible on homogeneous meshes, thus we only present results with definition  $h_m$  for this test case. For each configuration, we present: an horizontal cut along the axis  $y = 0.5, 0 \leq x \leq 1$  on figure 3.24, a vertical cut along the axis  $x = 0.5, 0 \leq y \leq 1$  on figure 3.23, a pressure profile with contour on figure 3.26 and a velocity profile with streamlines on figure 3.25.

First, regarding the cuts 3.24 and 3.23, we observe some differences between the stabilized simulation and the reference for  $Re = 5000$  and  $Re = 10000$ . The localization of the extremum is correct but not their values. We hypothesized that the mesh was too coarse to model the diffusion in the smaller scales accurately. We performed a new simulation on a thin mesh ( $h = 0.01$ ) and obtained better results. This example perfectly illustrates the issue of using too coarse meshes with stabilization methods: those techniques may give a stable numerical solution,

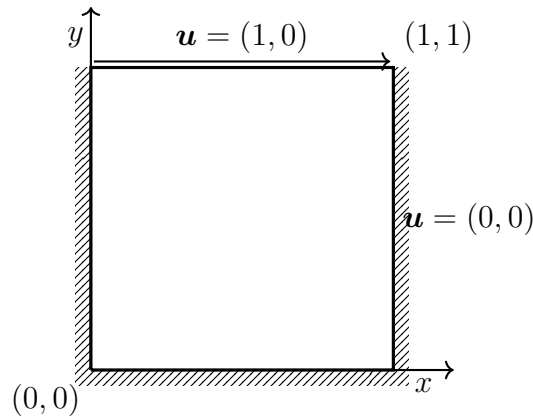
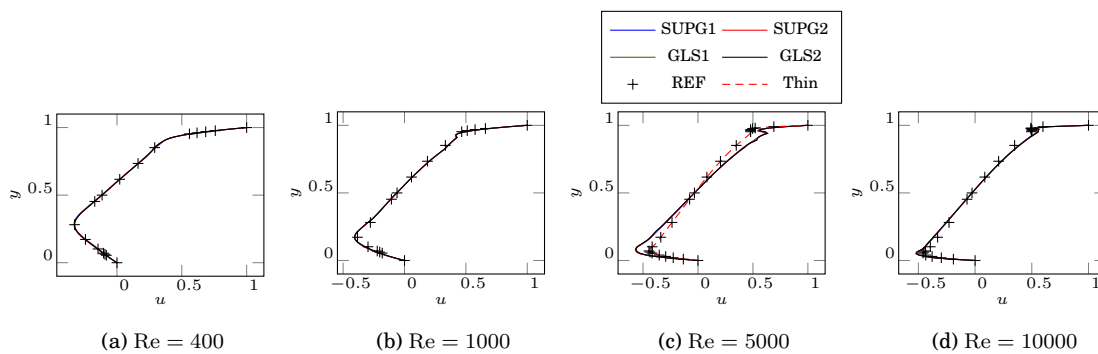


Figure 3.22 – Driven Cavity. Problem Setup

Figure 3.23 – Driven Cavity on Unstructured Homogeneous Mesh.  $u_x$  along a Vertical cut  $x = 0.5$ . Comparison of the two methods with different parameters

but will not properly model the diffusion in the small scales. In the next chapter, we present a solution to this issue, using turbulence models. Eventually, we can conclude in the efficiency of the two considered stabilization methods, as long as the mesh is not too coarse. The methods give comparable results for both parameters,  $\tau_h^1$  and  $\tau_h^2$ , as illustrated on figures 3.25 and 3.26.

## Fluid Layer

This second test case for stabilized Navier-Stokes equation was run on a flat plate geometry, proposed on NASA turbulence resources website<sup>1</sup>, see figure 3.27. This test case was initially designed to benchmark turbulence codes, but for now, we only want to test the robustness of our Navier-Stokes solver with high Reynolds numbers and anisotropic meshes. We only consider the fluid resolution, without turbulence model. The test case will be studied again, with turbulence, in a further section. The geometry of the model is a rectangle domain  $[-0.5, 2] \times [0, 1]$  with the following boundary conditions, see figure 3.28:

<sup>1</sup><https://turbmodels.larc.nasa.gov/>

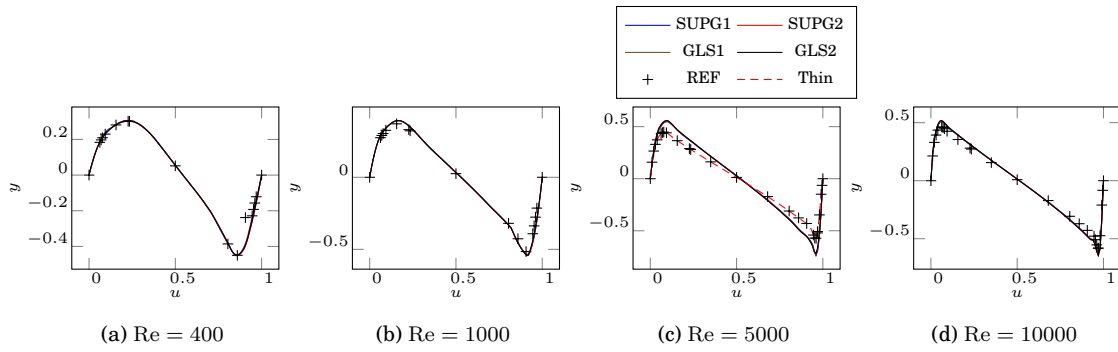


Figure 3.24 – Driven Cavity on Unstructured Homogeneous Mesh.  $u_y$  along a Horizontal cut  $y = 0.5$ . Comparison of the two methods with different parameters

- Inlet ( $x = 0.5, 0 \leq y \leq 1$ ):  $\mathbf{u} = (1, 0)$ ,
- Solid Wall ( $0 \leq x \leq 2, y = 0$ ): no-slip condition  $\mathbf{u} = (0, 0)$ ,
- Bottom ( $-0.5 \leq x \leq 0, y = 0$ ): slip condition  $\mathbf{u} \cdot \mathbf{n} = 0$ ,  $(\mathbf{T}(\mathbf{u}, p)\mathbf{n})_t = 0$ ,
- Top ( $-0.5 \leq x \leq 2, y = 1$ ): slip condition  $\mathbf{u} \cdot \mathbf{n} = 0$ ,  $(\mathbf{T}(\mathbf{u}, p)\mathbf{n})_t = 0$ ,
- Outlet ( $x = 1, 0 \leq y \leq 1$ ): free stream condition  $\mathbf{T}(\mathbf{u}, p)\mathbf{n} = (0, 0)$ .

We faced different issues with this simulation. The geometry presents a harsh singularity at the point  $(0, 0)$  as the boundary condition turns from slip to a no-slip condition. This discontinuity became annoying for high Reynolds numbers ( $\text{Re} \geq 10^4$ ). Another difficulty stands in the highly anisotropic mesh. We generated a simplex-mesh from the grid proposed by the NASA. All the rectangular cells were split into two triangles. In the boundary layer on the bottom, the cells have huge aspect ratios, and we had to adapt our stabilization parameters to fit with it. Initially, we performed simulations on the  $35 \times 35$  grid but this mesh was too coarse, and the solution was not acceptable. We finally present the more satisfying results obtained on the  $69 \times 49$  grid. The most challenging point of these simulations was the stability of the solution around the singularity at  $(0, 0)$ . To give a better view of the problem we present the velocity profile as a warp of the velocity field magnitude. This representation highlights the non-physical oscillations in the boundary layers, see figure 3.29. Note that those oscillations are not visible if you only display the streamlines of the velocity field. The pressure may also present significant fluctuations just after the discontinuity. We never succeeded in entirely removing those oscillations, but we aimed to keep them local near the singularity, see figure 3.30. We ran the simulations with high Reynolds number,  $\text{Re} = 10^5$ . From a physical point of view, the results have no interest since the model lacks diffusivity, but for now, we are only interested in the numerical stability of the different configurations: SUPG/GLS,  $\tau_1/\tau_2$ ,  $h_m/h_h/h_d$ .

From the results of figures 3.29 and 3.30, we clearly see the influence of the characteristic length. This test case confirms the practicality of the  $h_d$  version. This version will be now systematically used on anisotropic meshes. Again, we



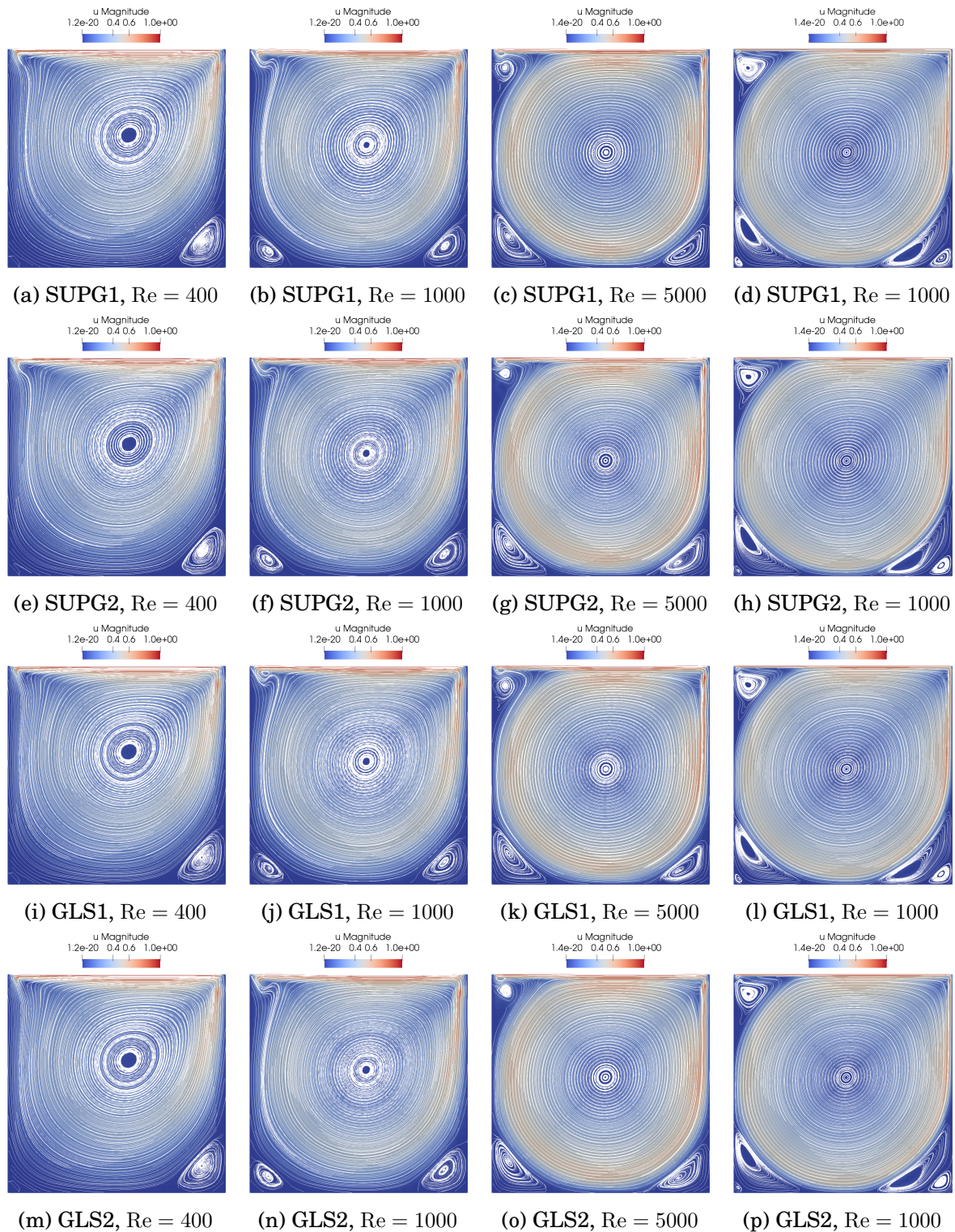


Figure 3.25 – Driven Cavity on Unstructured Homogeneous Mesh. Velocity Streamlines. Comparison of the two methods with different parameters

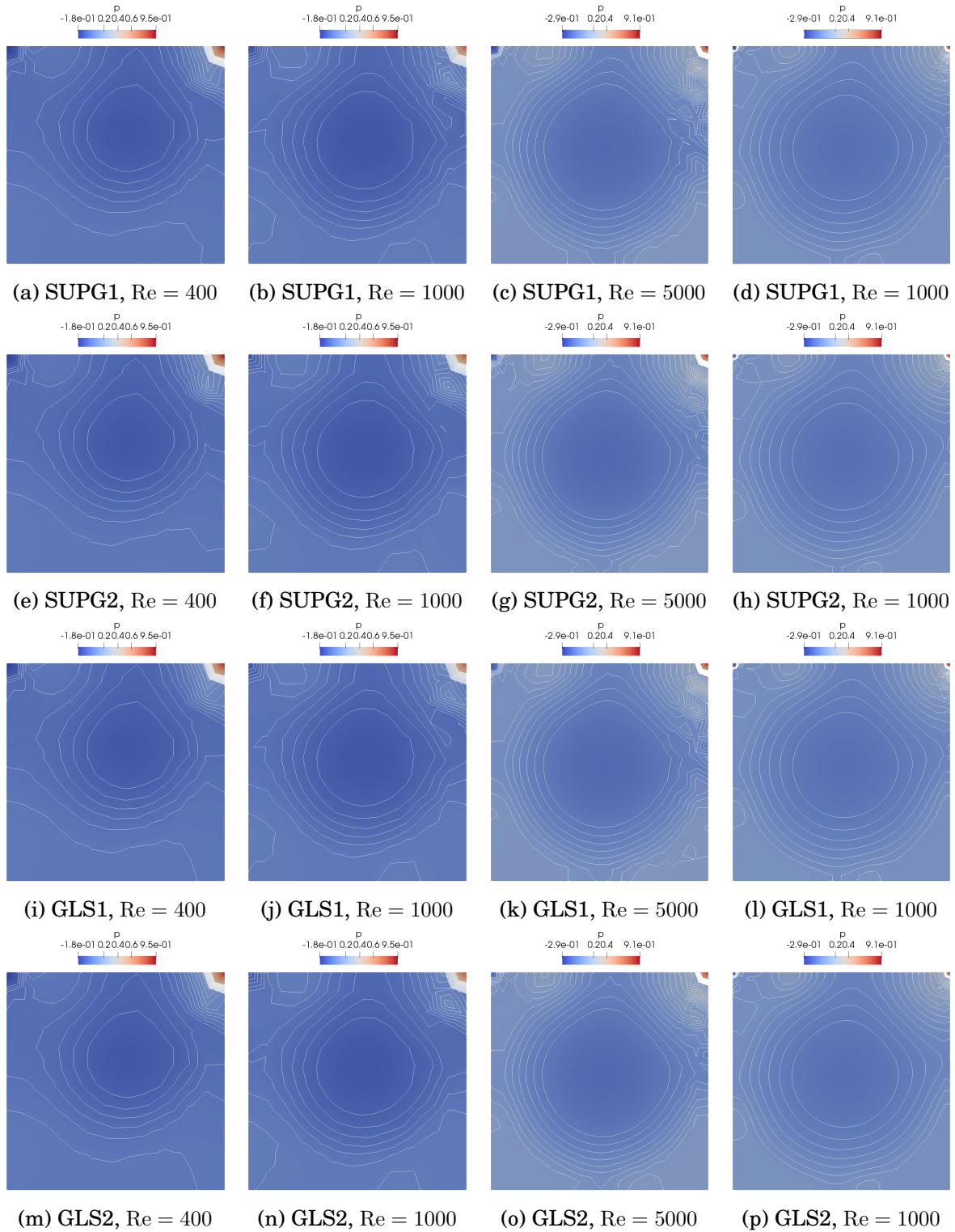


Figure 3.26 – Driven Cavity on Unstructured Homogeneous Mesh. Pressure Contours. Comparison of the two methods with different parameters

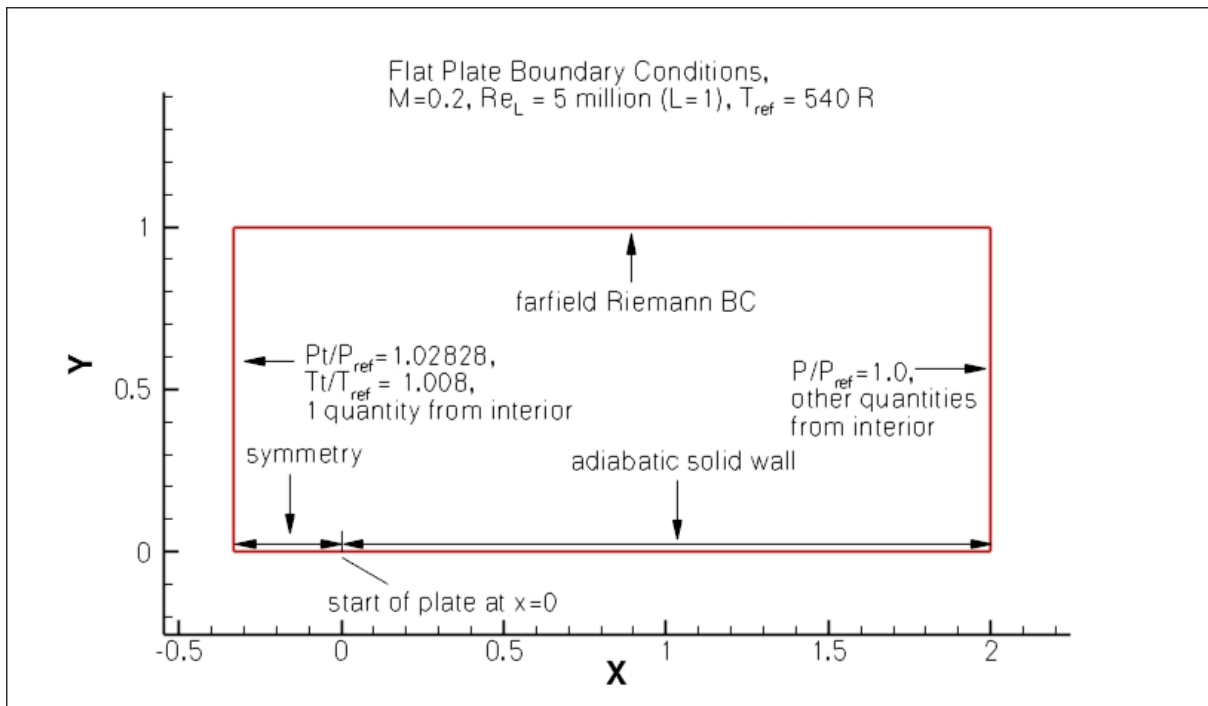


Figure 3.27 – Flat Plate Boundary Conditions, as presented on Nasa resource page.

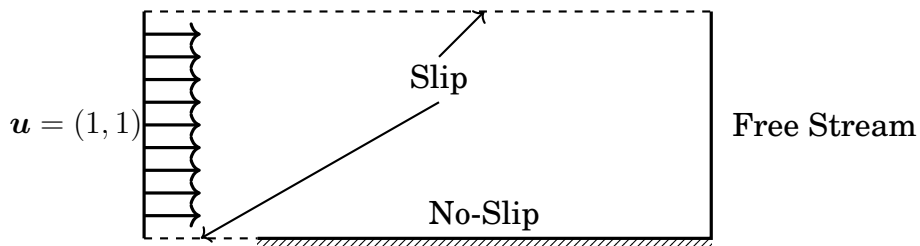


Figure 3.28 – Fluid Layer Setup

denote only small differences between the SUPG and GLS methods. However, the convergence of our non-linear solver is improved when we use the GLS, and this method will be preferred for high Reynolds simulations.

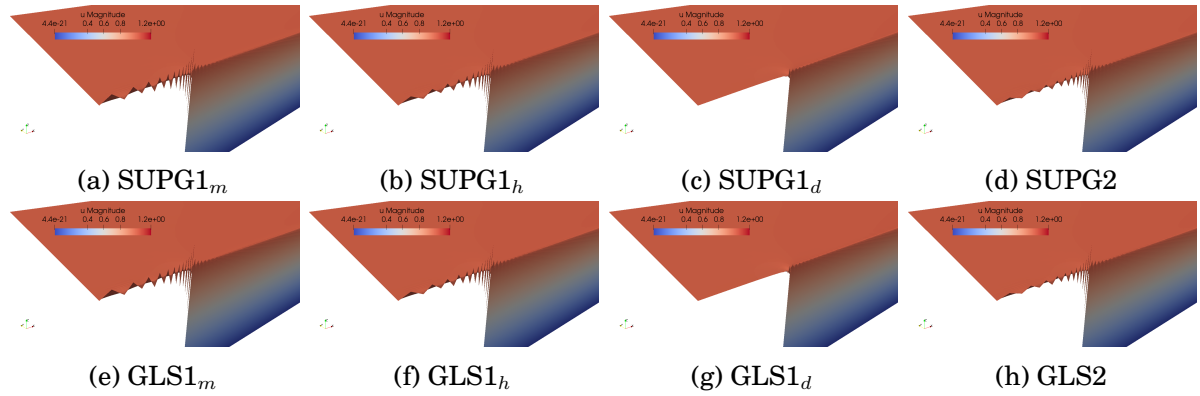


Figure 3.29 – Fluid Layer on structured anisotropic mesh with boundary structure on the bottom. Focus on the region around the singularity. Warp of the magnitude of the velocity field. Comparison of the two methods with different parameters

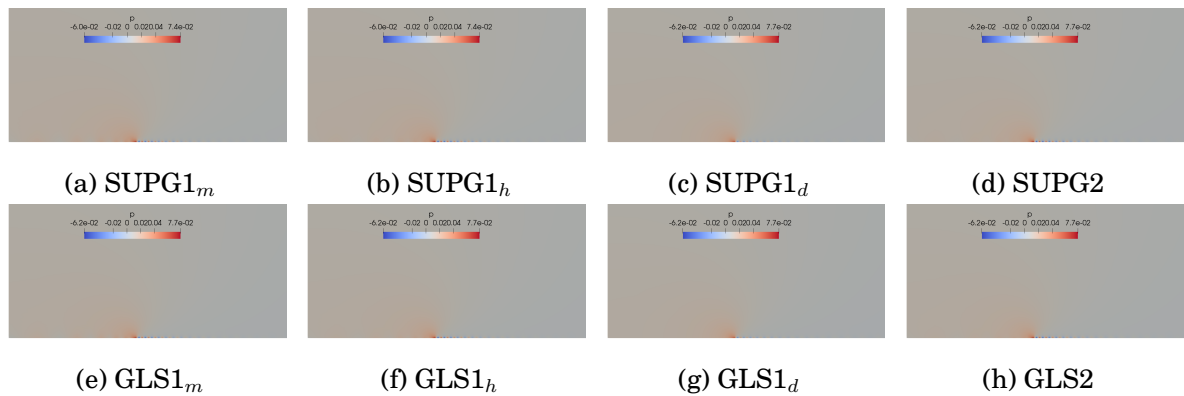


Figure 3.30 – Fluid Layer on structured anisotropic mesh with boundary structure on the bottom. Focus on the region around the singularity. Pressure profile. Comparison of the two methods with different parameters

## Conclusion

In this chapter, we presented the stabilization strategy we developed in our aerothermal framework. The SDM is very efficient and offers an appropriate response to the numerical instability of the Galerkin method. The numerical examples of the last section illustrate the importance of the stabilization parameter. The method is usually very robust but encounters some difficulties on anisotropic meshes or near geometric singularities. In that cases, the formulation  $h_d$  proposed in this chapter appears to be particularly adapted. It also provides accurate results on homogeneous isotropic meshes. However, the computation of this characteristic length is a bit more expensive than the  $h_m$  or  $h_h$  versions which do not depend on the convection field.

For advection-diffusion problems, SUPG and GLS usually give similar solutions and the SUPG will be used by default since it is a bit cheaper to compute. However, for advection-diffusion-reaction or the Navier-Stokes equations, GLS produces significantly smoother results. In the table 3.1, we propose a recap of the recommended configurations for each kind of problems.

In all this chapter we only presented and discussed the formulations for triangular meshes. This constraint is proper to the Feel++ library which does not provide the required tools for hypercubes. We expected to recover a more stable system, especially in the boundary layers, when working with hypercubes. Unfortunately, the needed implementation in Feel++ was not fulfilled, and this conjecture remains an open question.

Those stabilization issues were more challenging than initially expected. We spend a considerable amount of time on this topic, but once those stabilization methods were fully mastered, it finally opened the way to the turbulence modeling.

	Isotropic Mesh			Anisotropic Mesh		
	SDM	$\tau$	$h$	SDM	$\tau$	$h$
Advection Diffusion	SUPG	$\tau_1$ or $\tau_2$	$h_m$ or $h_h$	SUPG	$\tau_1$	$h_h$
Advection Diffusion Reaction	GLS	$\tau_3$	$h_m$ or $h_h$	GLS	$\tau_3$	$h_d$
Navier-Stokes	SUPG	$\tau_1$ or $\tau_2$	$h_m$ or $h_d$	GLS	$\tau_1$	$h_d$

Table 3.1 – Recap of recommended configurations for different kinds of problems.

## Turbulence Modeling

**Contents**

<b>1</b>	<b>Reynolds Averaged Navier-Stokes Method . . . . .</b>	<b>76</b>
1.1	Time-Averaged Equations . . . . .	76
1.2	The Boussinesq Approximation and Derived Models . . . . .	78
1.3	Resolution Strategy . . . . .	79
<b>2</b>	<b>Spalart Allmaras Model . . . . .</b>	<b>79</b>
2.1	Equations . . . . .	80
<b>3</b>	<b><math>k - \omega</math> SST model . . . . .</b>	<b>80</b>
3.1	The Equations . . . . .	81
3.2	Boundary and Initial Conditions . . . . .	82
<b>4</b>	<b>Validation: 2D Naca0012 Airfoil . . . . .</b>	<b>84</b>
4.1	Problem Description . . . . .	84

The behavior of a fluid is defined by the Reynolds number  $Re$ ,

$$Re = \frac{U_0 L}{\nu} \quad (4.1)$$

where  $U_0$  is the far-field velocity magnitude,  $L$  is the characteristic length of the problem (usually the smallest significant length) and  $\nu$  is the kinematic viscosity of the fluid. For low Reynolds number, the flow is laminar: the viscosity forces dominate, the model usually admits a steady state and is totally stable. When the

Reynolds number grows, the physic becomes less compliant: the stationary state is not guaranteed anymore, and we usually need stabilization methods if we do not want to pay the price of a very thin mesh. This kind of flow is not turbulent yet but can be defined as a transition state. If the flow is not stationary, he may recover a periodic behavior and is still deterministic. For higher Reynolds, the flow becomes turbulent: the behavior is chaotic with many scales of eddies.

The methodologies presented in the previous chapter provide an accurate framework for modelization of laminar and transitional flows but are inefficient in the simulation of turbulent problems. Even if the stabilization method allows to perform simulations with high Reynolds numbers, they do not take care of the energy which is supposed to be diffused in the smallest scales. We could theoretically solve all the different scales and perform what is know as a Direct Numerical Simulation (DNS) but for very high Reynolds number, the required time and space discretization would be ridiculously not practicable [57].

Here comes the necessity to introduce a turbulence model in order to modelize all the energy exchanges, at all scales, but with no proper resolution and a reasonable discretization. We mainly distinguish two families of turbulence model:

- The Large Eddy Simulation (LES) in which only the small structures are modeled. The larger structures are solved using Navier-Stokes equations. This method is particularly adapted to transient problems and was not studied in detail during this thesis.
- The Reynolds Average Navier-Stokes (RANS) method which consists in the averaging of the chaotic component of the flow. The next section is dedicated to this method and will present the main ingredients.

The field of turbulence modelization is very dense and complicated (especially when you want to tackle it with a mathematical background ). In this chapter, we only briefly remind the equations and the main ingredient of the presented models. For more details, we refer the reader to [41], and the further cited references.

# 1 Reynolds Averaged Navier-Stokes Method

## 1.1 Time-Averaged Equations

We detail here the averaging of the Navier-Stokes equations, in order to recover the RANS equations. The treatment of the energy equation will be discussed further.

## Reynolds Decomposition

The Reynolds Decomposition is used in order to separate the value of a quantity into its mean value and its fluctuations,

$$\mathbf{u} = \bar{\mathbf{u}} + \mathbf{u}', \quad (4.2)$$

where  $\bar{\mathbf{u}}$  is the mean value, also called steady or averaged component, and  $\mathbf{u}'$  is the fluctuations.

The Reynolds Averaged Navier-Stokes (RANS) equations are obtained by considering the Reynolds decomposition of the two fluid quantities

$$\begin{aligned} \mathbf{u} &= \bar{\mathbf{u}} + \mathbf{u}' \\ p &= \bar{p} + p' \end{aligned} \quad (4.3)$$

## Averaged Equations

After substitution of (4.3) into the Navier-Stokes equations, we obtain

$$\begin{cases} \rho(\bar{\mathbf{u}} + \mathbf{u}') \cdot \nabla(\bar{\mathbf{u}} + \mathbf{u}') + \nabla(\bar{p} + p') - \mu\Delta(\bar{\mathbf{u}} + \mathbf{u}') = \mathbf{f} \\ \nabla \cdot (\bar{\mathbf{u}} + \mathbf{u}') = 0. \end{cases} \quad (4.4)$$

This system of equation is then averaged. In the couple further formulas, we chose to use the tensor notation which was a bit more convenient.

$$\begin{cases} \overline{\rho(\bar{\mathbf{u}}_j + \mathbf{u}'_j) \cdot \frac{\partial(\bar{\mathbf{u}}_i + \mathbf{u}'_i)}{\partial x_j} + \frac{\partial(\bar{p} + p')}{\partial x_i} - \mu \frac{\partial^2(\bar{\mathbf{u}}_i + \mathbf{u}'_i)}{\partial x_j \partial x_j}} = \bar{\mathbf{f}}_i \\ \overline{\frac{\partial(\bar{\mathbf{u}}_i + \mathbf{u}'_i)}{\partial x_i}} = 0. \end{cases} \quad (4.5)$$

Given that  $\bar{\cdot}$  is a Reynolds operator, we may use the following properties,

$$\overline{a'} = 0, \quad \overline{\bar{a}} = \bar{a}, \quad \overline{a + b} = \bar{a} + \bar{b}, \quad \overline{ab} = \bar{a}\bar{b} + \overline{a'b'}, \quad (4.6)$$

to rewrite the previous equation (4.5),

$$\begin{cases} \rho \bar{\mathbf{u}}_j \frac{\partial \bar{\mathbf{u}}_i}{\partial x_j} + \overline{\rho \mathbf{u}'_j \frac{\partial \mathbf{u}'_i}{\partial x_j}} + \frac{\partial \bar{p}}{\partial x_i} - \mu \frac{\partial^2 \bar{\mathbf{u}}_i}{\partial x_j \partial x_j} = \bar{\mathbf{f}}_i \\ \frac{\partial \bar{\mathbf{u}}_i}{\partial x_i} = 0. \end{cases} \quad (4.7)$$

This formulation is very close to the classic incompressible Navier-Stokes equations, except for the term  $\overline{\mathbf{u}'_j \frac{\partial \mathbf{u}'_i}{\partial x_j}}$  which can be rewritten (using the continuity equation of the fluctuation) as  $\frac{\partial \overline{\mathbf{u}'_i \mathbf{u}'_j}}{\partial x_i}$ . The quantity  $\mathbf{R}_{ij} = \overline{\rho \mathbf{u}'_i \mathbf{u}'_j}$  is known as the **Reynolds**



**stress tensor** and represents the effects of velocity fluctuation on the averaged system. With this notation, the system (4.7) reads,

$$\begin{cases} \bar{u}_j \frac{\partial \bar{u}_i}{\partial x_j} + \frac{\partial \bar{p}}{\partial x_i} - \mu \frac{\partial^2 \bar{u}_i}{\partial x_j \partial x_j} + \rho \frac{\partial \mathbf{R}_{ij}}{\partial x_j} = \bar{f}_i \\ \frac{\partial \bar{u}_i}{\partial x_i} = 0. \end{cases} \quad (4.8)$$

Using the Reynolds decomposition in this averaged equations, we introduced 3 unknown quantities:  $\bar{u}$ ,  $\bar{p}$  and  $\mathbf{R}$ . Unfortunately, we have gained no additional equations. We now need to close this system with additional equations. That closure is provided by the different turbulence models. In what follow we focus on models derived from the Boussinesq approximation.

## 1.2 The Boussinesq Approximation and Derived Models

The key of the RANS approach is to find an appropriate model for the Reynolds stress tensor. A common method is to relate the Reynolds stresses to the mean velocity gradients

$$-\mathbf{R}_{ij} = \mu_t \frac{\partial u_i}{\partial x_j} - \frac{2}{3} \rho k \delta_{ij} \quad (4.9)$$

where the turbulence kinetic energy is defined as  $k = \frac{1}{2} \overline{u'_i u'_j}$  and  $\mu_t$  is the eddy viscosity.

The RANS models derived from the Boussinesq approximation provide closure equations in order to compute  $\mu_t$  and  $k$ . Many different models can be found in the literature. In the simplest models, the kinetic energy is neglected, and the turbulent viscosity is directly express from the velocity fields. Those algebraic models are usually not accurate enough. The most common models express the turbulent viscosity  $\mu_t$  and the kinetic energy  $k$  as solutions of one (Spalart-Allmaras, Baldwin-Barth, ...) or two equations ( $k - \epsilon$ ,  $k - \omega$ , SST, ...). Those models are reasonably complex to be solved efficiently and are sufficiently accurate in their dedicated fields of application. More complex models exist but were not studied.

In the context of this thesis, we were first interested in the Spalart-Allmaras model which is known as efficient and relatively simple. In a second time, we chose to implement a solver for the  $k - \omega$ -SST model which appears to be more adapted to the kind of problems we were studying. Those two models are detailed in the following sections.

### 1.3 Resolution Strategy

The resolution of the coupled RANS problems is integrated into the iterations of the non-linear solver for the aerothermal problem described in chapter 2. Turbulent problems requires both physical and  $\Psi$ tc continuations, the turbulence model is then activated when the molecular viscosity reaches a limit value  $\mu^*$  in the physical continuation process. Then, the turbulent viscosity  $\mu_t$  is updated at each iteration in the non-linear solver (usually during the update of the Jacobian in Newton algorithm). The actual computation of  $\mu_t$  depends on the choice of the model and will be detailed for each of them in the following sections.

We summarized the resolution procedure in the algorithm

---

**Algorithm 1:** Resolution Procedure for RANS Model using Newton Algorithm

---

```

/* We denote by  $\mathbf{W}$  the composite vector  $\mathbf{W} = (\mathbf{u}_h, p_h, T_h)$  */
1 Choose  $\mu^*$ ; // Limit viscosity
2 Choose  $N$ ; // Number of steps in the physical continuation
3 Choose  $\delta_0$ ; // Initial  $\Psi$ tc pseudo time-step
4 Choose  $\varepsilon$ ; // Tolerance for Newton algorithm
5 for  $0 \leq n \leq N$  do // Physical continuation loop
6    $\mu_n \leftarrow \mu^{\frac{n}{N}}$ ;
7    $\kappa_n \leftarrow \kappa^{\frac{n}{N}}$ ;
8    $\delta \leftarrow \delta_0$ ;
9   repeat // Newton iterations
10    if  $\mu_n \leq \mu^*$  then
11     | Update  $\mu_t^k$  with  $\mathbf{W}^k$ ; // Solve turbulence model PDE
12    else
13     |  $\mu_t \leftarrow 0$ ;
14    end
15    Solve  $(\frac{1}{\delta} \mathbf{D} + \mathcal{F}'(\mathbf{W}^k; \mu^n, \kappa_n, \mu_t^k)) \mathbf{S} = -\mathcal{F}(\mathbf{W}^k; \mu^n, \kappa_n, \mu_t^k)$ ;
16     $\mathbf{W}^k \leftarrow \mathbf{W}^k + \mathbf{S}$ ;
17    Update  $\delta$ ; // See section 3
18  until  $\|\mathcal{F}(\mathbf{W}^k; \mu_n, \kappa_n, \mu_t)\|_2 \leq \varepsilon$ ;
19 end
20 return  $\mathbf{W}^n$ 

```

---

## 2 Spalart Allmaras Model

The Spalart Allmaras (SA) [96] model is probably one of the most used RANS model. He is very popular thanks to its relative simplicity and accuracy. In this

model, the kinetic energy  $k$  is neglected and the eddy viscosity  $\mu_t$  is computed through the resolution of a non-linear PDE.

## 2.1 Equations

The model is described by one transport equation. The calibration of the many closure constants and more details can be read in [95, 50].

The eddy viscosity is given by  $\mu_t = \rho \bar{\nu} f_{v1}$ , where  $\bar{\nu}$  is the solution of

$$\underbrace{\mathbf{u} \cdot \nabla \bar{\nu}}_{\text{Convection}} - \underbrace{\frac{1}{\sigma} [\nabla \cdot ((\nu + \bar{\nu}) \nabla \bar{\nu}) + c_{b2} (\nabla \bar{\nu} \cdot \nabla \bar{\nu})]}_{\text{Diffusion}} - \underbrace{c_{b1} (1 - f_{t2}) S \bar{\nu}}_{\text{Production}} + \underbrace{[c_{w1} f_w - \frac{c_{b1}}{\kappa^2} f_{t2}] \left(\frac{\bar{\nu}}{d}\right)^2}_{\text{Destruction}} = 0, \quad (4.10)$$

with the following notations,

$$\begin{aligned} \chi &= \frac{\bar{\nu}}{\nu}, & S &= \sqrt{(\nabla \times \mathbf{u}) \cdot (\nabla \times \mathbf{u})} + \frac{\bar{\nu}}{\kappa^2 d^2} f_{v2}, \\ f_{v1} &= \frac{\chi^3}{\chi^3 + c_{v1}^3}, & f_{v2} &= 1 - \frac{\chi}{1 + \chi f_{v1}}, & f_{t2} &= c_{t3} \exp(-c_{t4} \chi^2), \\ f_w &= g \left( \frac{1 + c_{w3}^6}{g^6 + c_{w3}^6} \right)^{1/6}, & g &= r + c_{w2} (r^6 - r), & r &= \min \left( \frac{\bar{\nu}}{S \kappa^2 d^2}, 10 \right), \end{aligned} \quad (4.11)$$

where  $d$  is the distance from the closest surface. We furthermore consider the constant parameters

$$\begin{aligned} c_{b1} &= 0.1355, & c_{b2} &= 0.622, & c_{w1} &= \frac{c_{b1}}{\kappa^2} + \frac{1 + c_{b2}}{\sigma}, & c_{w2} &= 0.2, & c_{w3} &= 2 \\ c_{v1} &= 7.1, & c_{t3} &= 1.2, & c_{t4} &= 0.5, & \kappa &= 0.41, & \sigma &= 2/3. \end{aligned} \quad (4.12)$$

The problem is closed with the following boundary conditions,

$$\left\{ \begin{array}{ll} \bar{\nu} = 0, & \text{on } \Gamma_W \text{ (Walls)} \\ \bar{\nu} = 3\nu, & \text{on } \Gamma_I \text{ (Inlets)} \\ \nabla \bar{\nu} \cdot \mathbf{n} = 0, & \text{on } \Gamma_O \text{ (Outlets)}. \end{array} \right. \quad (4.13)$$

## 3 $k - \omega$ SST model

The  $k - \omega$  SST model is known to be very robust and adaptive. The model was first proposed in [70] and [71]. The model was then improved in [72]. We chose to work on this formulation, known as the SST-2003 model. Note that many other versions of the model were proposed, for instance in [94].

### 3.1 The Equations

The turbulent eddy viscosity is computed by

$$\mu_t = \frac{\rho a_1 k}{\max(a_1 \omega, \mathcal{S} F_2)} \quad (4.14)$$

where the turbulent kinetic energy  $k$  and the specific dissipation  $\omega$  are solutions of the coupled transport equations

$$\begin{cases} \frac{\partial k}{\partial t} + \mathbf{u} \nabla k + \beta^* k \omega - \nabla \cdot \left[ \frac{\mu + \sigma_k \mu_t}{\rho} \nabla k \right] = P_k, \\ \frac{\partial \omega}{\partial t} + \mathbf{u} \nabla \omega + \beta \omega^2 - \nabla \cdot \left[ \frac{\mu + \sigma_w \mu_t}{\rho} \nabla \omega \right] = \frac{\gamma}{\mu_t} P_k + 2(1 - F_1) \frac{\sigma_{w2}}{\omega} \nabla k \cdot \nabla \omega, \end{cases} \quad (4.15)$$

with the following closure relations

$$F_1 = \tanh(\arg_1^4), \quad \arg_1 = \min \left[ \max \left( \frac{\sqrt{k}}{\beta^* \omega y}, \frac{500 \mu}{\rho y^2 \omega} \right), \frac{4 \sigma_{w2} k}{CD_{k\omega} y^2} \right] \quad (4.16)$$

$$F_2 = \tanh(\arg_2^2), \quad \arg_2 = \max \left( \frac{2\sqrt{k}}{\beta^* \omega y}, \frac{500 \mu}{\rho y^2 \omega} \right) \quad (4.17)$$

$$P_k = \min(\bar{\tau} : \nabla \mathbf{u}, 10 \beta^* k \omega), \quad CD_{k\omega} = \max \left( 2 \rho \frac{\sigma_{w2}}{\omega} \nabla k \nabla \omega, 10^{-10} \right), \quad (4.18)$$

$$\mathcal{S} = \sqrt{2(\bar{S} : \bar{S})}. \quad (4.19)$$

Each closure constant is a blend of an inner and an outer constant, via

$$\phi = F_1 \phi_1 + (1 - F_1) \phi_2 \quad (4.20)$$

where  $\phi_1$  represents constant 1 and  $\phi_2$  represents constant 2. The closure constants are

$$\begin{aligned} \sigma_{k1} &= 0.85, & \sigma_{\omega1} &= 0.5, & \beta_1 &= 0.075, & \gamma_1 &= 5/9 \\ \sigma_{k2} &= 1.0, & \sigma_{\omega2} &= 0.856, & \beta_2 &= 0.0828, & \gamma_2 &= 0.44 \\ \beta^* &= 0.09, & \kappa &= 0.41, & a_1 &= 0.31. \end{aligned} \quad (4.21)$$

We also remind some usual notations.  $\bar{\tau}$  is the Reynolds stress tensor, defined, for incompressible flow, by

$$\bar{\tau}_{ij} = 2\mu_t \bar{S}_{ij} - \frac{2}{3} \rho k \delta_{ij} \quad (4.22)$$

where  $\delta_{ij}$  is the Kronecker delta and  $\bar{S}$  is given by stress tensor,

$$\bar{S}_{ij} = \frac{1}{2} \left( \frac{\partial u_i}{\partial x_j} + \frac{\partial u_j}{\partial x_i} \right). \quad (4.23)$$

We can also rewrite this tensors as

$$\bar{\tau} = 2\mu_t \bar{S} - \frac{2}{3} \rho k \mathbb{I}_d, \quad \bar{S} = \frac{1}{2} (\nabla \mathbf{u} + (\nabla \mathbf{u})^T) \quad (4.24)$$

### 3.2 Boundary and Initial Conditions

As initial values for  $k$  and  $\omega$ , we choose

$$k_0 = \left( \frac{\mu_0}{\rho l_0} \right)^2, \quad \omega_0 = \frac{k^{1/2}}{l_0} \quad (4.25)$$

where  $\mu_0 = \mathcal{O}(\mu)$  is a constant initial viscosity and  $l_0$  is a default length such that  $l_0 \in [l_{\min}, l_{\max}]$ ,  $l_{\min}$  corresponding to the size of the smallest admissible eddies.

At the inflow boundary the prescribed value are

$$k = c_{bc}|u|^2, \quad \omega = \frac{k^{1/2}}{l_0}, \quad \text{on } \Gamma_{\text{in}}, \quad (4.26)$$

where  $c_{bc} \in [0.003, 0.01]$  is an empirical constant.

At the walls, we impose slip condition on the velocity

$$\mathbf{u} \cdot \mathbf{n} = 0, \quad \text{on } \Gamma_w. \quad (4.27)$$

#### Wall functions

We denote by  $t_w$  the tangential component of the force exerted by the viscous stress tensor  $\bar{\varepsilon}(\mathbf{u}) = \frac{\mu + \mu_t}{\rho} (\nabla \mathbf{u} + \nabla \mathbf{u}^T)$  on  $\Gamma_w$

$$\mathbf{t}_w = \mathbf{n} \cdot \bar{\varepsilon}(\mathbf{u}) - ((\mathbf{n} \cdot \varepsilon(\mathbf{u})) \cdot \mathbf{n})\mathbf{n}. \quad (4.28)$$

We consider an internal boundary  $\Gamma_y$  located at a distance  $y$  from the physical wall  $\Gamma_w$ . Assuming that we are in the log-layer region, the logarithm wall laws leads to the boundary conditions

$$\mathbf{t}_w = -u_\tau \frac{\mathbf{u}}{|\mathbf{u}|}, \quad k = \frac{u_\tau^2}{\sqrt{\beta^*}}, \quad \omega = \frac{u_\tau}{\sqrt{\beta^*} \kappa y} \quad (4.29)$$

where  $\beta^* = 0.09$ ,  $\kappa = 0.41$  is the von Karman constant. The friction velocity satisfy the equation

$$\frac{|\mathbf{u}|}{u_\tau} = \frac{1}{\kappa} \log y^+ + \beta \quad (4.30)$$

assuming that the local Reynolds number  $y^+ = \frac{\rho u_\tau y}{\mu}$  is in the range  $11.06 \leq y^+ \leq 300$ . In theory we have to impose the wall function on the inner boundary  $\Gamma_y$  and the boundary layer of width  $y$  should be removed from the domain  $\Omega$ . In practice, the wall function is applied on the first degree of freedom at the interior of the domain but we can not guarantee that this first dof will not fall into the viscous layer where the current model would not be valid anymore.

To avoid this issue, we chose a solution proposed in [39]. We assume that the first node always corresponds to a fixed value of  $y^+$ . Since the deviation to the wall

is assumed to remain very small, we can neglect the width of this layer and then we avoid a costly mesh adaptation for each iteration. The wall function are then directly applied on  $\Gamma_w$ , just as if the nodes of this boundary were shifted in the normal direction by the distance  $y$ .

The accuracy of this model depends of the choice of  $y$  which has to be chosen as small as possible. [39] explains that the smallest acceptable value for  $y$  corresponds to meeting point between the log-layer and the viscous sub-layer. At this point  $y^+$  verifies the both equations

$$\frac{|\mathbf{u}|}{u_\tau} = \frac{1}{\kappa} \log y^+ + \beta, \quad (4.31)$$

$$\frac{|\mathbf{u}|}{u_\tau} = y^+. \quad (4.32)$$

We can then identify the corresponding value,  $\tilde{y}^+$ , which verifies

$$\tilde{y}^+ = \frac{1}{\kappa} \log \tilde{y}^+ + \beta, \quad (4.33)$$

for the default value of  $\kappa = 0.41$  and  $\beta = 5.2$ , the non linear equation(4.33) can be solved iteratively and we find  $\tilde{y}^+ = 11,06$ . All the nodes of  $\Gamma_w$  are now supposed to be at a distance  $y$  of the physical wall such that  $y^+ = \tilde{y}^+$ .

On this boundary  $\Gamma_w$  we can now evaluate the friction velocity  $u_\tau$

$$u_\tau = \frac{|\mathbf{u}|}{\frac{1}{\kappa} \log \tilde{y}^+ + \beta} = \frac{|\mathbf{u}|}{\tilde{y}^+} \quad (4.34)$$

on the other hand, by considering (4.29) we have the relation  $u_\tau = (\beta^*)^{\frac{1}{4}} \sqrt{k}$ , [39] proposes to set

$$u_\tau = \max \left( (\beta^*)^{\frac{1}{4}} \sqrt{k}, \frac{|\mathbf{u}|}{\tilde{y}^+} \right), \quad \mathbf{t}_w = -\frac{u_\tau}{\tilde{y}^+} \mathbf{u}. \quad (4.35)$$

Resulting of (4.35) we have natural boundary condition to plug in the variational formulation for the RANS equation

$$\int_{\Gamma_w} \mathbf{t}_w \cdot \mathbf{v} \, d\Gamma = - \int_{\Gamma_w} \frac{u_\tau}{\tilde{y}^+} \mathbf{u} \cdot \mathbf{v} \, d\Gamma \quad (4.36)$$

when the condition (4.27) will be imposed weakly with term

$$\int_{\Gamma_w} \gamma_h (\mathbf{u} \cdot \mathbf{n}) \cdot (\mathbf{v} \cdot \mathbf{n}) \, d\Gamma = 0 \quad (4.37)$$

with  $\gamma_h = \gamma \max \left( \frac{\mu + \mu_t}{h_e}, |\mathbf{u}| \right)$ , where  $\gamma$  is a penalty constant and  $h_e$  is the characteristic length of the cell, see 3 for the evaluation of  $h_e$ .

Using (4.29) we now evaluate the eddy viscosity on  $\Gamma_w$ , assuming that the model for this value of  $y^+$  will act as a classic  $k - \omega$  model,

$$\mu_t = \rho \frac{k}{\omega} = \rho u_\tau \kappa y = \kappa \tilde{y}^+ \mu \quad (4.38)$$

This value is automatically satisfied if we strongly impose the value of  $k$  and  $\omega$  on  $\Gamma_w$ , but this would result in a one way coupling for the boundary conditions between the fluid equations and the  $k - \omega$  system. To avoid this problem, we chose to weakly impose the boundary condition for  $k$  and  $\omega$ ,

$$\nabla k \cdot \mathbf{n} = -\frac{\partial k}{\partial y} = 0, \quad (4.39)$$

$$\nabla \omega \cdot \mathbf{n} = -\frac{\partial \omega}{\partial y} = \frac{u_\tau}{\sqrt{\beta^* \kappa y^2}} = \frac{\omega}{y} \quad (4.40)$$

we can then replace the unknown  $y = \frac{\tilde{y}^+ \mu}{\rho u_\tau}$ , and we obtain the natural weak formulation for the boundary condition on  $\Gamma_w$

$$\int_{\Gamma_w} \frac{\mu + \sigma_k \mu_t}{\rho} \nabla k \cdot \mathbf{n} w \, d\Gamma = 0, \quad (4.41)$$

$$\int_{\Gamma_w} \frac{\mu + \sigma_\omega \mu_t}{\rho} \nabla \omega \cdot \mathbf{n} w \, d\Gamma = \int_{\Gamma_w} \frac{\mu + \sigma_\omega \mu_t}{\mu} \frac{u_\tau}{\tilde{y}^+} \omega w \, d\Gamma. \quad (4.42)$$

Since these two quantities are not imposed strongly, we have to impose the exact value for the eddy viscosity  $\mu_t$  using the relation (4.34) and the production term  $P_k$  on  $\Gamma_k$ ,

$$P_k = 10\beta^* \rho \omega k = 10 \frac{u_\tau^3}{\kappa y} = 10 \frac{\rho u_\tau^4}{\mu_t} = 10 \frac{\rho u_\tau}{\kappa \tilde{y}^+ \mu_t} \quad (4.43)$$

## 4 Validation: 2D Naca0012 Airfoil

In this section we propose to validate our implementation of the presented turbulence models on a classic 2D benchmark. This validation case is proposed on the NASA Turbulence Modeling Resource page<sup>1</sup>

### 4.1 Problem Description

The geometry is a 2D cut of a NACA0012 airfoil, in sufficiently large domain to neglect the side effects. The airfoil surface is defined by the curve

$$y = \pm 0.5946891(0.2982227\sqrt{x} - 0.1271252x - 0.3579079x^2 + 0.2919849x^3 - 0.1051746x^4). \quad (4.44)$$

The figure 4.1 is an overview of the domain, discretized with a grid proposed by the NASA. The problem is described for compressible code but with low mach number  $M = 0.15$ . The Reynolds number is  $\text{Re} = 6 \cdot 10^6$ . It corresponds to a viscosity  $\mu = 0.16 \cdot 10^{-6}$  with a density  $\rho = 1$ , a farfield velocity  $V = 1$  and the characteristic length of the airfoil  $L = 1$ .

<sup>1</sup><https://turbmodels.larc.nasa.gov/index.html>

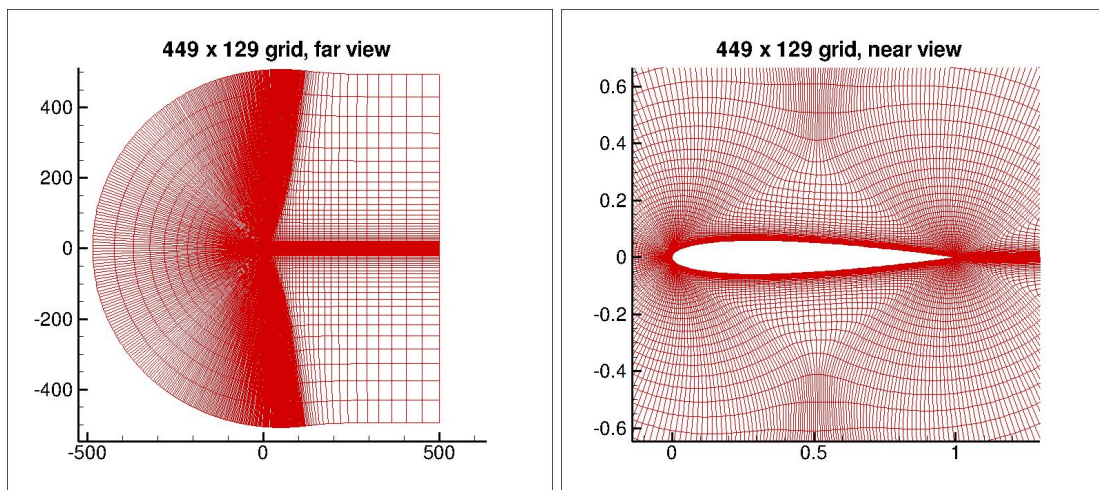


Figure 4.1 – NACA0012 : Grid  $449 \times 129$  from the NASA website





## **Part II**

# **Model Order Reduction**



---

The Reduced Basis Method

---

**Contents**


---

<b>1</b>	<b>Method and Notations : Application to Coercive Problems .</b>	<b>90</b>
1.1	Reduced Basis for Coercive Problems . . . . .	92
1.2	Error Estimation . . . . .	95
<b>2</b>	<b>Reduced Basis for Saddle Point Problems . . . . .</b>	<b>98</b>
2.1	Problem formulation . . . . .	99
2.2	Error Estimation . . . . .	102
<b>3</b>	<b>Numerical Applications : RB for Stokes Problem . . . . .</b>	<b>105</b>
3.1	Problem Formulation . . . . .	105
3.2	Numerical Results . . . . .	108

---

The efficient resolution of parametrized PDEs is of great interest for both industry and academic research. With the development of super calculators and suitable methods, it is now possible to perform simulations of complex 3D models with high accuracy. However, the needs of high-speed simulations are also growing. For instance in the context of risk analysis or uncertainty quantification where thousands of simulations are needed. For those many-query methods, it is not realistic to perform a full order resolution (FEM, FVM, ... ) of an industrial problem, for each new set of parameters. A solution is to replace the full order model with a reduced one. Different methods of reduction exist: surrogate models, proper generalized decomposition, proper orthogonal decomposition, reduced basis, ...

In this work, we chose to focus on the Reduced Basis Method (RBM) which presents some interesting characteristics:

- Rigorous a posteriori error estimators
- An efficient Offline/Online decomposition.

Those features are provided by two necessary ingredients: a Galerkin projection onto a low-dimensional space spanned by properly selected basis functions and an affine parametric dependency.

The development of the RBM started in the late 70's, initially for non-linear mechanics problems [76]. The method was then extended to many different kinds of problems [2, 32, 78, 84]. The RBM knew in the early 2000's an important improvement with the development of the Certified Reduced Basis Method (CRBM) [80, 103, 81, 83, 88]. The main improvement of the CRBM is the use of efficient error estimators, evaluated through an offline/online procedure. Those error estimators allow quantifying the error made with the reduced model. This error is usually measured on a quantity of interest, called the output of the model. A direct consequence of this fast error evaluation is the possibility to explore the parameter space efficiently and to optimally select the parameters used to build the reduced space  $X_N$ .

This chapter is split into three sections. The basics of the method are introduced in the first section using the convenient example of the linear coercive problems. In the second section, we will introduce some specific features, required for the reduction of saddle point models. In the last part we detail the algorithm applied to aerothermal simulations.

## 1 Method and Notations : Application to Coercive Problems

We first introduce the reduced basis method (RBM) in a very general context.

Let  $\Omega \in \mathbb{R}^d$ ,  $d = 2, 3$  be a bounded domain and  $X \subset H^1(\Omega)$  a Hilbert space on this domain, we denote by  $(\cdot, \cdot)_X$  and  $\|\cdot\|_X = \sqrt{(\cdot, \cdot)_X}$  respectively the inner product and associated norm on  $X$ . We also introduce a parameter space  $\mathcal{D} \in \mathbb{R}^P$  of dimension  $P \geq 1$ . Let us consider, for a first time, a very generic problem, depending on a parameter  $\boldsymbol{\mu} \in \mathcal{D}$ , as finding  $u(\boldsymbol{\mu}) \in X$  such that

$$\mathcal{F}(u(\boldsymbol{\mu}); \boldsymbol{\mu}) = 0. \quad (5.1)$$

In the context of parametrized PDE, the resolution of this problem using the FEM as discretization method will lead to the resolution of the discrete problem: find

$u_h(\boldsymbol{\mu}) \in X_h$  such that

$$\mathcal{F}_h(u_h(\boldsymbol{\mu}); \boldsymbol{\mu}) = 0 \quad (5.2)$$

where  $X_h \subset X$  is a finite element function space of high dimension  $\mathcal{N}$ , as introduced in section 2.

The resolution of this discrete problem can be very costly, even with an optimized method and comfortable computation capabilities. The idea of the reduced basis method, and more generally of the different model order reduction methods is to build a low (relatively to  $\mathcal{N}$ ) dimensional space in which the resolution of problem (5.2) will give an acceptable approximation.

Note that when we use the RBM, we make a non-trivial assumption on the manifold

$$\mathcal{S}(\mathcal{D}) = \{u(\boldsymbol{\mu}) | \mathcal{F}(u(\boldsymbol{\mu}); \boldsymbol{\mu}) = 0, \boldsymbol{\mu} \in \mathcal{D}\}, \quad (5.3)$$

which has to be approximable by a finite-dimensional function space  $X_N$  of relatively small dimension  $N$ . The Kolmogorov  $N$ -width determines the feasibility of such an approximation. We will not detail this notion in this work, but we assume that, for the considered problems, the Kolmogorov  $N$ -width will decay sufficiently rapidly and then it makes sense to consider a reduced model.

In the context of RBM, this space  $X_N$  is usually constructed as

$$X_N = \text{span}\{u_h(\boldsymbol{\mu}), \boldsymbol{\mu} \in S_N\} \quad (5.4)$$

where  $S_N \subset \mathcal{D}$  is a sample of parameters chosen in  $\mathcal{D}$ . The construction of this sample can be more or less optimized and will determine the quality of the reduced basis space  $X_N$ . For instance, we can imagine choosing  $S_N$  as a random or an equidistributed subset of  $\mathcal{D}$ , but more efficient techniques have been developed using, for instance, a greedy algorithm to optimally select the parameters in  $S_N$ .

## Parameters Selection through Greedy Algorithm

The greedy algorithm builds a sample  $S_N$  in order to minimize a quantity  $\Delta_N(\boldsymbol{\mu})$  on  $\mathcal{D}$ . For instance, this quantity could be the error,  $\tilde{\Delta}_N(\boldsymbol{\mu}) = \|u_h(\boldsymbol{\mu}) - u_N(\boldsymbol{\mu})\|$ , between the FEM solution  $u_h(\boldsymbol{\mu})$  and the reduced solution  $u_N(\boldsymbol{\mu})$ . But the algorithm requires the evaluation of the maximum of this quantity on the entire parameter space  $\mathcal{D}$ ,  $\tilde{\Delta}_N = \max_{\boldsymbol{\mu} \in \mathcal{D}} \tilde{\Delta}_N(\boldsymbol{\mu})$ . Such maximum is obviously not computable.

Hence we introduce two significant simplifications. First, we compute this maximum on a discrete subset  $\Xi \subset \mathcal{D}$  (potentially very large). Then, we do not use the real error  $\tilde{\Delta}_N(\boldsymbol{\mu})$  but error indicators  $\Delta_N(\boldsymbol{\mu})$  which can be evaluated for many values  $\boldsymbol{\mu}$  for a minimal cost. Alternatively, we might be more interested in the error measured for a specific quantity of interest: the output of the model. The greedy construction remains the same using an output-based error bound. The nature of these error indicators is problem-dependent and will be detailed further for the different kinds of considered PDE.

We present in the algorithm 2 a generic version of the greedy algorithm. The general idea of the algorithm does not depend on the nature of the problems.

---

**Algorithm 2: Generic Greedy Algorithm**


---

```

1 Choose  $\Xi \subset \mathcal{D}$ ; // Super Sampling
2 Choose  $N_{\max}$ ; // Maximum size of  $X_N$ 
3 Choose  $\delta_{tol} \in (0, 1)$ ; // Algorithm tolerance
4 Choose  $\mu_1 \in \mathcal{D}$ ; // Arbitrary chosen first parameter
5 Set  $N \leftarrow 0, S_N \leftarrow \{\}, X_N \leftarrow \{\}$ ; // Initialization
6 repeat
7    $N \leftarrow N + 1$ ;
8    $S_N \leftarrow S_{N-1} \cup \{\mu_N\}$ ;
9   Solve  $\mathcal{F}_h(u_h(\mu); \mu)$ ;
10   $X_N \leftarrow X_{N-1} \oplus \text{span}\{u_h(\mu_N)\}$ ;
11  for  $\mu \in \Xi$  do
12    | Compute error indicator,  $\Delta_N(\mu)$ ;
13  end
14   $\mu_{N+1} = \arg \max_{\mu \in \Xi} \Delta_N(\mu)$ ;
15 until  $\Delta_N(\mu_{N+1}) < \delta_{tol}$  or  $N \geq N_{\max}$ ;

```

---

## 1.1 Reduced Basis for Coercive Problems

We detail in this section the use of the RBM for coercive problems. The problem is defined as in the previous section. We also introduce a parametrized bilinear form  $a(\cdot, \cdot; \mu) : X \times \mathcal{D} \rightarrow \mathbb{R}$  and a parametrized linear form  $f(\cdot; \mu) : X \times \mathcal{D} \rightarrow \mathbb{R}$ . The problem (5.1) can be written as the variational formulation of a parametrized PDE, then we want to find  $u(\mu) \in X$  such that

$$a(u(\mu), v; \mu) = f(v; \mu), \quad \forall v \in X. \quad (5.5)$$

We also do some assumptions on this parametrized weak formulation, for any value of  $\mu \in \mathcal{D}$ :

- (a)  $a(\cdot, \cdot; \mu)$  is symmetric and positive definite. We can then define a scalar product and the associated norm

$$(\cdot, \cdot)_\mu = a(\cdot, \cdot; \mu), \quad \|\cdot\|_\mu = \sqrt{(\cdot, \cdot)_\mu}. \quad (5.6)$$

In practice we will essentially use the scalar product  $(\cdot, \cdot)_{\bar{\mu}}$  and norm  $\|\cdot\|_{\bar{\mu}}$  associated to a reference parameter  $\bar{\mu}$ .

- (b)  $a(\cdot, \cdot; \mu)$  is continuous and coercive with respect to the norm  $\|\cdot\|_X$ , i.e., there exist a positive coercivity constant  $\alpha(\mu) > 0$  and a finite continuity constant

$\gamma(\boldsymbol{\mu}) < \infty$  such that

$$a(v, v; \boldsymbol{\mu}) \geq \alpha(\boldsymbol{\mu}) \|v\|_X^2, \quad a(w, v; \boldsymbol{\mu}) \leq \gamma(\boldsymbol{\mu}) \|w\|_X \|v\|_X, \quad \forall v, w \in X \quad (5.7)$$

These constants are defined respectively as

$$\alpha(\boldsymbol{\mu}) = \inf_{v \in X} \frac{a(v, v; \boldsymbol{\mu})}{\|v\|_X^2}, \quad \gamma(\boldsymbol{\mu}) = \sup_{w \in X} \sup_{v \in X} \frac{a(w, v; \boldsymbol{\mu})}{\|v\|_X \|w\|_X}, \quad \forall \boldsymbol{\mu} \in \mathcal{D} \quad (5.8)$$

(c) The forms  $a(\cdot, \cdot; \boldsymbol{\mu})$  and  $f(\cdot; \boldsymbol{\mu})$  have affine decompositions

$$a(w, v; \boldsymbol{\mu}) = \sum_{q=1}^{Q_a} \theta_q^a(\boldsymbol{\mu}) a_q(w, v), \quad (5.9)$$

$$f(v; \boldsymbol{\mu}) = \sum_{q=1}^{Q_f} \theta_q^f(\boldsymbol{\mu}) f_q(v) \quad (5.10)$$

where each form

$$a_q : X \times X \rightarrow \mathbb{R}, \quad f_q : X \rightarrow \mathbb{R}$$

is independent of parameter  $\boldsymbol{\mu}$  and each

$$\theta_q^a : \mathcal{D} \rightarrow \mathbb{R}, \quad \theta_q^f : \mathcal{D} \rightarrow \mathbb{R}$$

is a function which only depends on  $\boldsymbol{\mu}$ .

This affine decomposition is essential to provide an efficient offline/online strategy. This point is detailed in subsection 1.1.

### Offline Procedure: Construction of the RB Space

We construct the reduced basis space using the greedy algorithm as proposed in section 1. We assume that we have a good error bound  $\Delta_N(\boldsymbol{\mu})$ , defined on a quantity of interest. For coercive problems, the evaluation of this error bound is detailed in section 1.2.

The first parameter  $\boldsymbol{\mu}_1$  to initialize the greedy procedure is usually arbitrarily chosen. With  $\boldsymbol{\mu}_1$  we obtain the first snapshot  $u_h(\boldsymbol{\mu}_1)$  which is not directly used as basis vector, but normalized, with respect to the specific norm  $\|\cdot\|_{\bar{\boldsymbol{\mu}}}$ , as  $\xi_1 = \frac{u_h}{\|u_h\|_{\bar{\boldsymbol{\mu}}}}$ . The reduced basis space  $X_1$  is then defined as  $X_1 = \text{span}\{\xi_1\}$ .

After this initialization, at any rank  $N$ , we use the greedy algorithm (2) to chose the next parameter  $\boldsymbol{\mu}_{N+1}$  and then the next snapshot  $u_h(\boldsymbol{\mu}_{N+1})$ . The next basis vector  $\xi_{N+1}$  is obtained using the Gram-Schmidt process on the set of vectors



$(\xi_1, \dots, \xi_N, u(\boldsymbol{\mu}_{N+1}))$ , still with respect to the norm  $\|\cdot\|_{\bar{\mu}}$ . And eventually, the reduced space is defined as

$$X_{N+1} = \text{span}\{\xi_n, 1 \leq n \leq N + 1\}. \quad (5.11)$$

The recursive construction ends when either the expected error or the max size  $N$  is reached, as proposed in the algorithm (2). This way, a nested set of reduced basis spaces is built,  $X_1 \subset X_2 \subset \dots \subset X_N$ . Note that even if the orthonormalization process on the basis is not mandatory, it allows controlling the conditioning of the reduced matrix obtained after the Galerkin projection. Otherwise, this matrix can be very ill-conditioned, and it generates numerical instability.

The reduced solution  $u_N(\boldsymbol{\mu})$ , on  $X_N$ , can then be written as

$$u_N(\boldsymbol{\mu}) = \sum_{i=1}^N u_{Ni}(\boldsymbol{\mu})\xi_i \quad (5.12)$$

the coefficient  $u(\boldsymbol{\mu})_{Ni}$  are evaluated during the online phase of the RBM, detailed in next paragraph.

### Online Evaluation : Computation of $u_N(\boldsymbol{\mu})$

Once the reduced space is built we can use this structure to solve the problem introduced in (5.5) on  $X_N$ , for any new parameter  $\boldsymbol{\mu} \in \mathcal{D}$ . The reduced solution  $u_N(\boldsymbol{\mu})$  is then the Galerkin projection of the continuous solution  $u(\boldsymbol{\mu})$  on the reduced basis  $\{\xi_i, 1 \leq i \leq N\}$ . This projection implies the resolution of the linear problem

$$A_N(\boldsymbol{\mu})u_N = F_N(\boldsymbol{\mu}) \quad (5.13)$$

where

$$(A_N(\boldsymbol{\mu}))_{ij} = a(\xi_i, \xi_j; \boldsymbol{\mu}), \quad (F_N(\boldsymbol{\mu}))_i = f(\xi_i; \boldsymbol{\mu}), \quad \forall 1 \leq i, j \leq N. \quad (5.14)$$

The size  $N \times N$  of this problem is supposed to be very small compared to the finite element dimension  $\mathcal{N}$ . The order of the model is reduced as expected. However the assembly of the reduced matrix  $A_N(\boldsymbol{\mu})$  and the reduced vector  $F_N(\boldsymbol{\mu})$  still involve the finite element dimension through the evaluation of the quantities  $a(\xi_i, \xi_j; \boldsymbol{\mu})$  and  $f(\xi_i; \boldsymbol{\mu})$ . To be efficient the online part has to be totally independent of the full dimension  $\mathcal{N}$ . The assumption of affine decomposition introduced in (5.9) allows to decompose the evaluation of  $A_N(\boldsymbol{\mu})$  and  $F_N(\boldsymbol{\mu})$  in a online/offline procedure, with a  $\mathcal{N}$ -independent online part.

Using (5.9) we can rewrite (5.14)

$$\begin{aligned}
 (A_N(\boldsymbol{\mu}))_{ij} &= a(\xi_i, \xi_j; \boldsymbol{\mu}), & (F_N(\boldsymbol{\mu}))_i &= f(\xi_i; \boldsymbol{\mu}), \\
 &= \sum_{q=1}^{Q_a} \theta_q^a(\boldsymbol{\mu}) a_q(\xi_i, \xi_j), & &= \sum_{q=1}^{Q_f} \theta_q^f(\boldsymbol{\mu}) f_q(\xi_i) \\
 &= \sum_{q=1}^{Q_a} \theta_q^a(\boldsymbol{\mu}) \underbrace{(A_N^q)_{ij}}_{\text{precomputed}}, & &= \sum_{q=1}^{Q_f} \theta_q^f(\boldsymbol{\mu}) \underbrace{(F_N^q)_i}_{\text{precomputed}}, \quad \forall 1 \leq i, j \leq N.
 \end{aligned} \tag{5.15}$$

Thanks to the affine decomposition we can now precompute the matrices  $A_N^q$  and the vectors  $F_N^q$  during the offline phase, since they do not depend of the parameter  $\boldsymbol{\mu}$ . Hence the online phase only consist in a sum of  $Q_a$   $N \times N$ -matrices, a sum of  $Q_f$   $N$ -size vectors and the resolution of the  $N \times N$  linear system.

During the offline phase, we have to introduce a new step in the construction of the nested RB spaces: the update of the precomputed  $A_N^q$  and  $F_N^q$ . Thanks to the nested structure of the RB spaces  $X_N$ , we only have to evaluate the last row, and the last column of the  $A_N^q$  matrices, the other values are similar to  $A_{N-1}^q$ . The same argument allows only to compute the last value of the vectors  $F_N^q$ .

## 1.2 Error Estimation

The computation of an efficient error bound for coercive problems has been widely studied [80, 21] during last decades, we remind in this section the main ingredients. In the context of the parametrized PDE, we may be interested either in the solution itself or by the evaluation of a particular quantity of interest (output). In this latter situation, we need some assumption on the considered quantity,  $s(\boldsymbol{\mu})$ .

We assume the existence of a parametrized linear form  $l(\cdot; \boldsymbol{\mu}) : X \times \mathcal{D} \rightarrow \mathbb{R}$  such that the output  $s(\boldsymbol{\mu})$  can be expressed as

$$s(\boldsymbol{\mu}) = l(u(\boldsymbol{\mu}); \boldsymbol{\mu}). \tag{5.16}$$

We furthermore assume that the form  $l$  admits an affine decomposition

$$l(v; \boldsymbol{\mu}) = \sum_{q=1}^{Q_l} \theta_q^l(\boldsymbol{\mu}) l_q(v) \tag{5.17}$$

with the conditions stated in (5.9).

For an efficient error bound of the error  $|s_h(\boldsymbol{\mu}) - s_N(\boldsymbol{\mu})|$  we need to introduce a dual problem associated to the linear form  $l(\cdot; \boldsymbol{\mu})$ , as find  $\psi(\boldsymbol{\mu}) \in X$  such that

$$a(v, \psi(\boldsymbol{\mu}); \boldsymbol{\mu}) = -l(v; \boldsymbol{\mu}), \quad \forall v \in X. \tag{5.18}$$

For this dual problem we build a dual reduced space

$$X_N^{du} = \text{span}\{\xi_n^{du}, 1 \leq n \leq N\} \quad (5.19)$$

where the basis vectors  $\xi_n^{du}$  are generated through a Gram-Schmidt algorithm on the snapshots  $\psi_h(\boldsymbol{\mu}_n)$  as we did for the construction of the primal reduced space  $X_N$ .

Once this dual problem and reduced space are introduced, we can define the first main ingredient for the error bound: the dual norm of the residuals associated with the primal and dual problem, respectively  $\epsilon_N(\boldsymbol{\mu})$  and  $\epsilon_N^{du}(\boldsymbol{\mu})$ , defined as

$$\epsilon_N(\boldsymbol{\mu}) = \sup_{v \in X_h} \frac{r_N(v; \boldsymbol{\mu})}{\|v\|_{\bar{\mu}}} = \|\hat{e}_N(\boldsymbol{\mu})\|_{\bar{\mu}} \quad (5.20)$$

$$\epsilon_N^{du}(\boldsymbol{\mu}) = \sup_{v \in X_h} \frac{r_N^{du}(v; \boldsymbol{\mu})}{\|v\|_{\bar{\mu}}} = \|\hat{e}_N^{du}(\boldsymbol{\mu})\|_{\bar{\mu}} \quad (5.21)$$

where

$$r_N(v; \boldsymbol{\mu}) = f(v; \boldsymbol{\mu}) - a(u_N(\boldsymbol{\mu}), v; \boldsymbol{\mu}) \quad (5.22)$$

$$r_N^{du}(v; \boldsymbol{\mu}) = l(v; \boldsymbol{\mu}) + a(v, \psi_N(\boldsymbol{\mu}); \boldsymbol{\mu}), \quad \forall v \in X_h \quad (5.23)$$

are the residuals, respectively of the primal and dual problems.  $\hat{e}_N(\boldsymbol{\mu}) \in X_h$  and  $\hat{e}_N^{du}(\boldsymbol{\mu}) \in X_h$  are their Riesz representations which satisfy

$$(\hat{e}_N(\boldsymbol{\mu}), v)_{\bar{\mu}} = r_N(v; \boldsymbol{\mu}) \quad (5.24)$$

$$(\hat{e}_N^{du}(\boldsymbol{\mu}), v)_{\bar{\mu}} = r_N^{du}(v; \boldsymbol{\mu}), \quad \forall v \in X_h \quad (5.25)$$

the evaluation of these Riesz representations through a efficient offline/online procedure is detailed in the appendix 3.1.

The second ingredient for the error bound are lower bound  $\alpha_{LB}(\boldsymbol{\mu})$  and an upper bound  $\alpha_{UB}(\boldsymbol{\mu})$  for the coercivity constant introduce in (5.8)

$$\alpha_{LB}(\boldsymbol{\mu}) \leq \alpha(\boldsymbol{\mu}) \leq \alpha_{UB}(\boldsymbol{\mu}), \quad \forall \boldsymbol{\mu} \in \mathcal{D}. \quad (5.26)$$

The evaluation of these bounds has to be independent of the FE dimension  $\mathcal{N}$ . The successive constraint method (SCM) [47, 89, 101] has been developed specifically to provide an efficient computation of such bounds for coercive or inf-sup problems. In the context of coercive problems, we can also use the Min- $\theta$  approach [63, 104]. We did not detail these methods in this report and refer the interested reader to the cited publications.

We can now introduce the error bounds for the errors on the reduced solution  $\Delta_N(\boldsymbol{\mu})$  and for the output  $\Delta_N^s(\boldsymbol{\mu})$  defined as

$$\Delta_N(\boldsymbol{\mu}) = \frac{\epsilon_N(\boldsymbol{\mu})}{(\alpha_{LB}(\boldsymbol{\mu}))^{1/2}}, \quad (5.27)$$

$$\Delta_N^{du}(\boldsymbol{\mu}) = \frac{\epsilon_N^{du}(\boldsymbol{\mu})}{(\alpha_{LB}(\boldsymbol{\mu}))^{1/2}}, \quad (5.28)$$

$$\Delta_N^s(\boldsymbol{\mu}) = \Delta_N(\boldsymbol{\mu})\Delta_N^{du}(\boldsymbol{\mu}), \quad (5.29)$$

and satisfying

$$\|u_h(\boldsymbol{\mu}) - u_N(\boldsymbol{\mu})\|_{\bar{\boldsymbol{\mu}}} \leq \Delta_N(\boldsymbol{\mu}) \quad (5.30)$$

$$|s_h(\boldsymbol{\mu}) - s_N(\boldsymbol{\mu})| \leq \Delta_N^s(\boldsymbol{\mu}). \quad (5.31)$$

For a study of effectiveness of these error estimators, we refer the reader to [80, 89, 9]. Once we have an efficient procedure for the evaluation of the error bounds, it can be used within the greedy algorithm 2 in order to guide the parameters selection. The evaluation of the reduced output using a dual problem (during the online part) is then made through

$$s_N(\boldsymbol{\mu}) = l(u_N(\boldsymbol{\mu}); \boldsymbol{\mu}) - r_N(\psi_N(\boldsymbol{\mu}); \boldsymbol{\mu}), \quad \forall \boldsymbol{\mu} \in \mathcal{D} \quad (5.32)$$

**Remark 6** (Error bound for compliant problems). *A compliant problem is defined when the linear form defining the output also defines the right hand side of the variational formulation (5.5) :  $l(\cdot; \boldsymbol{\mu}) = f(\cdot; \boldsymbol{\mu})$ . In this case we don't need the resolution of the dual problem. The error bound on the output is simply defined as  $\Delta_N^s(\boldsymbol{\mu}) = (\Delta_N(\boldsymbol{\mu}))^2$*

**Remark 7** (Problem with no output). *When we do not need the evaluation of a particular output or when we need to evaluate many outputs, we will directly use the error bound on the solution  $\Delta_N(\boldsymbol{\mu})$  instead of an error bound on the output. In that case, we will not use the dual problem but we lose the quadratic convergence of the reduced solution. Another solution is proposed in [109] for a generalization to vector valued quantities of interest.*

The two ingredients are evaluated through an efficient offline/online strategy. The coercivity constant can be computed using the Successive Constraint Method (SCM) [47, 89, 101]. The dual norm of the residual is partly precomputed offline, using the affine decomposition. This procedure is detailed in appendix 3.1.

## Summarized Offline Algorithm for Coercive Problems

We complete the generic algorithm 2 with the different procedures introduced in this section. Some steps might not be necessary, for instance, if we do not need the

dual problem.

---

**Algorithm 3:** Offline Procedure for Coercive Problems

---

```

1 Choose  $\Xi \subset \mathcal{D}$ ; // Super Sampling
2 Choose  $N_{\max}$ ; // Maximum size of  $X_N$ 
3 Choose  $\delta_{tol}$ ; // Algorithm tolerance
4 Choose  $\mu_1 \in \mathcal{D}$ ; // Arbitrary chosen first parameter
5 Set  $N \leftarrow 0, S_N \leftarrow \{\}, X_N \leftarrow \{\}, X_N^{du} \leftarrow \{\}$ ; // Initialization
6 repeat
7    $N \leftarrow N + 1$ ;
8    $S_N \leftarrow S_{N-1} \cup \mu_N$ ;
9   Solve  $a(u, v; \mu_N) = f(v; \mu_N)$ ; // (5.5)
10   $\xi_N \leftarrow \text{Gram-Schmidt}(\{\xi_1, \dots, \xi_{N-1}, u(\mu_N)\})$ ;
11   $X_N \leftarrow X_{N-1} \oplus \text{span}\{\xi_N\}$ ;
12  Solve dual problem  $a(v, \psi; \mu_N) = -l(v; \mu)$ ; // (5.17)
13   $\xi_N^{du} \leftarrow \text{Gram-Schmidt}(\{\xi_1^{du}, \dots, \xi_{N-1}^{du}, \psi(\mu_N)\})$ ;
14   $X_N^{du} \leftarrow X_{N-1}^{du} \oplus \text{span}\{\xi_N^{du}\}$ ;
15  Precompute  $A_N^q$  and  $F_N^q$ ; // (5.15)
16  Precompute  $(\hat{f}_q, \hat{f}_{q'}), C_{qq',N}^{fa}, C_{qq',N}^{aa}$ ; // (43)
17  for  $\mu \in \Xi$  do
18    | Compute error bound,  $\Delta_N(\mu)$  or  $\Delta_N^s(\mu)$ ;
19  end
20   $\mu_{N+1} = \arg \max_{\mu \in \Xi} \Delta_N(\mu)$  or  $\mu_{N+1} = \arg \max_{\mu \in \Xi} \Delta_N^s(\mu)$ ;
21 until  $\Delta_N^{(s)}(\mu_{N+1}) < \delta_{tol}$  or  $N \geq N_{\max}$ ;

```

---

## 2 Reduced Basis for Saddle Point Problems

In this section, we present the RBM applied to saddle point problems, and we will especially focus on the Stokes equations as an illustration. The RBM is very well developed for coercive problems, but saddle point parametrized systems introduce new difficulties with the construction of stable reduced spaces and the development of efficient a posteriori error bounds [66, 90, 86]. The delicate problem of geometric parameters is particularly developed in [34] with a penalty method and later in [35, 87] with specific construction algorithm of the reduced space.

We propose a quick review of the reduced basis framework for saddle point problem. We present the problem formulation, the existing methods for the construction of a stable reduced space and the detailed error estimators. We provide more technical details about the integration of these features in the Feel++ framework in the implementation part III.

## 2.1 Problem formulation

Let  $V$  and  $Q$  two Hilbert spaces on a bounded domain  $\Omega$  with inner products  $(\cdot, \cdot)_V$ ,  $(\cdot, \cdot)_Q$  and associated norms  $\|\cdot\|_V = \sqrt{(\cdot, \cdot)_V}$ ,  $\|\cdot\|_Q = \sqrt{(\cdot, \cdot)_Q}$  respectively. We also define the product space  $Z = V \times Q$  with inner product  $((\mathbf{u}, p), (\mathbf{v}, q))_Z = (\mathbf{u}, \mathbf{v})_V + (p, q)_Q$  and norm  $\|\cdot\|_Z = \sqrt{(\cdot, \cdot)_Z}$ . We consider a parametrized saddle point problem: for any  $\boldsymbol{\mu} \in \mathcal{D}$ , find  $\mathbf{u}(\boldsymbol{\mu}) \in V$  and  $p(\boldsymbol{\mu}) \in Q$  such that

$$\begin{aligned} a(\mathbf{u}(\boldsymbol{\mu}), \mathbf{v}; \boldsymbol{\mu}) + b(\mathbf{v}, p(\boldsymbol{\mu}); \boldsymbol{\mu}) &= f(\mathbf{v}; \boldsymbol{\mu}), & \forall \mathbf{v} \in V \\ b(\mathbf{u}(\boldsymbol{\mu}), q; \boldsymbol{\mu}) &= g(q; \boldsymbol{\mu}) & \forall q \in Q \end{aligned} \quad (5.33)$$

where  $a(\cdot, \cdot; \boldsymbol{\mu}) : V \times V \rightarrow \mathbb{R}$ ,  $b(\cdot, \cdot; \boldsymbol{\mu}) : V \times Q \rightarrow \mathbb{R}$  and  $f(\cdot; \boldsymbol{\mu}) : V \rightarrow \mathbb{R}$ ,  $g(\cdot; \boldsymbol{\mu}) : Q \rightarrow \mathbb{R}$  are bilinear and linear bounded forms respectively, for any  $\boldsymbol{\mu} \in \mathcal{D}$ .

We furthermore assume that :

- (a)  $a(\cdot, \cdot; \boldsymbol{\mu})$  and  $b(\cdot, \cdot; \boldsymbol{\mu})$  are continuous

$$\gamma_a(\boldsymbol{\mu}) = \sup_{\mathbf{u} \in V} \sup_{\mathbf{v} \in V} \frac{a(\mathbf{u}, \mathbf{v}; \boldsymbol{\mu})}{\|\mathbf{u}\|_V \|\mathbf{v}\|_V} < \infty, \quad \forall \boldsymbol{\mu} \in \mathcal{D}, \quad (5.34)$$

$$\gamma_b(\boldsymbol{\mu}) = \sup_{\mathbf{v} \in V} \sup_{q \in Q} \frac{b(\mathbf{v}, q; \boldsymbol{\mu})}{\|\mathbf{v}\|_V \|q\|_Q} < \infty, \quad \forall \boldsymbol{\mu} \in \mathcal{D}, \quad (5.35)$$

- (b)  $a(\cdot, \cdot; \boldsymbol{\mu})$  is coercive on  $V$

$$\alpha_a(\boldsymbol{\mu}) = \inf_{\mathbf{v} \in V} \frac{a(\mathbf{v}, \mathbf{v}; \boldsymbol{\mu})}{\|\mathbf{v}\|_V^2} > 0, \quad \forall \boldsymbol{\mu} \in \mathcal{D}, \quad (5.36)$$

- (c)  $b(\cdot, \cdot; \boldsymbol{\mu})$  satisfies the inf-sup condition, this can be expressed either using the Brezzi constant,

$$\beta_{Br}(\boldsymbol{\mu}) = \inf_{q \in Q} \sup_{\mathbf{v} \in V} \frac{b(\mathbf{v}, q; \boldsymbol{\mu})}{\|\mathbf{v}\|_V \|q\|_Q} > 0, \quad \forall \boldsymbol{\mu} \in \mathcal{D}, \quad (5.37)$$

or the Babuška constant

$$\beta_{Ba}(\boldsymbol{\mu}) = \inf_{(\mathbf{u}, p) \in Z} \sup_{(\mathbf{v}, q) \in Z} \frac{a(\mathbf{u}, \mathbf{v}; \boldsymbol{\mu}) + b(\mathbf{v}, p; \boldsymbol{\mu}) + b(\mathbf{u}, q; \boldsymbol{\mu})}{\|(\mathbf{u}, p)\|_Z \|(\mathbf{v}, q)\|_Z} > 0, \quad \forall \boldsymbol{\mu} \in \mathcal{D}, \quad (5.38)$$

- (d) The bilinear forms  $a(\cdot, \cdot; \boldsymbol{\mu})$ ,  $b(\cdot, \cdot; \boldsymbol{\mu})$  and the linear forms  $f(\cdot; \boldsymbol{\mu})$  and  $g(\cdot; \boldsymbol{\mu})$  admit affine decompositions

$$a(\mathbf{u}, \mathbf{v}; \boldsymbol{\mu}) = \sum_{k=1}^{Q_a} \theta_k^a(\boldsymbol{\mu}) a_k(\mathbf{u}, \mathbf{v}), \quad f(\mathbf{v}; \boldsymbol{\mu}) = \sum_{k=1}^{Q_f} \theta_k^f(\boldsymbol{\mu}) f_k(\mathbf{v}) \quad (5.39)$$

$$b(\mathbf{u}, q; \boldsymbol{\mu}) = \sum_{k=1}^{Q_b} \theta_k^b(\boldsymbol{\mu}) b_k(\mathbf{u}, q), \quad g(q; \boldsymbol{\mu}) = \sum_{k=1}^{Q_g} \theta_k^g(\boldsymbol{\mu}) g_k(q) \quad (5.40)$$

### Truth Approximation

As truth approximation of the continuous problem, we use the finite element solution  $(\mathbf{u}_h, p_h)$  of (5.33) on a regular triangulation  $\mathcal{T}_h$  of the domain  $\Omega$ . To that end, we introduce finite element function spaces  $V_h, Q_h$  and  $Z_h$  and their respective scalar products and norms. The coercivity and continuity properties of the forms are naturally inherited from the continuous formulation. The spaces  $V_h$  and  $Q_h$  are supposed to be properly defined to satisfy the Ladyzhenskaya–Babuška–Brezzi condition. More details about the finite element resolution are given in section 2.

### Offline Procedure : Construction of the RB Space

The classical procedure for the construction of the RB space, as seen in section 1.1, is not adapted for saddle point problems. This procedure would build the RB spaces  $V_N^0$  and  $Q_N$  as spans of the truth solutions  $\{\mathbf{u}(\boldsymbol{\mu}_n)\}$  and  $\{p(\boldsymbol{\mu}_n)\}$ ,  $1 \leq n \leq N$ , respectively,

$$V_N^0 = \text{span}\{\mathbf{u}(\boldsymbol{\mu}_n), 1 \leq n \leq N\} = \text{span}\{\boldsymbol{\zeta}_n, 1 \leq n \leq N\}, \quad (5.41)$$

$$Q_N = \text{span}\{p(\boldsymbol{\mu}_n), 1 \leq n \leq N\} = \text{span}\{\eta_n, 1 \leq n \leq N\}. \quad (5.42)$$

With such RB spaces, the online saddle-point system would not be well-posed, since it does not satisfy the inf-sup condition, see [5, 6],

$$\beta_N(\boldsymbol{\mu}) = \inf_{q_N \in Q_N} \sup_{\mathbf{v}_N \in V_N} \frac{b(\mathbf{v}_N, q_N; \boldsymbol{\mu})}{\|q_N\|_Q \|\mathbf{v}_N\|_V} > 0. \quad (5.43)$$

It is proved in [85, 90] that the space  $V_N^0 = \text{span}\{\mathbf{u}(\boldsymbol{\mu}_n), 1 \leq n \leq N\}$  can be enriched in order to obtain a stable pair  $(V_N^1, Q_N)$  in the sense of the inf-sup condition. This space  $V_N^1$  is built by adding the supremizer functions  $T^k \eta_n$  for any  $1 \leq k \leq Q_b$ ,  $1 \leq n \leq N$ .

$$V_N^1 = V_N^0 \oplus \text{span}\{T^k \eta_n, 1 \leq k \leq Q_b, 1 \leq n \leq N\}. \quad (5.44)$$

The supremizer function  $T^k q$  is defined for any  $q \in Q$  and any  $1 \leq k \leq Q_b$  as

$$(T^k q, \mathbf{v}_h)_V = b^k(\mathbf{v}_h, q), \quad \forall \mathbf{v}_h \in V_h. \quad (5.45)$$

In practice, it is widespread not to build the RB space  $V_N^1$  which might become very large, depending on  $Q_b$ . The usual alternative is to add the *global* supremizer  $T^{\boldsymbol{\mu}_n} \eta_n$ , for any parameters  $\boldsymbol{\mu}_n \in S_N$ ,

$$V_N = V_N^0 \oplus \text{span}\{T^{\boldsymbol{\mu}_n} \eta_n, 1 \leq n \leq N\}, \quad (5.46)$$

where  $T^\mu q$  is defined, for any  $\boldsymbol{\mu} \in \mathcal{D}$  and any  $q \in Q$  as

$$(T^\mu q, \mathbf{v}_h)_V = b(\mathbf{v}_h, q; \boldsymbol{\mu}), \quad \forall \mathbf{v}_h \in V_h. \quad (5.47)$$

With this definition of  $V_N$  it is no longer possible to prove the stability of the pair  $(V_N, Q_N)$  *a priori*. Nevertheless, different numerical tests conclude in the stability of this option, see [34, 35].

The basis  $\{\eta_n, 1 \leq n \leq N\}$  of the RB space  $Q_N$  is obtained, just like in the coercive case, by an orthonormalization of the snapshots  $p(\boldsymbol{\mu}_n), 1 \leq n \leq N$ . We apply the same process to the RB space  $V_N$  to recover an orthonormal basis  $\{\zeta_n, 1 \leq n \leq 2N\}$ . The complete offline procedure implemented for saddle point problems is summarized in algorithm 4. The nature of the error indicator  $\Delta_N(\boldsymbol{\mu})$  used in the greedy selection is discussed in section 2.2.

---

**Algorithm 4:** Offline Procedure for Saddle Point Problems

---

```

1 Choose  $\Xi \subset \mathcal{D}$  ; // Super Sampling
2 Choose  $N_{\max}$  ; // Maximum size of  $X_N$ 
3 Choose  $\delta_{tol}$  ; // Algorithm tolerance
4 Set  $N \leftarrow 0, S_N \leftarrow \{\}, V_N \leftarrow \{\}, Q_N \leftarrow \{\}$  ; // Initialization
5 repeat
6    $N \leftarrow N + 1$ ;
7    $S_N \leftarrow S_{N-1} \cup \{\boldsymbol{\mu}_N\}$ ;
8   Find  $(\mathbf{u}(\boldsymbol{\mu}), p(\boldsymbol{\mu}))$  solution of system (5.33);
9    $\eta_N \leftarrow \text{Gram-Schmidt}(\{\eta_1, \dots, \eta_{N-1}, p(\boldsymbol{\mu}_N)\})$ ;
10  Compute  $T^{\boldsymbol{\mu}_N} \eta_N$  ; // (5.47)
11   $\zeta_{2N-1}, \zeta_{2N} \leftarrow \text{Gram-Schmidt}(\{\zeta_1, \dots, \zeta_{2N-2}, \mathbf{u}(\boldsymbol{\mu}_N), T^{\boldsymbol{\mu}_N} \eta_N\})$ ;
12   $V_N \leftarrow V_{N-1} \oplus \text{span}\{\zeta_{2N-1}, \zeta_{2N}\}$ ;
13   $Q_N \leftarrow Q_{N-1} \oplus \text{span}\{\eta_N\}$ ;
14  Precompute online structures;
15  for  $\boldsymbol{\mu} \in \Xi$  do
16    | Compute error bound,  $\Delta_N(\boldsymbol{\mu})$ ;
17  end
18   $\boldsymbol{\mu}_{N+1} = \arg \max_{\boldsymbol{\mu} \in \Xi} \Delta_N(\boldsymbol{\mu})$ ;
19 until  $\Delta_N(\boldsymbol{\mu}_{N+1}) < \delta_{tol}$  or  $N \geq N_{\max}$ ;

```

---

Advanced algorithms exist in the literature to improve the convergence of the method. For instance, [35] and [87] propose specific enrichment of the RB space  $V_N$  to provide better error bounds. Unfortunately, we had no time to implement these methods properly, and we refer the readers to the cited references for more details.

**Online Evaluation : computation of  $(\mathbf{u}_N(\boldsymbol{\mu}), p_N(\boldsymbol{\mu}))$**

As for the coercive problem, we use a Galerkin projection of on the RB space  $Z_N = V_N \times Q_N$ . It implies the resolution of the linear system

$$\begin{bmatrix} \mathbf{A}_N(\boldsymbol{\mu}) & \mathbf{B}_N^T(\boldsymbol{\mu}) \\ \mathbf{B}_N(\boldsymbol{\mu}) & 0 \end{bmatrix} \begin{bmatrix} \mathbf{u}_N(\boldsymbol{\mu}) \\ p_N(\boldsymbol{\mu}) \end{bmatrix} = \begin{bmatrix} \mathbf{F}_N(\boldsymbol{\mu}) \\ \mathbf{G}_N(\boldsymbol{\mu}) \end{bmatrix} \quad (5.48)$$



where

$$\begin{aligned} (\mathbf{A}_N(\boldsymbol{\mu}))_{i,j} &= a(\boldsymbol{\zeta}_j, \boldsymbol{\zeta}_i; \boldsymbol{\mu}), & (\mathbf{F}(\boldsymbol{\mu}))_i &= f(\eta_i; \boldsymbol{\mu}), & 1 \leq i, j \leq 2N \\ (\mathbf{B}_N(\boldsymbol{\mu}))_{l,j} &= b(\boldsymbol{\zeta}_j, \eta_l; \boldsymbol{\mu}), & (\mathbf{G}_N(\boldsymbol{\mu}))_l &= g(\eta_l; \boldsymbol{\mu}), & 1 \leq l \leq N. \end{aligned} \quad (5.49)$$

Here again the affine assumption (5.39) is crucial in order to precompute these matrices and vectors.

$$\begin{aligned} (\mathbf{A}_N(\boldsymbol{\mu}))_{i,j} &= \sum_{k=1}^{Q_a} \theta_k^a(\boldsymbol{\mu}) \underbrace{a_k(\boldsymbol{\zeta}_j, \boldsymbol{\zeta}_i)}_{\text{precomputed}}, & (\mathbf{F}(\boldsymbol{\mu}))_i &= \sum_{k=1}^{Q_f} \theta_k^f \underbrace{f_k(\boldsymbol{\zeta}_i)}_{\text{precomputed}}, & 1 \leq i, j \leq 2N \\ (\mathbf{B}_N(\boldsymbol{\mu}))_{l,j} &= \sum_{k=1}^{Q_b} \theta_k^b \underbrace{b_k(\boldsymbol{\zeta}_j, \eta_l)}_{\text{precomputed}}, & (\mathbf{G}_N(\boldsymbol{\mu}))_l &= \sum_{k=1}^{Q_g} \theta_k^g \underbrace{g_k(\eta_l)}_{\text{precomputed}}, & 1 \leq l \leq N. \end{aligned} \quad (5.50)$$

These quantities are evaluated once during the offline procedure and then stored for online phases. During the online phase, the assembly of the linear system (5.49) only requires the additions of matrices (resp. vectors) of size  $\mathcal{O}(N^2)$  (resp.  $\mathcal{O}(N)$ ).

## 2.2 Error Estimation

Unfortunately, we did not have enough time to implement error estimators for the saddle-point problems: the computation of the dual norm of the residual is working, but we had no time to adapt the SCM algorithm for saddle-point problems. In the further numerical results, we will use the dual norm of the residual as an error indicator. Nevertheless, we remind in this section the principal error estimators found in the literature [21, 26, 80, 87].

We assume that the RB spaces  $V_N$  and  $Q_N$  are built and so we are able to evaluate a reduced basis solution  $(\mathbf{u}_N(\boldsymbol{\mu}), p_N(\boldsymbol{\mu}))$  for any  $\boldsymbol{\mu} \in \mathcal{D}$ . We are also interested in giving a sharp and efficient bound to the different errors, with respect to the truth model,

$$e_N^u(\boldsymbol{\mu}) = \mathbf{u}_h(\boldsymbol{\mu}) - \mathbf{u}_N(\boldsymbol{\mu}) \in V, \quad (5.51)$$

$$e_N^p(\boldsymbol{\mu}) = p_h(\boldsymbol{\mu}) - p_N(\boldsymbol{\mu}) \in Q, \quad (5.52)$$

$$e_N(\boldsymbol{\mu}) = (e_N^u(\boldsymbol{\mu}), e_N^p(\boldsymbol{\mu})) \in Z. \quad (5.53)$$

We also may be interested in the evaluation of the error for a linear output. We assume here that the output  $s(\boldsymbol{\mu})$  can be express as

$$s(\boldsymbol{\mu}) = l((\mathbf{u}(\boldsymbol{\mu}), p(\boldsymbol{\mu})); \boldsymbol{\mu}) = l^u(\mathbf{u}(\boldsymbol{\mu}); \boldsymbol{\mu}) + l^p(p(\boldsymbol{\mu}); \boldsymbol{\mu}) \quad (5.54)$$

where  $l : Z \rightarrow \mathbb{R}$ ,  $l^u : V \rightarrow \mathbb{R}$  and  $l^p : Q \rightarrow \mathbb{R}$  are linear forms which admit affine decompositions. In the particular case where  $l^u = f$  and  $l^p = g$  the problem is compliant. We denote by  $e_N^s(\boldsymbol{\mu})$  the error on the output  $s(\boldsymbol{\mu})$

$$e_N^s(\boldsymbol{\mu}) = s_h(\boldsymbol{\mu}) - s_N(\boldsymbol{\mu}). \quad (5.55)$$

For the construction of the different error bounds, we may need the following ingredients:

- Online-efficient upper and lower bounds for the continuity and coercivity constants introduced in (5.34)-(5.36),

$$\gamma_a^{LB}(\boldsymbol{\mu}) \leq \gamma_a(\boldsymbol{\mu}) \leq \gamma_a^{UB}(\boldsymbol{\mu}), \quad (5.56)$$

$$\alpha_a^{LB}(\boldsymbol{\mu}) \leq \alpha_a(\boldsymbol{\mu}) \leq \alpha_a^{UB}(\boldsymbol{\mu}). \quad (5.57)$$

Those bounds can be computed using the SCM, as mentioned in section 1.2.

- Online-efficient upper and lower bounds for the inf-sup constant, either the Brezzi (5.37) or the Babuška (5.38) version,

$$\beta_{Br}^{LB}(\boldsymbol{\mu}) \leq \beta_{Br}(\boldsymbol{\mu}) \leq \beta_{Br}^{UB}(\boldsymbol{\mu}), \quad (5.58)$$

$$\beta_{Ba}^{LB}(\boldsymbol{\mu}) \leq \beta_{Ba}(\boldsymbol{\mu}) \leq \beta_{Ba}^{UB}(\boldsymbol{\mu}). \quad (5.59)$$

These bounds can also be computed using the SCM.

- The dual norms of the residuals associated with the RB approximation

$$\|r_N^u(\cdot; \boldsymbol{\mu})\|_{V'} = \sup_{\mathbf{v} \in V} \frac{r_N^u(\mathbf{v}; \boldsymbol{\mu})}{\|\mathbf{v}\|_V}, \quad (5.60)$$

$$\|r_N^p(\cdot; \boldsymbol{\mu})\|_{Q'} = \sup_{q \in Q} \frac{r_N^p(q; \boldsymbol{\mu})}{\|q\|_Q}, \quad (5.61)$$

where,  $\forall \boldsymbol{\mu} \in \mathcal{D}$ ,

$$r_N^u(\mathbf{v}; \boldsymbol{\mu}) = f(\mathbf{v}; \boldsymbol{\mu}) - a(\mathbf{u}_N(\boldsymbol{\mu}), \mathbf{v}; \boldsymbol{\mu}) - b(\mathbf{v}, p_N(\boldsymbol{\mu}); \boldsymbol{\mu}), \quad (5.62)$$

$$r_N^p(q; \boldsymbol{\mu}) = g(q; \boldsymbol{\mu}) - b(\mathbf{u}_N(\boldsymbol{\mu}), q; \boldsymbol{\mu}). \quad (5.63)$$

We also define the total residual  $r_N((\mathbf{v}, q); \boldsymbol{\mu})$  as

$$r_N((\mathbf{v}, q); \boldsymbol{\mu}) = r_N^u(\mathbf{v}; \boldsymbol{\mu}) + r_N^p(q; \boldsymbol{\mu}), \quad \forall \mathbf{v} \in V, \forall q \in Q, \quad (5.64)$$

and its dual norm as

$$\|r_N(\cdot; \boldsymbol{\mu})\|_{Z'} = \sup_{(\mathbf{v}, q) \in Z} \frac{r_N((\mathbf{v}, q); \boldsymbol{\mu})}{\|(\mathbf{v}, q)\|_Z} = \sqrt{\|r_N^u(\cdot; \boldsymbol{\mu})\|_{V'}^2 + \|r_N^p(\cdot; \boldsymbol{\mu})\|_{Q'}^2}. \quad (5.65)$$

The efficient evaluation of this dual norms through an offline-online decomposition is detailed in appendix 3.

## Error Bounds for the Solutions

A first error bound  $\Delta_N(\boldsymbol{\mu})$  for Saddle Point problems, and especially Stokes equations, uses the Babuška inf-sup constant in order to bound the global error of the solutions  $e_N(\boldsymbol{\mu})$ , with respect to the truth model,

$$\|e_N(\boldsymbol{\mu})\|_Z \leq \Delta_N(\boldsymbol{\mu}) := \frac{\|r_N(\cdot; \boldsymbol{\mu})\|_{Z'}}{\beta_{Ba}^{LB}(\boldsymbol{\mu})}, \quad \forall \boldsymbol{\mu} \in \mathcal{D}. \quad (5.66)$$

This error bound is a direct extension of the error estimation for coercive problems.

Another bound is proposed in [35, 34], in order to measure separately the errors on the velocity field  $e_N^u(\boldsymbol{\mu})$  and the pressure field  $e_N^p(\boldsymbol{\mu})$ . This formulation requires more elements and uses the Brezzi inf-sup constant to define two bounds  $\Delta_N^u(\boldsymbol{\mu})$  and  $\Delta_N^p(\boldsymbol{\mu})$ , for any  $\boldsymbol{\mu} \in \mathcal{D}$ ,

$$\|e_N^u(\boldsymbol{\mu})\|_V \leq \Delta_N^u(\boldsymbol{\mu}) := \frac{\|r_N^u(\cdot; \boldsymbol{\mu})\|_{V'}}{\alpha_a^{LB}(\boldsymbol{\mu})} + \left(1 + \frac{\gamma_a^{UB}(\boldsymbol{\mu})}{\alpha_a^{LB}(\boldsymbol{\mu})}\right) \frac{\|r^p(\cdot; \boldsymbol{\mu})\|_{Q'}}{\beta_{Br}^{LB}(\boldsymbol{\mu})}, \quad (5.67)$$

$$\|e_N^p(\boldsymbol{\mu})\|_{Q'} \leq \Delta_N^p(\boldsymbol{\mu}) := \frac{\|r_N^u(\cdot; \boldsymbol{\mu})\|_{V'}}{\beta_{Br}^{LB}(\boldsymbol{\mu})} + \frac{\gamma_a^{UB}(\boldsymbol{\mu})}{\beta_{Br}^{LB}(\boldsymbol{\mu})} \Delta_N^u(\boldsymbol{\mu}). \quad (5.68)$$

From these two bounds, we can then define a new global error estimator for the solution  $(\mathbf{u}_N(\boldsymbol{\mu}), p_N(\boldsymbol{\mu}))$  as

$$\|e_N(\boldsymbol{\mu})\|_Z \leq \Delta_N^{Br}(\boldsymbol{\mu}) := \sqrt{(\Delta_N^u(\boldsymbol{\mu}))^2 + (\Delta_N^p(\boldsymbol{\mu}))^2}, \quad \forall \boldsymbol{\mu} \in \mathcal{D}. \quad (5.69)$$

For a sharpness study and a comparison of these different error bounds, we refer the readers to the cited references.

## Error Bound for Linear Output

In the convenient case of compliant problems, as in the coercive case, we can directly evaluate the error bound on the error  $e_N^s(\boldsymbol{\mu})$  [87],

$$|e_N^s(\boldsymbol{\mu})| \leq \Delta_N^{s,c}(\boldsymbol{\mu}) := 2 \frac{\|r_n(\cdot; \boldsymbol{\mu})\|_{Z'}}{\beta_{Ba}^{LB}(\boldsymbol{\mu})}, \quad (5.70)$$

this formulation is again a direct generalization of the coercive framework.

In the more general case of non-compliant problems, we introduce a dual problem, as find  $(\boldsymbol{\psi}(\boldsymbol{\mu}), \lambda(\boldsymbol{\mu})) \in V \times Q$  such that

$$\begin{cases} a(\mathbf{v}, \boldsymbol{\psi}(\boldsymbol{\mu}); \boldsymbol{\mu}) + b(\lambda(\boldsymbol{\mu}), \mathbf{v}; \boldsymbol{\mu}) = -l^u(\mathbf{v}; \boldsymbol{\mu}) & \forall \mathbf{v} \in V \\ b(q, \boldsymbol{\psi}(\boldsymbol{\mu}); \boldsymbol{\mu}) = -l^p(q; \boldsymbol{\mu}), & \forall q \in Q \end{cases} \quad (5.71)$$

where  $l^u$  and  $l^p$  are defined in (5.54). We denote by  $(\boldsymbol{\psi}_h(\boldsymbol{\mu}), \lambda_h(\boldsymbol{\mu}))$  the discrete solution of system (5.71) obtained using our truth model, with a discretization and a space definition adapted in order to satisfy the inf-sup condition. We denote by  $V_n^{du}$  and  $Q_N^{du}$  the dual RB spaces. We assume that these RB spaces are built during the greedy algorithm in a way similar to the coercive case. We can then define the dual residuals,

$$\begin{aligned} r_N^{u,du}(\mathbf{v}; \boldsymbol{\mu}) &= l^u(\mathbf{v}; \boldsymbol{\mu}) + a(\mathbf{v}, \boldsymbol{\psi}_N(\boldsymbol{\mu}); \boldsymbol{\mu}) + b(\lambda_N(\boldsymbol{\mu}), \mathbf{v}; \boldsymbol{\mu}), \quad \forall \mathbf{v} \in V, \\ r_N^{p,du}(q; \boldsymbol{\mu}) &= l^p(q; \boldsymbol{\mu}) + b(q, \boldsymbol{\psi}_N(\boldsymbol{\mu}); \boldsymbol{\mu}), \quad \forall q \in Q, \\ r_N^{du}(\mathbf{v}, q; \boldsymbol{\mu}) &= r_N^{u,du}(\mathbf{v}; \boldsymbol{\mu}) + r_N^{p,du}(q; \boldsymbol{\mu}), \end{aligned} \quad (5.72)$$

and the associated dual norms

$$\begin{aligned}
\|r_N^{\mathbf{u},du}(\cdot; \boldsymbol{\mu})\|_{V'} &= \sup_{\mathbf{v} \in V} \frac{r_N^{\mathbf{u},du}(\mathbf{v}; \boldsymbol{\mu})}{\|\mathbf{v}\|_V} \\
\|r_N^{p,du}(\cdot; \boldsymbol{\mu})\|_{Q'} &= \sup_{q \in Q} \frac{r_N^{p,du}(q; \boldsymbol{\mu})}{\|q\|_Q} \\
\|r_N^{du}(\cdot; \boldsymbol{\mu})\|_{Z'} &= \sup_{(\mathbf{v}, q) \in Z} \frac{r_N^{du}((\mathbf{v}, q); \boldsymbol{\mu})}{\|(\mathbf{v}, q)\|_Z} = \sqrt{\|r_N^{\mathbf{u},du}(\cdot; \boldsymbol{\mu})\|_{V'}^2 + \|r_N^{p,du}(\cdot; \boldsymbol{\mu})\|_{Q'}^2}.
\end{aligned} \tag{5.73}$$

The error bound for non-compliant linear output is then defined as

$$|e_N^s(\boldsymbol{\mu})| \leq \Delta_N^s(\boldsymbol{\mu}) := 2 \frac{\|r_N^{du}(\boldsymbol{\mu})\|_{Z'}^{1/2} \|r_N(\boldsymbol{\mu})\|_{Z'}^{1/2}}{\beta_{Ba}^{LB}(\boldsymbol{\mu})} \tag{5.74}$$

The evaluation of the dual norms of the residuals associated with each RB space is detailed in appendix 3.2.

### 3 Numerical Applications : RB for Stokes Problem

To validate our implementation of the RB framework for Saddle Point problems, we present, in this section, our results on a test case proposed in [35].

#### 3.1 Problem Formulation

We consider a Stokes flow in a 2D channel with an obstacle, see figure 5.1. The channel is a rectangle with fixed dimension  $(4 \times 1)$ . The obstacle is rectangular with variable width  $(\mu_1)$  and height  $(\mu_2)$ . The parameter  $\boldsymbol{\mu} = (\mu_1, \mu_2)$  is bounded such that  $\boldsymbol{\mu} \in \mathcal{D} = [0.1, 0.5] \times [0.1, 0.5]$ . Thus the parametrized domain is given by  $\tilde{\Omega}(\boldsymbol{\mu}) = [0, 4] \times [0, 1] \setminus \tilde{\mathcal{O}}(\boldsymbol{\mu})$  where  $\tilde{\mathcal{O}}(\boldsymbol{\mu}) = ]2 - \mu_1/2, 2 + \mu_1/2[ \times ]0, \mu_2[$  is the obstacle. We denote by  $\tilde{\Gamma}(\boldsymbol{\mu})$  the boundary of  $\tilde{\Omega}(\boldsymbol{\mu})$ ,  $\tilde{\Gamma}_{in}$  the inflow boundary,  $\tilde{\Gamma}_{out}$  the outflow boundary and by  $\tilde{\Gamma}_0(\boldsymbol{\mu}) = \tilde{\Gamma}(\boldsymbol{\mu}) \setminus (\tilde{\Gamma}_{in} \cup \tilde{\Gamma}_{out})$  the boundary on where we assume a no-slip condition. Here the  $\tilde{\cdot}$  symbol refers to the parametrized domain in opposition to the reference domain  $\Omega$ .

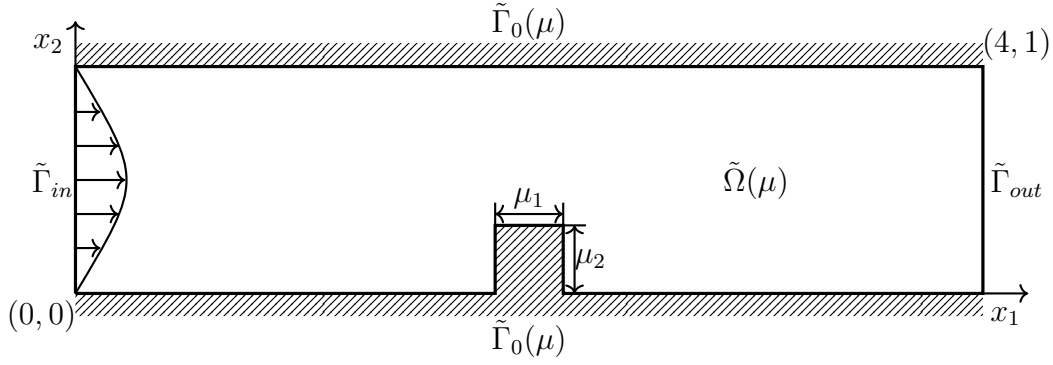


Figure 5.1 – 2D Geometry of the Problem

### Formulation on a Parametrized Domain $\tilde{\Omega}(\mu)$

The Stokes equations for the velocity  $\tilde{\mathbf{u}}(\mu) : \tilde{\Omega}(\mu) \rightarrow \mathbb{R}^2$  and the pressure  $\tilde{p}(\mu) : \tilde{\Omega}(\mu) \rightarrow \mathbb{R}$  are given by

$$\left\{ \begin{array}{l} \tilde{\nabla} \tilde{p}(\mu) - \nu \tilde{\Delta} \tilde{\mathbf{u}}(\mu) = 0, \quad \text{in } \tilde{\Omega}(\mu) \\ \tilde{\nabla} \cdot \tilde{\mathbf{u}}(\mu) = 0, \quad \text{in } \tilde{\Omega}(\mu), \\ \tilde{\mathbf{u}}(\mu) = 0, \quad \text{on } \tilde{\Gamma}_0(\mu), \\ (\nu \tilde{\nabla} \tilde{\mathbf{u}}(\mu) - \mathbb{I} \tilde{p}(\mu)) \cdot \tilde{\mathbf{n}} = 0, \quad \text{on } \tilde{\Gamma}_{out}, \\ \tilde{\mathbf{u}} = \boldsymbol{\phi}, \quad \text{on } \tilde{\Gamma}_{in}, \end{array} \right. \quad (5.75)$$

with  $\boldsymbol{\phi} = \begin{pmatrix} 4\tilde{x}_2(1 - \tilde{x}_2) \\ 0 \end{pmatrix}$  and  $\nu = 1$  the fluid viscosity.

We introduce the finite element space  $\tilde{Z}_h(\boldsymbol{\mu}) = \tilde{V}_h(\boldsymbol{\mu}) \times \tilde{Q}_h(\boldsymbol{\mu}) = \text{THH}_{ch}^k(\tilde{\Omega}(\boldsymbol{\mu}))$  and the usual associated scalar products and norms. The weak formulation on the parametrized domain  $\tilde{\Omega}(\boldsymbol{\mu})$  can be written as: find  $(\tilde{\mathbf{u}}_h, \tilde{p}_h) \in \tilde{Z}_h(\boldsymbol{\mu})$  such that

$$\tilde{a}(\tilde{\mathbf{u}}_h, \tilde{\mathbf{v}}_h) + \tilde{b}(\tilde{\mathbf{v}}_h, \tilde{p}_h) = \tilde{f}(\tilde{\mathbf{v}}_h), \quad \forall \tilde{\mathbf{v}}_h \in \tilde{V}_h(\boldsymbol{\mu}), \quad (5.76)$$

$$\tilde{b}(\tilde{\mathbf{u}}_h, \tilde{q}_h) = \tilde{g}(\tilde{q}_h), \quad \forall \tilde{q}_h \in \tilde{Q}_h(\boldsymbol{\mu}), \quad (5.77)$$

where

$$\tilde{a}(\mathbf{u}, \mathbf{v}) = \int_{\tilde{\Omega}(\boldsymbol{\mu})} \nu \tilde{\nabla} \mathbf{u} : \tilde{\nabla} \mathbf{v} \, d\tilde{\Omega}, \quad \tilde{f}(\mathbf{v}) = 0, \quad (5.78)$$

$$\tilde{b}(\tilde{\mathbf{u}}, \tilde{q}) = - \int_{\tilde{\Omega}(\boldsymbol{\mu})} q(\tilde{\nabla} \cdot \tilde{\mathbf{u}}) \, d\tilde{\Omega}, \quad \tilde{g}(\tilde{q}_h) = 0. \quad (5.79)$$

### Formulation on the Reference Domain $\Omega$

In order to apply the RBM to this problem, we perform the simulation on a reference domain  $\Omega = \tilde{\Omega}(\boldsymbol{\mu}_r)$ . We denote by  $Z_h = V_h \times Q_h = \text{THH}_{ch}^k(\Omega)$  the FE space on

the reference domain  $\Omega$ . We will use a mapping from  $\tilde{\Omega}(\boldsymbol{\mu})$  to  $\Omega$  in order to solve the problem (5.76) in the space  $Z_h$ .

We first introduce  $\forall \boldsymbol{\mu} \in \mathcal{D}$  a domain decomposition of  $\tilde{\Omega}(\boldsymbol{\mu})$ ,

$$\overline{\tilde{\Omega}(\boldsymbol{\mu})} = \bigcup_{k=1}^P \overline{\tilde{\Omega}^k(\boldsymbol{\mu})}, \quad \text{s.t.} \quad \tilde{\Omega}^k(\boldsymbol{\mu}) \cap \tilde{\Omega}^l(\boldsymbol{\mu}) = \emptyset, \quad \forall 1 \leq k < l \leq P. \quad (5.80)$$

This construction is also naturally verified for the reference domain  $\Omega$ ,

$$\overline{\Omega} = \bigcup_{k=1}^P \overline{\Omega^k}, \quad \text{s.t.} \quad \Omega^k \cap \Omega^l = \emptyset, \quad 1 \leq k < l \leq P. \quad (5.81)$$

We now define  $P$  continuous and bijective affine mappings  $\mathcal{G}^k(\cdot; \boldsymbol{\mu})$  such that  $\forall \boldsymbol{\mu} \in \mathcal{D}$ ,  $\mathcal{G}^k(\tilde{\Omega}^k(\boldsymbol{\mu})) = \Omega^k$ ,  $\forall 1 \leq k \leq P$  and  $\mathcal{G}^k(\tilde{\boldsymbol{x}}; \boldsymbol{\mu}) = \mathcal{G}^l(\tilde{\boldsymbol{x}}; \boldsymbol{\mu})$ ,  $\forall \tilde{\boldsymbol{x}} \in \overline{\tilde{\Omega}^k(\boldsymbol{\mu})} \cap \overline{\tilde{\Omega}^l(\boldsymbol{\mu})}$ ,  $\forall 1 \leq k < l \leq P$ ,

$$\mathcal{G}^k(\tilde{\boldsymbol{x}}; \boldsymbol{\mu}) = \mathbf{G}^k(\boldsymbol{\mu})\tilde{\boldsymbol{x}} + \mathbf{H}^k(\boldsymbol{\mu}), \quad 1 \leq k \leq P. \quad (5.82)$$

We propose such a decomposition for the domain  $\Omega = \tilde{\Omega}(\boldsymbol{\mu}_r)$  in figure 5.2. Note that the design of this decomposition will definitely change the affine decomposition, but not the number of terms (as long as the decomposition is adapted to the geometry and the parametrization).

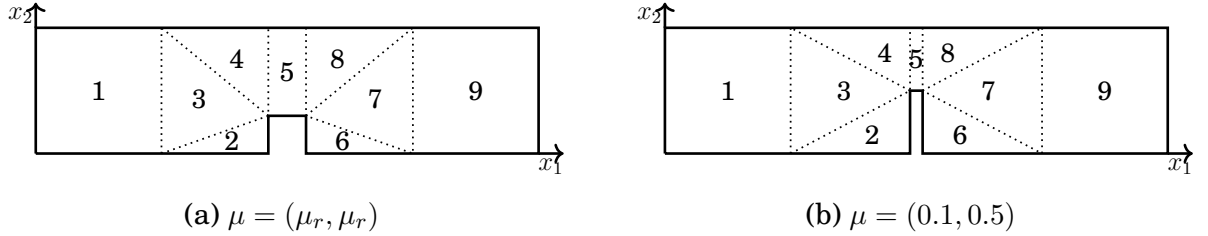


Figure 5.2 – Domain Decomposition used for the mapping

Using the usual composition rules for differentiation operators, we have the relations,

$$\begin{aligned} \partial/\partial x_i &= \mathbf{G}_{ij}^k \partial/\partial x_j \\ d\tilde{\Omega}(\boldsymbol{\mu}) &= J^k(\boldsymbol{\mu})d\Omega, \quad J^k(\boldsymbol{\mu}) = |\det(\mathbf{G}^k(\boldsymbol{\mu}))^{-1}|. \end{aligned} \quad (5.83)$$

We can then rewrite the weak formulation (5.76) on the reference domain  $\Omega$  as find  $(\mathbf{u}_h(\boldsymbol{\mu}), p_h(\boldsymbol{\mu})) \in Z_h$  such that

$$a(\mathbf{u}_h, \mathbf{v}_h) + b(\mathbf{v}_h, p_h) = f(\mathbf{v}_h), \quad \forall \mathbf{v}_h \in V_h, \quad (5.84)$$

$$b(\mathbf{u}_h, q_h) = g(q_h), \quad \forall q_h \in Q_h, \quad (5.85)$$

with

$$a(\mathbf{u}, \mathbf{v}; \boldsymbol{\mu}) = \sum_{k=1}^P G_{mj}^k(\boldsymbol{\mu}) G_{lj}^k J^k(\boldsymbol{\mu}) \int_{\Omega^k} \frac{\partial u_i}{\partial x_m} \frac{\partial v_i}{\partial x_l} d\Omega, \quad f(\mathbf{v}) = 0 \quad (5.86)$$

$$b(\mathbf{u}, q) = \sum_{k=1}^P -G_{ji}^k J^k(\boldsymbol{\mu}) \int_{\Omega^k} \frac{\partial u_i}{\partial x_j} q \, d\Omega, \quad g(q) = 0. \quad (5.87)$$

The detailed affine decomposition is presented in appendix 4.

### 3.2 Numerical Results

We use a FE space  $Z_h = V_h \times Q_h = \mathbb{TH}_{ch}^1$ . The discretization leads to  $\mathcal{N} = 66857$  velocity and pressure degrees of freedom. The mesh of the reference domain is presented in figure 5.3. We illustrate the results with some velocity and pressure profiles in figure 5.4. These solutions were computed on the reference domain and then warped onto the real domain for visualization.

The supersampling  $\Xi$  is an equidistributed grid of  $\mathcal{D}$  with 1000 points. We are interested in this case in the convergence of the errors

$$\begin{aligned} \|e_N^u(\boldsymbol{\mu})\|_2 &= \|\mathbf{u}_h(\boldsymbol{\mu}) - \mathbf{u}_N(\boldsymbol{\mu})\|_2, \\ \|e_N^p(\boldsymbol{\mu})\|_2 &= \|p_h(\boldsymbol{\mu}) - p_N(\boldsymbol{\mu})\|_2, \\ \|e_N(\boldsymbol{\mu})\|_2 &= \sqrt{\|e_N^u(\boldsymbol{\mu})\|_2^2 + \|e_N^p(\boldsymbol{\mu})\|_2^2}. \end{aligned} \quad (5.88)$$

We computed these quantities for 100 different values of  $\boldsymbol{\mu}$ . The figure 5.5 presents the evolution of the maximum, minimum and mean value of these errors, for  $N \in \llbracket 1, 50 \rrbracket$ . The reduced basis was built with 50 different parameters in  $S_N$ , using the dual norm of the residual as an error indicator in the greedy algorithm 4. We clearly observe the convergence of the method. We also notice an unexpected behavior with the error growing again near the value  $N = 20$ . We found no explanation for this “jump” for now, however these results validate our implementation of the method.

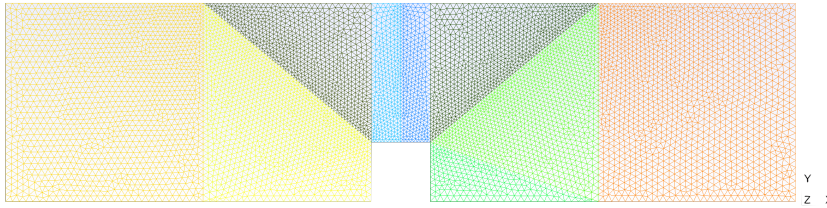


Figure 5.3 – Mesh Used for the FE simulation.

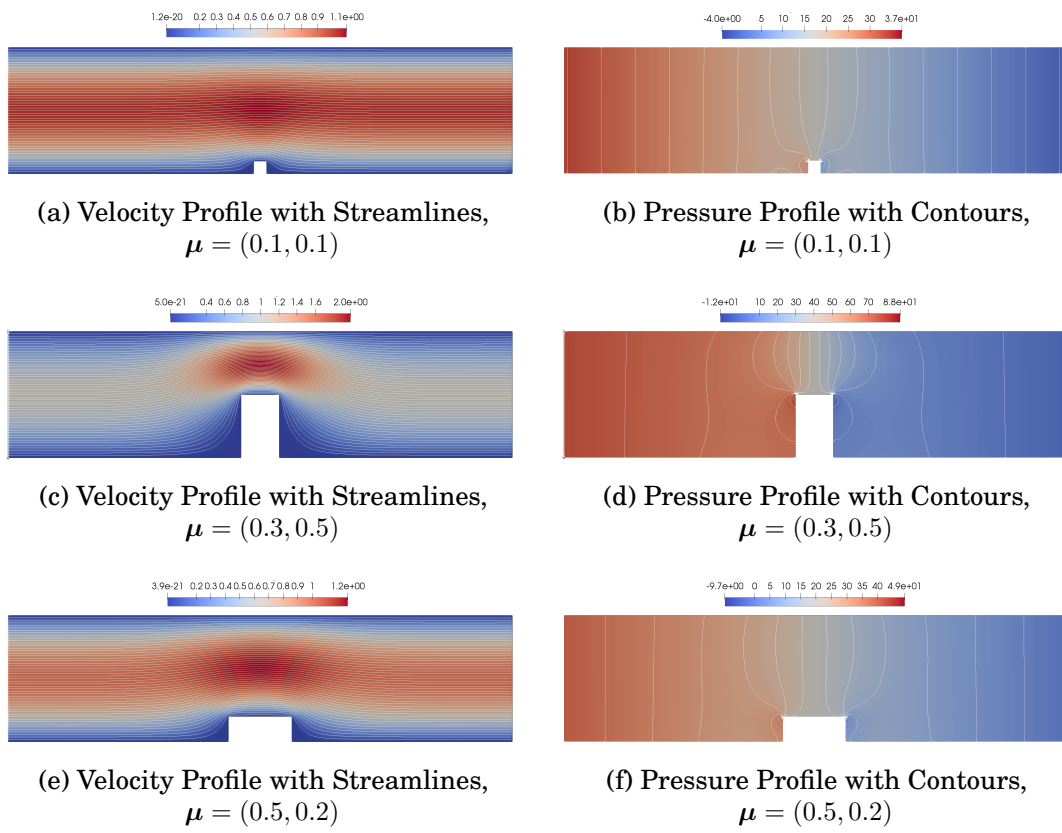


Figure 5.4 – Velocity and pressure profiles for 3 different parameters.



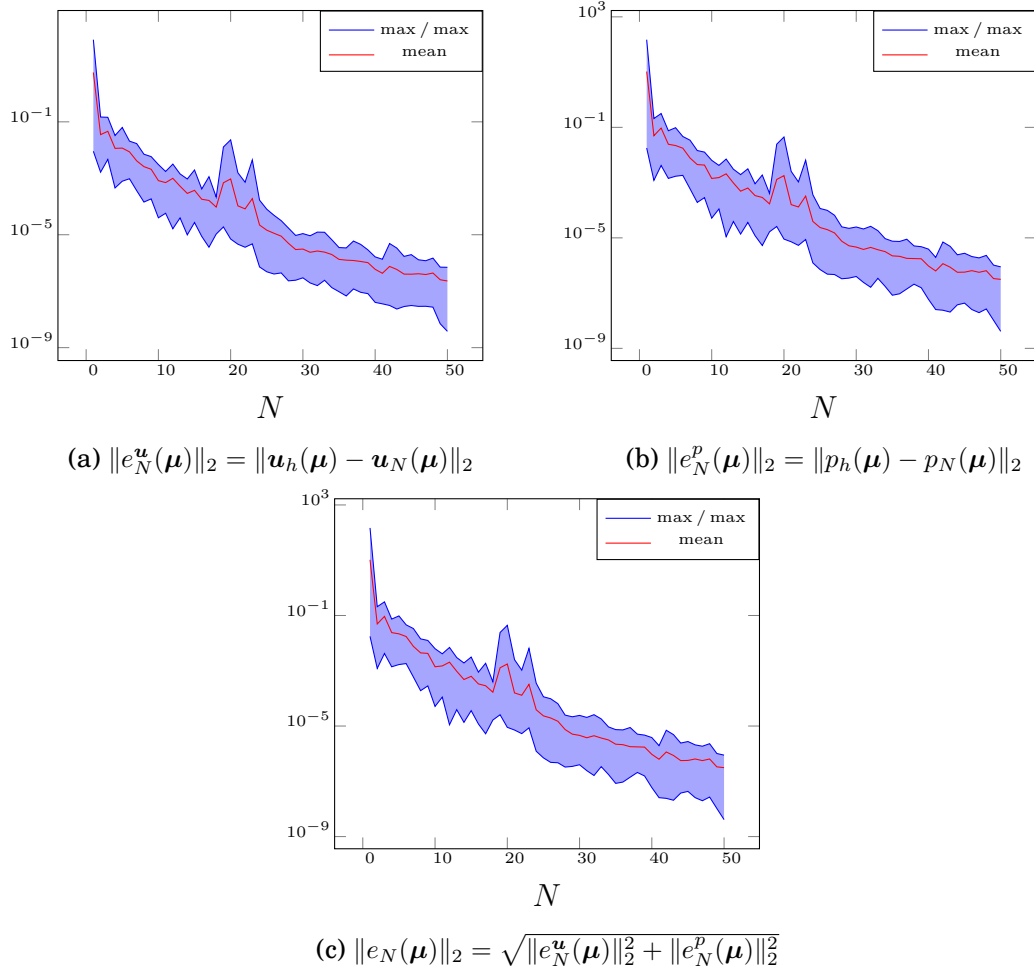


Figure 5.5 – Convergence of the errors with respect to the number of vectors in the reduced basis. Maximum, minimum and mean values evaluated from 100 approximations.

## Conclusion

In this chapter, we presented the main ingredients of the Reduced Basis Method. We started by introducing the method with the coercive problems. Then we detailed the algorithms required for the construction of a stable reduced space in the case of saddle-point problems. In the last section of this chapter, we tested the implementation of our RB saddle-point framework on a Stokes problem with geometric parameters.

The theoretical part of this chapter was essentially a reminder of techniques which are already well developed in the literature. However, it seemed important to rewrite these algorithms. Indeed, the RBM was already implemented in *Feel++* for coercive problems, but the framework for decoupled multi-physics problems (and saddle-point in particular) was entirely implemented during this thesis work. This implementation is detailed in part III.

With the validation of our framework for saddle-point problems, we can now consider the reduction of Stokes problems. It is a first step towards the reduction of Navier-Stokes and aerothermal problems. We now need tools to deal with the non-linearity. The next chapter will present solutions like the Empirical Interpolation Method which allows building an affine approximation of non-affine and non-linear operators.



---

## Non-Affine and Non-Linear Problems

---

### Contents

---

<b>1</b>	<b>The Empirical Interpolation Method . . . . .</b>	<b>114</b>
1.1	Reduced Basis for Non-Linear Problems . . . . .	115
<b>2</b>	<b>EIM for Discrete Operators . . . . .</b>	<b>117</b>
2.1	Online evaluation . . . . .	118
2.2	Operator on a Product Space . . . . .	120
<b>3</b>	<b>Simultaneous Empirical Interpolation and Reduced Basis Method for Non-linear Problems . . . . .</b>	<b>121</b>
<b>4</b>	<b>Preliminary Results . . . . .</b>	<b>122</b>
4.1	Problem Formulation . . . . .	124
4.2	Numerical Results . . . . .	124

---

We presented in the previous chapter the Reduced Basis Method for two different kinds of problems, all with the common assumptions: the affine decomposition and the linearity. In this chapter, we will present solutions to deal with complex problems which might be non-affine or non-linear.

One of the main advantages of the RBM is the offline/online procedure which allows very fast online  $\mathcal{N}$ -independent computations, once the expensive offline phase is done. However, it is only possible thanks to the affine decomposition of the operators as we detailed in the section 1.1. Many problems do not present an affine decomposition. For instance, we cannot recover an affine dependency for the

stabilization operators introduced in the chapter 3. Another limitation of the RBM as presented in the previous chapter is the delicate problem of non-linear operators. We need a solution to deal with the non-linearity. The iterative solvers introduced in the chapter 2 allows to solve the non-linearity. However, we want to keep an online resolution independent of the FE dimension  $\mathcal{N}$ , even for non-linear problems. The Empirical Interpolation Method (EIM) [3, 65, 23, 37] is an answer to both of these problems: the non-affine dependency and the non-linearity of the models. This method allows recovering an affine approximation which will be used within the RBM.

In the present chapter, we first remind the principle of this method. Then we propose a discrete version of the EIM which was implemented during this thesis. And finally, we present the Simultaneous EIM and RB (SER) construction which was also integrated into our discrete version of EIM. The last section is dedicated to a brief validation of our implementation.

## 1 The Empirical Interpolation Method

Assuming that a parametrized PDE is not affine in parameter  $\boldsymbol{\mu}$  mean that it depends at least on a non-affine - and potentially non-linear - function  $w(u(\boldsymbol{\mu}), \boldsymbol{x}, \boldsymbol{\mu}) : X \times \mathbb{R}^d \times \mathcal{D} \rightarrow \mathbb{R}$ . In this section we present how the Empirical Interpolation Method allows to recover an affine approximation  $w_M$  of this function such that

$$w(u(\boldsymbol{\mu}), \boldsymbol{x}, \boldsymbol{\mu}) \approx w_M(u(\boldsymbol{\mu}), \boldsymbol{x}, \boldsymbol{\mu}) = \sum_{m=1}^M \beta_m(u(\boldsymbol{\mu}), \boldsymbol{\mu}) q_m(\boldsymbol{x}). \quad (6.1)$$

The construction of this approximation requires an expensive offline phase which essentially consists in the selection of a set of parameters  $\bar{S}_M = \bar{\boldsymbol{\mu}}_m, 1 \leq m \leq M$ . From this set of parameters will be built a function space  $\bar{W}_M$  and a set of interpolation point  $\bar{T}_M$ . We detail this greedy algorithm in the next paragraphs.

### Offline Construction : Greedy Algorithm

Let  $\bar{\Xi} \subset \mathcal{D}$  a sample of the parameter space  $\mathcal{D}$ . The algorithm starts with an arbitrary choice of  $\bar{\boldsymbol{\mu}}_1 \in \bar{\Xi}$ , assuming that  $\bar{\xi}_1 = w(u(\bar{\boldsymbol{\mu}}_1), \cdot, \bar{\boldsymbol{\mu}}_1) \neq 0$ . The first interpolation point  $t_1$  is chosen as the maximum of the function  $\bar{\xi}_1$  over the domain  $\Omega$ . The first basis function  $q_1$  is then a normalization of the snapshot  $\bar{\xi}_1$ ,

$$t_1 = \arg \max_{\boldsymbol{x} \in \Omega} |\bar{\xi}_1(\boldsymbol{x})|, \quad q_1 = \frac{\bar{\xi}_1}{\bar{\xi}_1(t_1)}. \quad (6.2)$$

And then we set

$$\bar{S}_1 = \{\bar{\boldsymbol{\mu}}_1\}, \quad \bar{W}_1 = \text{span}\{\bar{\xi}_1\} = \text{span}\{q_1\}, \quad T_1 = \{t_1\} \quad (6.3)$$

Once this initialization is done, for any  $M \geq 2$ , the next parameters  $\bar{\boldsymbol{\mu}}_M$  are chosen to maximize the EIM approximation error,

$$\bar{\boldsymbol{\mu}}_M = \arg \max_{\bar{\boldsymbol{\mu}} \in \bar{\Xi}} \|w(u(\bar{\boldsymbol{\mu}}), \cdot, \bar{\boldsymbol{\mu}}) - w_{M-1}(u(\bar{\boldsymbol{\mu}}), \cdot, \bar{\boldsymbol{\mu}})\|_{\infty}. \quad (6.4)$$

From this new parameter  $\bar{\boldsymbol{\mu}}_M$  we define the residual, between the exact function  $w$  and its affine approximation  $w_M$  of rank  $M$ ,

$$r_M(\boldsymbol{x}) = w(u(\bar{\boldsymbol{\mu}}_M), \boldsymbol{x}, \bar{\boldsymbol{\mu}}_M) - w_M(u(\bar{\boldsymbol{\mu}}_M), \boldsymbol{x}, \bar{\boldsymbol{\mu}}_M). \quad (6.5)$$

The next interpolation point  $t_M$  is chosen to maximize this residual. The next basis function  $q_M$  is the normalization of the residual.

$$t_M = \arg \sup_{\boldsymbol{x} \in \Omega} |r_M(\boldsymbol{x})|, \quad q_M = \frac{r_M(\boldsymbol{x})}{|r_M(t_M)|}. \quad (6.6)$$

And we can set

$$\bar{S}_M = \bar{S}_{M-1} \cup \{\bar{\boldsymbol{\mu}}_M\}, \quad \bar{W}_M = \bar{W}_{M-1} \oplus \text{span}\{q_M\}, \quad T_M = T_{M-1} \cup \{t_M\}. \quad (6.7)$$

### Online : Evaluation of the $\beta_m(u, \boldsymbol{\mu})$

For a new parameter  $\boldsymbol{\mu}$ , and a function  $u \in X$ , the coefficients  $\beta_m(u, \boldsymbol{\mu})$  are evaluated in order to ensure the exactness of the affine approximation  $w_M$  at the  $M$  interpolation points  $\{t_m\}_{m=1}^M$ . In practice, the computation is made through the resolution of a  $M \times M$  system,

$$\boldsymbol{A}_M \boldsymbol{\beta}_M(u, \boldsymbol{\mu}) = \boldsymbol{F}_M(u, \boldsymbol{\mu}), \quad (6.8)$$

where  $(\boldsymbol{A}_M)_{ij} = q_j(t_i)$ ,  $(\boldsymbol{\beta}_M(u, \boldsymbol{\mu}))_i = \beta_i(u, \boldsymbol{\mu})$  and  $(\boldsymbol{F}_M(u, \boldsymbol{\mu}))_i = w(u, t_i, \boldsymbol{\mu})$ . Once this resolution is made, we can build the affine approximation

$$w(u, \boldsymbol{x}, \boldsymbol{\mu}) \approx w_M(u, \boldsymbol{x}, \boldsymbol{\mu}) = \sum_{m=1}^M \beta_m(u, \boldsymbol{\mu}) q_m(\boldsymbol{x}). \quad (6.9)$$

The considered matrix  $\boldsymbol{A}_M$  is independent of the parameter  $\boldsymbol{\mu}$  and can be assembled and stored during the offline phase. But the right hand side  $\boldsymbol{F}_M(\boldsymbol{\mu})$  has to be rebuilt for each new parameter  $\boldsymbol{\mu}$ , by evaluating the function  $w$  on the interpolation points.

## 1.1 Reduced Basis for Non-Linear Problems

We detailed above how to recover an affine approximation to efficiently use the RBM with non-affine problems. We are now interested in using the EIM with an iterative solver for non-linear problems.

Let consider a non-linear parametrized problem. We assume that we have an iterative algorithm - Picard or Newton for instance - in which the solution at the  $k$ -th iteration is denoted by  $u^k(\boldsymbol{\mu})$ . We suppose that a step of the considered iterative method can be express as find  $u^{k+1}(\boldsymbol{\mu})$  such that

$$a(u^{k+1}(\boldsymbol{\mu}), v; u^k(\boldsymbol{\mu}); \boldsymbol{\mu}) = f(v; u^k(\boldsymbol{\mu}); \boldsymbol{\mu}), \quad \forall v \in X, \quad (6.10)$$

where  $X$  is a suitable function space on domain  $\Omega$ .  $a(\cdot, \cdot; u; \boldsymbol{\mu}) : X \times X \rightarrow \mathbb{R}$  and  $l(\cdot; u; \boldsymbol{\mu}) : X \rightarrow \mathbb{R}$  are continuous bilinear and linear forms, for any parameter  $\boldsymbol{\mu} \in \mathcal{D}$  and any  $u \in X$ .

The RBM will not be applied directly on the parametrized problem (6.10), which is obviously not affine in parameters, but on an affine approximation given by the EIM, written in a generic formulation as

$$\sum_{q=1}^{Q_a^M} \theta_q^a(u^k(\boldsymbol{\mu}); \boldsymbol{\mu}) a_q(u^{k+1}, v) = \sum_{q=1}^{Q_f^M} \theta_q^f(u^k(\boldsymbol{\mu}); \boldsymbol{\mu}) f_q(v), \quad \forall v \in X \quad (6.11)$$

This affine approximation is obtained using once or several times the EIM algorithm, depending on the number of non-affine components of the original problem. Note that some of the coefficients  $\theta_q^a$  and  $\theta_q^f$  might not depend on the iterates  $u^k$  and remain linear. Some of these coefficients might also come from an affine part of the original problem, their evaluation is then analytical and does not require the use of EIM.

Once we have an affine decomposition thanks to the EIM, we can perform a usual RB offline/online procedure. Offline, the resolution is made using our favorite iterative solver until convergence,  $u^k(\boldsymbol{\mu}) \rightarrow u(\boldsymbol{\mu})$ . This solution -snapshot- will be used to build the RB space  $X_N$ , as detailed in chapter 5.

The online resolution will require some non-linear iterations as well, which can be written, using the same affine approximation: find  $u_N^{k+1}(\boldsymbol{\mu})$  such that

$$\sum_{q=1}^{Q_a^M} \theta_q^a(u_N^k(\boldsymbol{\mu}); \boldsymbol{\mu}) a_q(u_N^{k+1}, v) = \sum_{q=1}^{Q_f^M} \theta_q^f(u_N^k(\boldsymbol{\mu}); \boldsymbol{\mu}) f_q(v), \quad \forall v \in X_N. \quad (6.12)$$

We rewrite this equation under a matricial formulation after the projection on the RB space  $W_N = \text{span}\{\xi_n, 1 \leq n \leq N\}$ ,

$$\sum_{j=1}^N \sum_{q=1}^{Q_a^M} \sum_{n=1}^N \theta_q^a(u_N^k(\boldsymbol{\mu}); \boldsymbol{\mu}) (\mathbf{A}_N^q)_{ij} u_N^{k+1} = \sum_{q=1}^{Q_f^M} \theta_q^f(u_N^k(\boldsymbol{\mu}); \boldsymbol{\mu}) (\mathbf{F}_N^q)_i, \quad 1 \leq i \leq N \quad (6.13)$$

where,  $\forall 1 \leq i, j \leq N$ ,

$$(\mathbf{A}_N^q)_{ij} = a_q(\xi_j, \xi_i) \quad (\mathbf{F}_N^q)_j = f_q(\xi_j) \quad \forall 1 \leq q \leq Q_a^M, \quad \forall 1 \leq q' \leq Q_f^M \quad (6.14)$$

The matrices  $\mathbf{A}_N^q$  and the vectors  $\mathbf{F}_N^q$  are precomputed once during the online phase. The online resolution then consists in several resolutions of the  $N \times N$  system

$$\left[ \sum_{q=1}^{Q_a^M} \theta_q^a(u_N^k(\boldsymbol{\mu}); \boldsymbol{\mu}) \mathbf{A}_N^q \right] u_N^{k+1}(\boldsymbol{\mu}) = \left[ \sum_{q=1}^{Q_f^M} \theta_q^f(u_N^k(\boldsymbol{\mu}); \boldsymbol{\mu}) \mathbf{F}_N^q \right]. \quad (6.15)$$

The reduced solution  $u_N(\boldsymbol{\mu})$  is obtained after convergence of the non-linear solver,  $u_N^k(\boldsymbol{\mu}) \rightarrow u_N(\boldsymbol{\mu})$ .

## 2 EIM for Discrete Operators

We saw in the previous section that the Empirical Interpolation Method is a required enabler for RBM applied to non-affine or non-linear problems. In the Feel++ framework, this method was already implemented, see [105], but it required an analytical expression of the non-affine component  $w$ . Also it was not adapted for the stabilization operators presented in chapter 3. To tackle this problem, we implemented a discrete version of the algorithm, known in the literature as multi-component EIM or Matrix EIM [97]. This method is also called Discrete Empirical Interpolation Method (DEIM) [16], when a Proper Orthogonal Decomposition (POD) is used to select the EIM vectors. The method and the global algorithm are identical to the one described in previous section but the formalism and the implementation, detailed in chapter 10, are different, and then we chose to dedicate a section to this particular version of the EIM.

The algorithm aims at building an affine approximation of a parametrized tensor  $\mathbf{T}(u(\boldsymbol{\mu}), \boldsymbol{\mu})$ . The algorithm as presented in this section is independent of the order of the tensor,  $r = 1, 2, \dots$  and can be used either for discrete linear or bilinear operators. The affine approximation of tensor  $\mathbf{T}(u(\boldsymbol{\mu}), \boldsymbol{\mu})$  can be express as

$$\mathbf{T}(u(\boldsymbol{\mu}), \boldsymbol{\mu}) \approx \mathbf{T}_M(u(\boldsymbol{\mu}), \boldsymbol{\mu}) = \sum_{m=0}^M \beta_m(\boldsymbol{\mu}) \boldsymbol{\Phi}^m \quad (6.16)$$

where  $\beta_m(\boldsymbol{\mu})$  are scalar functions,  $\beta_m : \mathcal{D} \rightarrow \mathbb{R}$ , and  $\boldsymbol{\Phi}^m$  are order  $r$  tensors independent of  $\boldsymbol{\mu}$ . The affine approximation  $\mathbf{T}_M$  lives in a space  $\bar{W}_M = \text{span}\{\boldsymbol{\Phi}^m\}$  of low dimension  $M \ll \mathcal{N}$ . The space  $\bar{W}_M$  is built during an offline phase, using a greedy construction.

$\bar{W}_1$  is built by choosing an arbitrary first parameter  $\bar{\boldsymbol{\mu}}_1$ . For this parameter we compute  $\mathbf{T}(u(\bar{\boldsymbol{\mu}}_1), \bar{\boldsymbol{\mu}}_1)$  and look for the index  $i_1$  of the maximum entry of the tensor  $\mathbf{T}(u(\bar{\boldsymbol{\mu}}_1), \bar{\boldsymbol{\mu}}_1)$

$$i_1 = \arg \max_{i \in \mathbb{I}} |\mathbf{T}(u(\bar{\boldsymbol{\mu}}_1), \bar{\boldsymbol{\mu}}_1)_i|. \quad (6.17)$$



We can then initialize the greedy algorithm with

$$\Phi_1 = \frac{1}{|\mathbf{T}(u(\bar{\boldsymbol{\mu}}_1), \bar{\boldsymbol{\mu}}_1)_{i_1}|} \mathbf{T}(u(\bar{\boldsymbol{\mu}}_1), \bar{\boldsymbol{\mu}}_1), \quad \bar{W}_1 = \text{span}\{\Phi_1\}, \quad \mathcal{I}_1 = \{i_1\}, \quad \bar{S}_1 = \{\bar{\boldsymbol{\mu}}_1\}. \quad (6.18)$$

**Remark 8.** *The nature of the indices  $i$  depends on the order  $r$  of the tensor. An index can be defined, in a general way as  $i \in \mathbb{I} = \llbracket 1 : \mathcal{N} \rrbracket^r$ .*

We assume now that  $\bar{W}_{M-1}$  is built. The next parameters  $\bar{\boldsymbol{\mu}}_M$  is chosen to maximize the norm of the residual,  $\mathbf{R}_{M-1}(\boldsymbol{\mu})$ ,

$$\mathbf{R}_{M-1}(\boldsymbol{\mu}) = \mathbf{T}(u(\boldsymbol{\mu}), \boldsymbol{\mu}) - \mathbf{T}_{M-1}(u(\boldsymbol{\mu}), \boldsymbol{\mu}), \quad \bar{\boldsymbol{\mu}}_M = \arg \max_{\boldsymbol{\mu} \in \bar{\Xi}} \|\mathbf{R}_{M-1}(\boldsymbol{\mu})\|_\infty. \quad (6.19)$$

Then we can select the next interpolation index  $i_M$  as the maximum entry of the tensor  $\mathbf{R}_{M-1}(\bar{\boldsymbol{\mu}}_M)$

$$i_M = \arg \max_{i \in \mathbb{I}} |(\mathbf{R}_{M-1}(\bar{\boldsymbol{\mu}}_M))_i|, \quad (6.20)$$

and we set

$$\begin{aligned} \Phi_M &= \frac{1}{|(\mathbf{R}_{M-1}(\bar{\boldsymbol{\mu}}_M))_{i_M}|} \mathbf{R}_{M-1}(\bar{\boldsymbol{\mu}}_M), & \bar{W}_M &= \bar{W}_{M-1} \oplus \text{span}\{\Phi_M\}, \\ \mathcal{I}_M &= \mathcal{I}_{M-1} \cup \{i_M\}, & \bar{S}_M &= \bar{S}_{M-1} \cup \{\bar{\boldsymbol{\mu}}_M\}. \end{aligned} \quad (6.21)$$

The complete offline methodology is summarized in algorithm 5.

**Remark 9.** *We focused in this chapter on the treatment of the non-linearity of the problems. However, the EIM proposed in this section is also very convenient to automatically provide the affine decomposition of complex affine operators. For instance, it is widespread to have huge and complex affine decomposition when we are working with geometric parameters. EIM allows to recover this decomposition and to avoid the painful calculus session. We successfully tested this use on the Stokes test case 3.*

## 2.1 Online evaluation

The coefficients  $\beta_m(u, \boldsymbol{\mu})$  are defined to ensure the exactness of the interpolation in entries corresponding to the indices  $i_m \in \mathcal{I}_M$ . Then we can evaluate these coefficients by solving the  $M \times M$  system

$$\mathbf{A}_M \boldsymbol{\beta}_M(u, \boldsymbol{\mu}) = \mathbf{F}_M(u, \boldsymbol{\mu}), \quad (6.22)$$

where  $(\mathbf{A}_M)_{lj} = (\Phi_j)_{i_l}$ ,  $(\boldsymbol{\beta}_M(u, \boldsymbol{\mu}))_l = \beta_l(u, \boldsymbol{\mu})$  and  $(\mathbf{F}_M(u, \boldsymbol{\mu}))_l = (\mathbf{T}(u, \boldsymbol{\mu}))_{i_l}$ .

The matrix  $\mathbf{A}_M$  is built during the offline phase of the algorithm since it does not depend on the parameter  $\boldsymbol{\mu}$ . However, the assembly of the vector  $\mathbf{F}_M(\boldsymbol{\mu})$  is

**Algorithm 5:** EIM offline algorithm for discrete operators

---

```

1 Choose  $\bar{\Xi} \subset \mathcal{D}$ ; // Super Samplings
2 Choose  $M_{\max} \in \mathbb{N}$ ; // Max size for  $\bar{W}_M$ 
3 Choose  $\delta_{EIM}$ ; // Algorithm tolerances
4 Choose  $\bar{\mu}_1 \in \bar{\Xi}$ ; // Initialization of EIM construction
5 Set  $i_M \leftarrow \arg \max_{i \in \mathbb{I}} |\mathbf{T}(u_h(\bar{\mu}_1), \bar{\mu}_1)_i|$ ,  $\Phi_M \leftarrow \frac{\mathbf{T}(u_h(\bar{\mu}_1), \bar{\mu}_1)}{|(\mathbf{T}(u_h(\bar{\mu}_1), \bar{\mu}_1))_{i_1}|}$ ;
6 Set  $M \leftarrow 1$ ,  $\bar{S}_1 \leftarrow \{\bar{\mu}_1\}$ ,  $W_1 \leftarrow \{\Phi_1\}$ ,  $\mathcal{I}_1 \leftarrow \{i_1\}$ ;
7 repeat
8    $\bar{\mu}_{M+1} \leftarrow \arg \max_{\mu \in \bar{\Xi}} \underbrace{\|\mathbf{T}(u_N(\mu), \mu) - \mathbf{T}_M(u_N(\mu), \mu)\|}_{\mathbf{R}_M(\mu)}_\infty$ ;
9   if  $\|\mathbf{R}_M(\bar{\mu}_{M+1})\|_\infty \leq \delta_{EIM}$  then
10     break;
11   else
12      $M \leftarrow M + 1$ ;
13      $i_M \leftarrow \arg \max_{i \in \mathbb{I}} |(\mathbf{R}_{M-1}(\bar{\mu}_M))_i|$ ;
14      $\Phi_M \leftarrow \frac{\mathbf{R}_{M-1}(\bar{\mu}_M)}{|(\mathbf{R}_{M-1}(\bar{\mu}_M))_{i_M}|}$ ;
15      $\bar{S}_M \leftarrow \bar{S}_{M-1} \cup \{\bar{\mu}_M\}$ ,  $\bar{W}_M \leftarrow \bar{W}_{M-1} \oplus \text{span}\{\{\}\Phi_M\}$ ,  $\mathcal{I}_M = \mathcal{I}_{M-1} \cup \{i_M\}$ ;
16   end
17 until  $M \geq M_{\max}$ ;

```

---

not trivial if we want to guarantee the  $\mathcal{N}$ -independence of the online phase. We have to assemble this vector for each new evaluation of the coefficients  $\beta_m(u, \mu)$ . Since  $\mathbf{T}(u, \mu)$  depends non-affinely of the parameter  $\mu$ , we cannot pre-compute anything. A trivial solution would be to reassemble the *full* tensor  $\mathbf{T}(u, \mu)$  and then to extract the entries corresponding to the indices in  $\mathcal{I}_M$ . However, this assembly depends directly on the FEM dimension  $\mathcal{N}$ , and is potentially very costly. To remain independent of this high dimension  $\mathcal{N}$  we have to assemble the tensor only on the necessary degrees of freedom. We detail this procedure in the next paragraphs and its implementation in chapter 10.

Thanks to the compact support nature of the basis functions, we only need to save a limited number of data to rebuild a reduced tensor  $\mathbf{T}_{\mathcal{I}_M}(u(\mu), \mu)$ . We first define the set  $D_{\mathcal{I}_M}$  as the set of degrees of freedom (DoF) associated with the indices in  $\mathcal{I}_M$ . Then we define  $N_{\mathcal{I}_M}$  the set of nodes in the triangulation  $\mathcal{T}_h$  of the domain  $\Omega$ , associated to the DoFs in  $D_{\mathcal{I}_M}$ . From this set of nodes  $N_{\mathcal{I}_M}$ , we can extract an interpolation mesh  $\mathcal{T}_{\mathcal{I}_M}$  out of the original  $\mathcal{T}_h$ . This sub-mesh is composed with the elements of  $\mathcal{T}_h$  containing at least one node of  $N_{\mathcal{I}_M}$ . The interpolation mesh is then of very low dimension  $\mathcal{O}(M) \ll \mathcal{N}$ . We construct and store this mesh at the end of the EIM greedy algorithm when the space  $\bar{W}_M$  is complete. From this interpolation mesh, we can build an interpolation function space  $X_{\mathcal{I}_M} = X_h(\mathcal{T}_{\mathcal{I}_M})$ . This low dimensional function space will be used to assemble the tensor  $\mathbf{T}_{\mathcal{I}_M}(u(\mu), \mu)$ , for any  $\mu \in \mathcal{D}$ , only on the needed DoFs.

## 2.2 Operator on a Product Space

An issue appeared when we applied EIM to discrete operators on product of spaces: when the physics corresponding to the different spaces do not have the same order of magnitude, the use of EIM on the whole operator may totally “miss” the non-affine dependency of one (or more) block in the composite operator. We propose here to build multiple EIM, one for each block, to correctly decompose each sub-operator.

Let  $W = X^1, \dots, X^{N_s}$ ,  $N_s$  Hilbert spaces. We also introduce  $N_s$  FE space  $W_h = X_h^1, \dots, X_h^{N_s}$  such that  $X_h^i \subset X^i$ ,  $\forall 1 \leq i \leq N_s$ . The respective dimension of each FE space  $X_h^i$  is denoted by  $\mathcal{N}_i$  and  $\mathcal{N} = \sum_{i=1}^{N_s} \mathcal{N}_i$ . We consider the discrete parametrized operator  $\mathbf{T}(w(\boldsymbol{\mu}); \boldsymbol{\mu})$ , typically obtained by the projection of a parametrized system of PDE on  $W_h$ . To set these ideas, we take here the example of a bilinear discrete operator  $\mathbf{T}(\cdot; \cdot) : W \times \mathcal{D} \rightarrow \mathbb{R}^{\mathcal{N} \times \mathcal{N}}$ . The assembled matrix  $\mathbf{T}(w(\boldsymbol{\mu}); \boldsymbol{\mu})$  can be written under a block form,

$$\mathbf{T}(w(\boldsymbol{\mu}); \boldsymbol{\mu}) = \begin{bmatrix} \mathbf{M}^{11}(w(\boldsymbol{\mu}); \boldsymbol{\mu}) & \cdots & \mathbf{M}^{1N_s}(w(\boldsymbol{\mu}); \boldsymbol{\mu}) \\ \vdots & & \vdots \\ \mathbf{M}^{N_s 1}(w(\boldsymbol{\mu}); \boldsymbol{\mu}) & \cdots & \mathbf{M}^{N_s N_s}(w(\boldsymbol{\mu}); \boldsymbol{\mu}) \end{bmatrix}. \quad (6.23)$$

From this point we can consider two options. The first one is to perform an EIM algorithm on the whole matrix/tensor  $\mathbf{T}$ . The other one is to perform multiple EIM algorithm on each submatrices  $\mathbf{M}_{ij}$ . With this second solution we ensure that all the different physics of the system have a proper affine approximation of maximum rank  $M$ . On the other hand, performing an EIM on each block can be very costly and not necessary if the physics are of comparable order of magnitude.

The affine approximation of maximum rank  $M$  for  $\mathbf{T}$  is then express, in a block version,

$$\mathbf{T}_M(w(\boldsymbol{\mu}); \boldsymbol{\mu}) = \begin{bmatrix} \mathbf{M}_M^{11}(w(\boldsymbol{\mu}); \boldsymbol{\mu}) & \cdots & \mathbf{M}_M^{1N_s}(w(\boldsymbol{\mu}); \boldsymbol{\mu}) \\ \vdots & & \vdots \\ \mathbf{M}_M^{N_s 1}(w(\boldsymbol{\mu}); \boldsymbol{\mu}) & \cdots & \mathbf{M}_M^{N_s N_s}(w(\boldsymbol{\mu}); \boldsymbol{\mu}) \end{bmatrix} \quad (6.24)$$

where, for any  $1 \leq i, j \leq N_s$ ,

$$\mathbf{M}_M^{ij}(w; \boldsymbol{\mu}) = \sum_{m=1}^M \beta_m^{ij}(w; \boldsymbol{\mu}) \boldsymbol{\Phi}_m^{ij}. \quad (6.25)$$

**Remark 10.** A lot of the  $\beta_m^{ij}$  in this generic formulation might be zero in practice. Also, we do not need to run the EIM on each block but only on blocks with non-affine / non-linear contributions.

### 3 Simultaneous Empirical Interpolation and Reduced Basis Method for Non-linear Problems

We detailed in the two previous section the EIM algorithm allowing to recover an affine approximation for non-affine and non-linear problems. During the greedy algorithm, an important part is the selection of the parameters  $\bar{\mu}_m \in \bar{S}_M$ . For any rank  $M \geq 2$ , the new parameter  $\bar{\mu}_M$  is chosen as

$$\bar{\mu}_M = \arg \max_{\mu \in \bar{\Xi}} \|\mathbf{R}_{M-1}(\mu)\|_\infty = \arg \max_{\mu \in \bar{\Xi}} \|\mathbf{T}(u(\mu), \mu) - \mathbf{T}_{M-1}(u(\mu), \mu)\|_\infty. \quad (6.26)$$

The evaluation of  $\bar{\mu}_M$  then requires to compute  $\mathbf{T}(u(\mu), \mu)$  for any  $\mu \in \bar{\Xi}$ . In the case of non-linear problems, this implies to solve a non-linear FEM problem for each  $\mu \in \bar{\Xi}$ . Depending of the dimension of the problem, this is often really costly and not reachable for too large sampling  $\bar{\Xi}$ . A solution, proposed in [22], is to use the Simultaneous Empirical Interpolation and Reduced Basis (SER) method in order to build alternatively the RB space  $X_N = \text{span}\{\xi_n\}$  and the EIM space  $\bar{W}_M = \text{span}\{\Phi_m\}$ . The RB approximation  $u_N(\mu)$  is then used in the evaluation of the residual  $\mathbf{R}_M(\mu)$  instead of the FEM solution  $u_h(\mu)$ .

The SER algorithm was already interfaced with EIM in the Feel++ library. However, we interfaced the method with the discrete version described in the previous section. We present the theory of the method in this section. The implementation work is detailed in chapter 10.

We consider a non-linear parametrized problem which requires an affine approximation for the tensor  $\mathbf{T}(u(\mu), \mu)$ . Since the reduced basis methodology needs an affine decomposition, SER starts with the initialization of the EIM algorithm, for  $M = 1$ . During this initialization, the FEM problem is solved once to build the first EIM basis vector  $\Phi_1$ , see (6.18). With this first EIM basis vector, we have a first affine approximation  $\mathbf{T}_1(u(\mu), \mu) = \beta_1(u(\mu), \mu)\Phi_1, \forall \mu \in \mathcal{D}$ . This approximation is totally inaccurate but will be used as a starting point to build the first RB vector  $\xi_1$ , using the classical RBM detail in chapter 5.

Once the first RB basis is computed, we dispose of a -very bad- first RB approximation  $u_1(\mu)$ , for any  $\mu \in \mathcal{D}$ . This RB approximation will be used to evaluate the EIM residual  $\mathbf{R}_{1,1}(\mu)$  and chose the next parameter  $\bar{\mu}_2$ ,

$$\bar{\mu}_2 = \arg \max_{\mu \in \bar{\Xi}} \|\mathbf{R}_{1,1}(\mu)\|_\infty = \arg \max_{\mu \in \bar{\Xi}} \|\mathbf{T}(u_1(\mu), \mu) - \mathbf{T}_1(u_1(\mu), \mu)\|_\infty \quad (6.27)$$

This iterative construction is summarized in the algorithm 20, for the sake of clarity, we used an elementary version of the offline RB construction but any advanced features presented in chapter 5 could be included in the loop.

SER algorithm iterates between the construction of a new EIM basis vector  $\Phi_M$  and a new RB vector  $\xi_N$ . The important features of this coconstruction are

- The use of the RB approximation  $u_N(\boldsymbol{\mu})$  in the evaluation of  $\mathbf{R}_{M,N}(\boldsymbol{\mu}) = \mathbf{T}(u_N(\boldsymbol{\mu}), \boldsymbol{\mu}) - \mathbf{T}_M(u_N(\boldsymbol{\mu}), \boldsymbol{\mu})$ . This residual is used in the greedy algorithm to select the next parameter. This major improvement changes the number of FEM resolutions needed from  $N_{\text{EIM+RB}} = \text{card}(\bar{\Xi}) + N$  to  $N_{\text{SER}} = N + M$ , where  $N$  is the dimension of the RB space and  $M$  the dimension of the EIM space.
- The update of the affine decomposition throughout the RB greedy algorithm.
- The simultaneously growing accuracy of both approximations  $u_N(\boldsymbol{\mu})$  and  $\mathbf{T}_M(u, \boldsymbol{\mu})$ . This growing accuracy allows recovering a very good convergence after some iterations in the SER algorithm.

### Difference with the Original Algorithm

SER algorithm was first presented in [24] and [22]. In these references, the RB approximation is used in the EIM greedy algorithm to select the parameters but also to build the EIM basis vectors. After some numerical experiments, we decided to use the RB approximation only during the selection process. The basis vectors are still built using the FEM solutions. This modification of the SER algorithm greatly increases the accuracy of the EIM expansion. However, the number of FEM solutions to be computed changes from  $N_{\text{SER}^*} = N + 1$  to  $N_{\text{SER}} = N + M$ .

This increase is not negligible. However, in the original SER algorithm (denoted  $\text{SER}^*$  hereafter), the multi-levels  $\text{SER}^*(l)$  is usually used. It consists of performing multiple passes in SER algorithm, in order to improve the RB approximation. Each new level is built using the RB approximation of the previous level. With this procedure, the number of FEM evaluations is now  $N_{\text{SER}^*(l)} = lN + 1$ .

With our version of SER, the multi-levels method becomes useless since we already used a truth solution for the construction. We compare the convergence of SER and  $\text{SER}^*$  in the following preliminary results.

## 4 Preliminary Results

To illustrate the method and also to validate our implementation of the discrete SER algorithm, we propose some numerical results on a 2D non-linear and non-affinely parametrized benchmark, proposed in [37]. The discrete version of EIM does not have a great interest in this test case since we have an analytic expression of the non-linear term. However, this benchmark is excellent to validate our implementation. The interest of the discrete EIM will be illustrated in the chapter 7 with the treatment of the stabilization operators.

**Algorithm 6:** SER Algorithm

---

```

1 Choose  $\Xi \subset \mathcal{D}$ ,  $\bar{\Xi} \subset \mathcal{D}$ ; // Super Samplings
2 Choose  $N_{\max} \in \mathbb{N}$ ,  $M_{\max} \in \mathbb{N}$ ; // Max size for  $X_N$  and  $\bar{W}_M$ 
3 Choose  $\delta_{EIM}$ ,  $\delta_{RB}$ ; // Algorithm tolerances
4 Choose  $\mu_1 \in \Xi$ ; // Initialization of RB construction
5 Set  $N \leftarrow 0$ ,  $S_0 \leftarrow \{\}$ ,  $X_0 \leftarrow \emptyset$ ;
6 Choose  $\bar{\mu}_1 \in \bar{\Xi}$ ; // Initialization of EIM construction
7 Set  $i_M \leftarrow \arg \max_{i \in \mathbb{I}} |\mathbf{T}(u_h(\bar{\mu}_1), \bar{\mu}_1)_i|$ ,  $\Phi_M \leftarrow \frac{\mathbf{T}(u_h(\bar{\mu}_1), \bar{\mu}_1)}{|(\mathbf{T}(u_h(\bar{\mu}_1), \bar{\mu}_1))_{i_1}|}$ ;
8 Set  $M \leftarrow 1$ ,  $\bar{S}_1 \leftarrow \{\bar{\mu}_1\}$ ,  $W_1 \leftarrow \{\Phi_1\}$ ,  $\mathcal{I}_1 \leftarrow \{i_1\}$ ;
9 Set  $\text{state}_{EIM} \leftarrow \text{todo}$ ,  $\text{state}_{RB} \leftarrow \text{todo}$ ;
10 repeat
11   if  $\text{state}_{RB} = \text{todo}$  then // RB Iteration
12      $N \leftarrow N + 1$ ;
13      $S_N \leftarrow S_{N-1} \cup \{\mu_N\}$ ;
14      $u_h(\mu_N) \leftarrow$  FE solution using current EIM affine approximation;
15      $X_N \leftarrow X_{N-1} \oplus \text{span}\{u_h(\mu_N)\}$ ;
16      $\mu_{N+1} = \arg \max_{\mu \in \Xi} \Delta_N(\mu)$ ;
17     if  $\Delta_N(\mu_{N+1}) \leq \delta_{RB}$  or  $N \geq N_{\max}$  then
18        $\text{state}_{RB} \leftarrow \text{finished}$ ;
19     end
20   end
21   if  $\text{state}_{EIM} = \text{todo}$  then
22      $\bar{\mu}_{M+1} \leftarrow \arg \max_{\mu \in \bar{\Xi}} \underbrace{\|\mathbf{T}(u_N(\mu), \mu) - \mathbf{T}_M(u_N(\mu), \mu)\|_{\infty}}_{\mathbf{R}_{M,N}(\mu)}$ ;
23     if  $\|\mathbf{R}_{M,N}(\bar{\mu}_{M+1})\|_{\infty} \leq \delta_{EIM}$  or  $M \geq M_{\max}$  then
24        $\text{state}_{EIM} \leftarrow \text{finished}$ ;
25     else
26        $M \leftarrow M + 1$ ;
27        $i_M \leftarrow \arg \max_{i \in \mathbb{I}} |(\mathbf{R}_{M-1}(\bar{\mu}_M))_i|$ ;
28        $\Phi_M \leftarrow \frac{\mathbf{R}_{M-1}(\bar{\mu}_M)}{|(\mathbf{R}_{M-1}(\bar{\mu}_M))_{i_M}|}$ ;
29        $\bar{S}_M \leftarrow \bar{S}_{M-1} \cup \{\bar{\mu}_M\}$ ,  $\bar{W}_M \leftarrow \bar{W}_{M-1} \oplus \text{span}\{\Phi_M\}$ ,  $\mathcal{I}_M = \mathcal{I}_{M-1} \cup \{i_M\}$ ;
30     end
31   end
32 until  $\text{state}_{EIM} = \text{finished}$  and  $\text{state}_{RB} = \text{finished}$ ;

```

---

## 4.1 Problem Formulation

The domain  $\Omega = [0, 1]^2$  is the unit square. The parameter space is  $\mathcal{D} = [0.01, 10]^2$ . We consider the following elliptic and non linear equation,  $\forall \boldsymbol{\mu} = (\mu_1, \mu_2) \in \mathcal{D}$

$$-\Delta u(\boldsymbol{\mu}) + \mu_1 \frac{e^{\mu_2 u(\boldsymbol{\mu})} - 1}{\mu_2} = 100 \sin(2\pi x) \sin(2\pi y), \quad (6.28)$$

and  $u(\boldsymbol{\mu}) = 0$  on  $\Gamma = \partial\Omega$ .

For this case, the output  $s$  is the average of  $u$  on the whole domain,

$$s(\boldsymbol{\mu}) = \int_{\Omega} u(\boldsymbol{\mu}) \, d\Omega. \quad (6.29)$$

We introduce a discretization of  $\Omega$  and a FE space  $X_h = \mathbb{L}_{ch}^k(\Omega)$ . We can then write the weak formulation,  $\forall \boldsymbol{\mu} \in \mathcal{D}$ , as find  $u_h(\boldsymbol{\mu}) \in X_h$  such that

$$\int_{\Omega} \nabla u_h(\boldsymbol{\mu}) \cdot \nabla v_h \, d\Omega + \int_{\Omega} \mu_1 \frac{e^{\mu_2 u_h(\boldsymbol{\mu})} - 1}{\mu_2} v_h \, d\Omega = \int_{\Omega} 100 \sin(2\pi x) \sin(2\pi y) v_h \, d\Omega, \quad \forall v_h \in X_h \quad (6.30)$$

The term  $\mu_1 \frac{e^{\mu_2 u_h(\boldsymbol{\mu})} - 1}{\mu_2}$  will be approximated by an EIM approximation

$$\mu_1 \frac{e^{\mu_2 u_h(\boldsymbol{\mu})} - 1}{\mu_2} = \sum_{m=1}^M \beta_m(\boldsymbol{\mu}) g^m. \quad (6.31)$$

This affine approximation will be used to build the RB approximation  $u_N(\boldsymbol{\mu})$ .

## 4.2 Numerical Results

The simulation is made with  $X_h = \mathbb{L}_{ch}^2(\Omega)$ . The mesh is basically a  $33 \times 33$  grid of simplex. We present a profile of the solution for an arbitrary chosen  $\boldsymbol{\mu}$ , on figure 6.1. The non-linearity is solved using a fixed point algorithm during the offline and the online phases. We choose a uniform grid  $15 \times 15$  as sampling for the EIM greedy algorithm. A thinner discretization would be very costly, at least for the classical EIM greedy algorithm. We used the dual norm of the residual as an error indicator in the RB greedy selection.

We compare the truth approximation  $u_h(\boldsymbol{\mu})$ , computed with the exact formulation, and the RB approximation  $u_N(\boldsymbol{\mu})$  computed using the EIM affine approximation. We are particularly interested in the quantities

$$\begin{aligned} e_N^u(\boldsymbol{\mu}) &= \frac{\|u_h(\boldsymbol{\mu}) - u_N(\boldsymbol{\mu})\|_0}{\|u_h(\boldsymbol{\mu})\|_0}, \\ e_N^s(\boldsymbol{\mu}) &= \frac{|s_h(\boldsymbol{\mu}) - s_N(\boldsymbol{\mu})|}{|s_h(\boldsymbol{\mu})|}. \end{aligned} \quad (6.32)$$

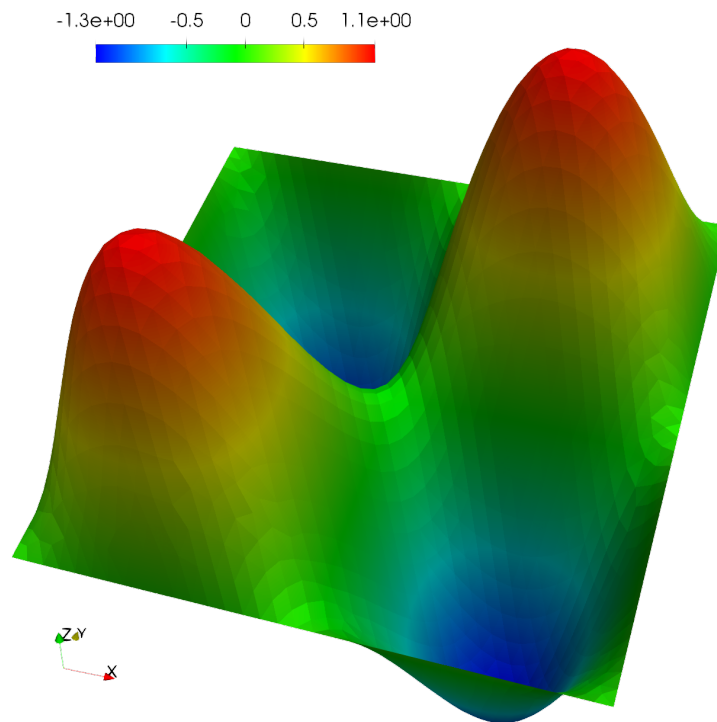


Figure 6.1 – Example of the solution for problem (6.28)  
Profile of the FE solution  $u_h(\mu)$  with  $\mu = (1.91, 2.77)$ .

This comparison is made for 100 different parameters, chosen randomly in  $\mathcal{D}$ . We study the convergence of these errors with  $M = 25$  EIM basis vectors. We also considered three different configurations: EIM with classical greedy, SER and SER\* (see paragraph 32). We plotted the evolution of these errors for each configuration on figures 6.2 and 6.3. We also compare the maximal errors for each method on figure 6.4.

The variant SER\*, which uses the RB approximation to build the EIM basis vector, is apparently less accurate than our version of the SER algorithm. This method will not be used anymore. The convergence is slightly slower than the convergence of the EIM solution, but eventually, the errors are comparable.

The results of this benchmark validate our implementation of SER with our new EIM. The convergence study is very promising concerning the use of SER. It is now time to test these methods on more complicated cases. We propose an interesting application in the next chapter, with the reduction of the stabilization operators introduced in 3.



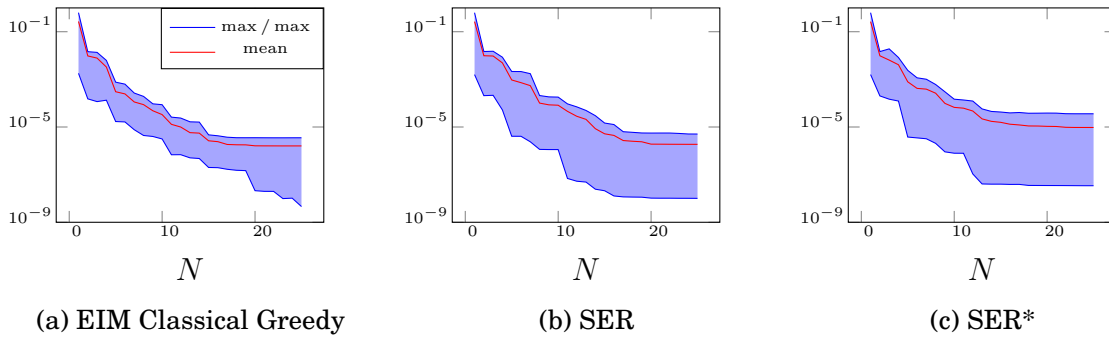


Figure 6.2 – Convergence of the errors  $e_N^u(\boldsymbol{\mu}) = \frac{\|u_h(\boldsymbol{\mu}) - u_N(\boldsymbol{\mu})\|_0}{\|u_h(\boldsymbol{\mu})\|_0}$  with respect to the number of vectors in the reduced basis. Maximum, minimum and mean values evaluated from 100 approximations.  $M_{\max} = 25$

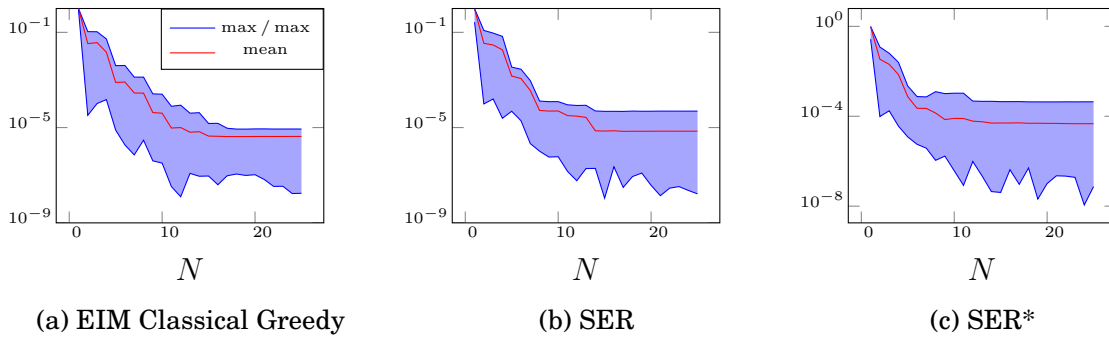


Figure 6.3 – Convergence of the errors  $e_N^s(\boldsymbol{\mu}) = \frac{|s_h(\boldsymbol{\mu}) - s_N(\boldsymbol{\mu})|}{|s_h(\boldsymbol{\mu})|}$  with respect to the number of vectors in the reduced basis. Maximum, minimum and mean values evaluated from 100 approximations.  $M_{\max} = 25$

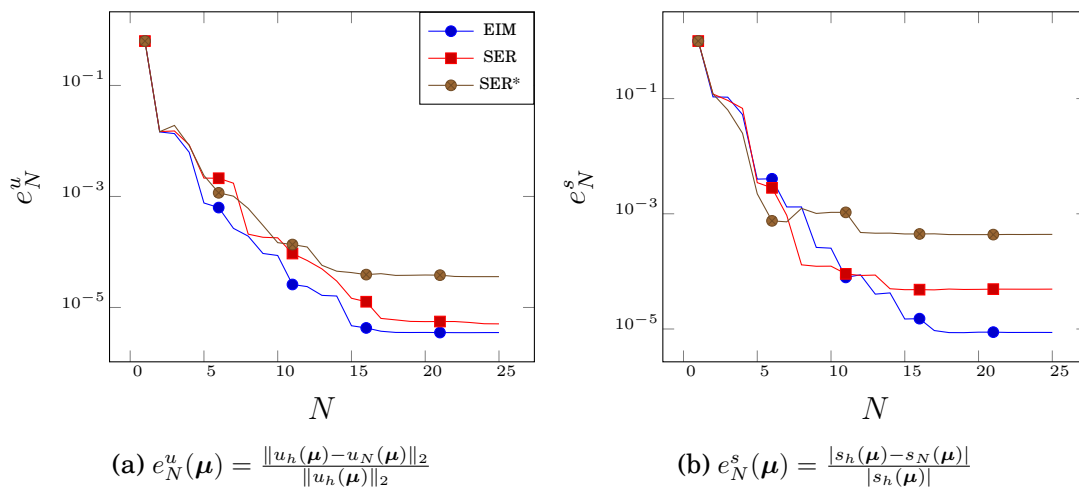


Figure 6.4 – Comparison of the maximum error over 100 approximations, for different methods.  $M_{\max} = 25$

## Conclusion

In this chapter, we detailed the methods which allow applying the RBM to non-linear and non-affine problems. The EIM is the key to recover an affine approximation and then to keep an efficient offline/online procedure. We presented a discrete version of EIM, recently implemented in the Feel++ RB framework. This new implementation of the method allows recovering an affine approximation of vectors or matrices and then to tackle more general problems. The EIM can also separate the non-linear dependency from the spatial one. It allows to use iterative solvers during the online phase and to preserve the  $\mathcal{N}$ -independence. The offline phase of EIM for non-linear problems is greatly improved using the SER construction. We proposed a slight modification of the original algorithm to improve the accuracy of the approximation provided by the coconstruction.



---

# The Reduced Basis Method Applied to Aerothermal Problems

---

## Contents

---

<b>1</b>	<b>Reduction of the Aerothermal Model . . . . .</b>	<b>130</b>
1.1	Offline Procedure: Construction of the RB Space . . . . .	131
1.2	Online Evaluation . . . . .	132
1.3	Treatment of the Stabilization Terms . . . . .	136
1.4	Reduction of a Turbulent Model . . . . .	137
<b>2</b>	<b>OPUS Test-Case: Application to the Cooling of Electronic Components . . . . .</b>	<b>139</b>
2.1	Problem Description . . . . .	139
2.2	Numerical Results . . . . .	143
<b>3</b>	<b>CHORUS Test-Case: Application to Aerothermal Simulations in an Airplane . . . . .</b>	<b>150</b>
3.1	Problem Description . . . . .	150
3.2	Numerical Results . . . . .	153

---

The reduction of an aerothermal problems includes some non-trivial issues. In this chapter, we detail the choices we made and the numerical solutions we used to build a reduced version of the framework presented in the first part.

First of all, we have to deal with the Navier-Stokes equations. [78, 49] proposed a first stable approximation for the incompressible Navier-Stokes flow. Rigorous a

posteriori error estimator were first developed later in [102, 21] and more recently in [68, 82]. Using the quadratic non-linearity of the Navier-Stokes equations, the RBM introduced in the chapter 5 can be adapted with no major difficulties. We usually build a stable RB space - in the sense of the Brezzi-Rappaz-Raviart (BRR) condition - by enriching the velocity space as detailed in the section 2.

The RB method for trilinear problems has already been used for natural convection problems [91, 26, 56]. However, the construction of the RB space is usually made monolithically. We propose in this section to build a product space  $V_N \times Q_N \times X_N$  to improve the stability and the accuracy of the reduced approximation.

The introduction of the stabilization operators will be treated with the EIM method presented in the previous chapter. This solution was already used for convection dominated problems using DEIM [75]. We proposed here to extend the idea to the whole aerothermal system. Other approaches were proposed for the stabilization of reduced convection dominated problems. In [64] the stabilization in the online stage is made independently of the offline one by adding an appropriate vanishing viscosity in the high RB modes. In [36, 107], the SUPG stabilization method is used with Proper Orthogonal Decomposition (POD) to reduce transient convection dominated problems.

In the section 1, we detail the construction of our reduced aerothermal framework, using the RBM for trilinear multi-physic problems and the EIM for the stabilization operators. Then, the following sections are dedicated to numerical applications. The first one is a simulation with no natural convection in the neighborhood of an electronic component submitted to cooling air flow. The second application is a simulation of mixed natural and forced convection in an airplane cabin proposed by Airbus Group.

## 1 Reduction of the Aerothermal Model

We consider the problem describe in chapter 2, and we remind here the system of equation,

$$\left\{ \begin{array}{l} \rho \mathbf{u} \cdot \nabla \mathbf{u} + \nabla p - 2 \nabla \cdot (\mu \mathbf{S}(\mathbf{u})) = -\rho \beta (T - T_0) \mathbf{g}, \\ \nabla \cdot \mathbf{u} = 0, \\ \mathbf{u} \cdot \nabla T - \nabla \cdot (\kappa \nabla T) = 0, \\ \text{+Boundary Conditions 2.14.} \end{array} \right. \quad (7.1)$$

We also introduce the Hilbert spaces  $V$ ,  $Q$  and  $X$  with the scalar products  $(\cdot, \cdot)_V$ ,  $(\cdot, \cdot)_Q$ ,  $(\cdot, \cdot)_X$  and associated norms  $\|\cdot\|_V$ ,  $\|\cdot\|_Q$ ,  $\|\cdot\|_X$ , respectively. We define the product space  $Z = V \times Q$  with the scalar product  $((\mathbf{u}, p), (\mathbf{v}, q))_Z = (\mathbf{u}, \mathbf{v})_V + (p, q)_Q$  and associated norm  $\|(\cdot, \cdot)\|_Z = \sqrt{\|\cdot\|_V^2 + \|\cdot\|_Q^2}$ . We finally note the full composite space  $W = V \times Q \times X$  with scalar product  $(\cdot, \cdot)_W$  and norm  $\|\cdot\|_W$ .

We supposed that the problem 7.1 is parametrized by a parameter  $\mu \in \mathcal{D} \in \mathbb{R}^P$

of dimension  $P \geq 1$ . The discretization of this system, using the finite element method (see chapter 2, section 2), leads to the weak problem: for any  $\boldsymbol{\mu} \in \mathcal{D}$ , find  $(\mathbf{u}_h(\boldsymbol{\mu}), p_h(\boldsymbol{\mu}), T_h(\boldsymbol{\mu})) \in V_h \times Q_h \times X_h$  such that

$$\begin{aligned} c(\mathbf{u}_h, \mathbf{u}_h, \mathbf{v}_h; \boldsymbol{\mu}) + d_u(\mathbf{u}_h, \mathbf{v}_h; \boldsymbol{\mu}) - b(\mathbf{v}_h, p_h; \boldsymbol{\mu}) + b(\mathbf{u}_h, q_h; \boldsymbol{\mu}) &= f(T_h, \mathbf{v}_h; \boldsymbol{\mu}), \\ a(T_h, \mathbf{u}_h, S_h; \boldsymbol{\mu}) + d_T(T_h, S_h; \boldsymbol{\mu}) &= g(S_h; \boldsymbol{\mu}), \\ \forall (\mathbf{v}_h, q_h, S_h) &\in V_h \times Q_h \times X_h, \end{aligned} \quad (7.2)$$

where  $\forall \boldsymbol{\mu} \in \mathcal{D}$ ,

- $c : V \times V \times V \rightarrow \mathbb{R}$  and  $a : X \times V \times X \rightarrow \mathbb{R}$  are continuous trilinear forms
- $d_u : V \times V \rightarrow \mathbb{R}$ ,  $b : V \times Q \rightarrow \mathbb{R}$ ,  $f : X \times V \rightarrow \mathbb{R}$  and  $d_T : X \times X \rightarrow \mathbb{R}$  are continuous bilinear forms
- $g : X \rightarrow \mathbb{R}$  is a continuous linear form.

We also assume that these forms can be written using an affine decomposition,  $\forall \boldsymbol{\mu} \in \mathcal{D}$ ,

$$\begin{aligned} c(\mathbf{u}, \mathbf{v}, \mathbf{w}; \boldsymbol{\mu}) &= \sum_{k=1}^{Q_c} \theta_k^c(\boldsymbol{\mu}) c_k(\mathbf{u}, \mathbf{v}, \mathbf{w}), & a(T, \mathbf{u}, S; \boldsymbol{\mu}) &= \sum_{k=1}^{Q_a} \theta_k^a(\boldsymbol{\mu}) a_k(T, \mathbf{u}, S), \\ d_u(\mathbf{u}, \mathbf{v}; \boldsymbol{\mu}) &= \sum_{k=1}^{Q_{d_u}} \theta_k^{d_u}(\boldsymbol{\mu}) d_{u,k}(\mathbf{u}, \mathbf{v}), & d_T(T, S; \boldsymbol{\mu}) &= \sum_{k=1}^{Q_{d_T}} \theta_k^{d_T}(\boldsymbol{\mu}) d_{T,k}(T, S), \\ b(\mathbf{u}, q; \boldsymbol{\mu}) &= \sum_{k=1}^{Q_b} \theta_k^b(\boldsymbol{\mu}) b_k(\mathbf{u}, q), & f(T, \mathbf{v}; \boldsymbol{\mu}) &= \sum_{k=1}^{Q_f} \theta_k^f(\boldsymbol{\mu}) f_k(T, \mathbf{v}) \\ g(\mathbf{v}; \boldsymbol{\mu}) &= \sum_{k=1}^{Q_g} \theta_k^g(\boldsymbol{\mu}) g_k(\mathbf{v}) \end{aligned} \quad (7.3)$$

## Truth Approximation

As truth approximation of this problem, we will use the finite element solution. The resolution of this system of equation is widely described in the part I. Thereafter we assume that  $\forall \boldsymbol{\mu} \in \mathcal{D}$  we dispose of a FE solution  $(\mathbf{u}_h(\boldsymbol{\mu}), p_h(\boldsymbol{\mu}), T_h(\boldsymbol{\mu})) \in V \times Q \times X$  of the system 7.1.

### 1.1 Offline Procedure: Construction of the RB Space

The reduced basis space will be built on a sampling  $S_N$ , as described in the chapter 5. This sampling can be iteratively determined using a greedy algorithm, if

we have an error indicator for our problem, or chosen randomly otherwise. In both cases, the aim of the offline procedure is to iteratively construct three sets of spaces  $V_1 \subset \dots \subset V_N$ ,  $Q_1 \subset \dots \subset Q_N$  and  $X_1 \subset \dots \subset X_N$ .

The first spaces  $V_1$ ,  $Q_1$  and  $X_1$  are built as

$$\begin{aligned} V_1 &= \text{span}\{\mathbf{u}_h(\boldsymbol{\mu}_1), \mathbf{T}^{\boldsymbol{\mu}_1} \eta_1\} = \text{span}\{\zeta_1, \zeta_2\}, \\ Q_1 &= \text{span}\{p_h(\boldsymbol{\mu}_1)\} = \text{span}\{\eta_1\}, \\ X_1 &= \text{span}\{T_h(\boldsymbol{\mu}_1)\} = \text{span}\{\xi_1\}. \end{aligned} \quad (7.4)$$

As mentioned in the chapter 5, the actual basis  $\{\zeta_n\}$ ,  $\{\eta_n\}$  and  $\{\xi_n\}$  are obtained from an orthonormalization of the respective snapshots  $\{\mathbf{u}_h(\boldsymbol{\mu}_n), \mathbf{T}^{\boldsymbol{\mu}_n} \eta_n\}$ ,  $\{p_h(\boldsymbol{\mu}_n)\}$  and  $\{T_h(\boldsymbol{\mu}_n)\}$ .  $\mathbf{T}^{\boldsymbol{\mu}}$  is the supremizer operator, defined in the chapter 5, section 2. However, the random selection of the  $\boldsymbol{\mu}_n$  or the greedy selection might provide snapshots that already live in the RB spaces. For instance, if the difference between the parameters  $\boldsymbol{\mu}_{n-1}$  and  $\boldsymbol{\mu}_n$  is only acting on the fluid temperature and not on the fluid dynamic, the new snapshots  $(\mathbf{u}_h(\boldsymbol{\mu}_n), p_h(\boldsymbol{\mu}_n))$  will not be relevant to enrich the space  $V_n$  and  $Q_n$ . It is then important to eliminate these useless contributions.

This selection phase is made during the orthonormalization process. When a new solution  $\mathbf{u}_h(\boldsymbol{\mu}_n)$  is computed for the parameter  $\boldsymbol{\mu}_n$ , we first evaluate its component orthogonal to the space  $V_{n-1}$ ,

$$\zeta_n^* = \mathbf{u}_h(\boldsymbol{\mu}_n) - \sum_{k=1}^{N_{n-1}^u} (\mathbf{u}_h(\boldsymbol{\mu}_n), \zeta_k)_V \zeta_k, \quad (7.5)$$

where  $N_{n-1}^u$  is the size of the space  $V_{n-1}$ . At this point, if the norm,  $\|\zeta_n^*\|_V$ , is too small, this solution is not normalized and not added to the space  $V_n$ . We do the same verification during the orthonormalization of the bases of the spaces  $Q_n$  and  $X_n$ . This way, we ensure that the vectors of each reduced basis remain linearly independent. A direct consequence of this construction is the variable sizes of the RB spaces. From now on, we denote by  $N$  the size of the sampling  $S_N$ . We denote respectively by  $N_N^u$ ,  $N_N^p$  and  $N_N^T$  the dimensions of the RB spaces  $V_N$ ,  $Q_N$  and  $X_N$ . We have the immediate, but notable relations:  $N_N^u \leq 2N$ ,  $N_N^p \leq N$  and  $N_N^T \leq N$ .

We present a summarized version of the selective Gram-Schmidt procedure in the algorithm 7. We also propose a greedy selection for the aerothermal problems in the algorithm 8.

## 1.2 Online Evaluation

The idea of the online evaluation is similar to the coercive or saddle point problems presented in the chapter 5. We perform a Galerkin projection of our continuous problem on the reduced space  $W_N = V_N \times Q_N \times X_N$ . The main difference with the current system of equations is the treatment of the non-linearity. We

**Algorithm 7: Selective Gram-Schmidt Orthonormalization**


---

```

Input:  $\{\xi_1, \dots, \xi_N\}, u_n$ 
Output:  $\{\xi_1, \dots, \xi_N\}$ 
1  Choose  $\delta_{tol}$ ; // Algorithm tolerance
2   $\xi^* \leftarrow u_n$ ;
3  for  $1 \leq k \leq N$  do
4  |  $\xi^* \leftarrow \xi^* - (u_n, \xi_k)_{X_h}$ ;
5  end
6  if  $\|\xi^*\|_X \geq \delta_{tol}$  then
7  |  $N \leftarrow N + 1$ ;
8  | return  $\{\xi_1, \dots, \xi_N\} \cup \{\frac{1}{\|\xi^*\|_X} \xi^*\}$ 
9  else
10 | return  $\{\xi_1, \dots, \xi_N\}$ 
11 end

```

---

**Algorithm 8: Offline Procedure for Aerothermal Problems**


---

```

1  Choose  $\Xi \subset \mathcal{D}$ ; // Super Sampling
2  Choose  $N_{max}$ ; // Maximum size of  $X_N$ 
3  Choose  $\delta_{tol}$ ; // Algorithm tolerance
4  Set  $S_N \leftarrow \{\}, V_N \leftarrow \{\}, Q_N \leftarrow \{\}, X_N \leftarrow \{\}$ ; // Initialization
5  Set  $N_0^u \leftarrow 0, N_0^p \leftarrow 0, N_0^T \leftarrow 0$ ;
6  repeat
7  |  $S_N \leftarrow S_{N-1} \cup \{\mu_N\}, N_N^u \leftarrow N_{N-1}^u, N_N^p \leftarrow N_{N-1}^p, N_N^T \leftarrow N_{N-1}^T$ ;
8  | Find  $(u_h(\mu_N), p_h(\mu_N), T_h(\mu_N))$  solution of system (7.1);
9  |  $\{\eta_1, \dots, \eta_{N^p}\} \leftarrow$  Selective Gram-Schmidt( $\{\eta_1, \dots, \eta_{N^p}\}, p_h(\mu_N)$ );
10 | if  $p_h(\mu_N)$  added then
11 | |  $\{\zeta_1, \dots, \zeta_{N^u}\} \leftarrow$  Selective Gram-Schmidt( $\{\zeta_1, \dots, \zeta_{N^u}\}, T^\mu \eta_{N^p}$ );
12 | end
13 |  $\{\zeta_1, \dots, \zeta_{N^u}\} \leftarrow$  Selective Gram-Schmidt( $\{\zeta_1, \dots, \zeta_{N^u}\}, u_h(\mu_N)$ );
14 |  $\{\xi_1, \dots, \xi_{N^T}\} \leftarrow$  Selective Gram-Schmidt( $\{\xi_1, \dots, \xi_{N^T}\}, T_h(\mu_N)$ );
15 |  $V_N \leftarrow \text{span}\{\zeta_1, \dots, \zeta_{N^u}\}$ ;
16 |  $Q_N \leftarrow \text{span}\{\eta_1, \dots, \eta_{N^p}\}$ ;
17 |  $X_N \leftarrow \text{span}\{\xi_1, \dots, \xi_{N^T}\}$ ;
18 | Precompute online structures;
19 | for  $\mu \in \Xi$  do
20 | | Compute error bound,  $\Delta_N(\mu)$ ;
21 | end
22 |  $\mu_{N+1} = \arg \max_{\mu \in \Xi} \Delta_N(\mu)$ ; // Or can be chosen randomly
23 until  $\Delta_N(\mu_{N+1}) < \delta_{tol}$  or  $N \geq N_{max}$ ;

```

---



choose to use a Newton iterative method to deal with this non-linearity. Note that the method described after is easily adaptable to other iterative methods.

We assume that the RB spaces  $V_N$ ,  $Q_N$  and  $X_N$  have been built during the offline phase, with the respective orthonormal basis  $\mathcal{B}_N^u = \{\zeta_i\}_{i=1}^{N_N^u}$ ,  $\mathcal{B}_N^p = \{\eta_i\}_{i=1}^{N_N^p}$  and  $\mathcal{B}_N^T = \{\xi_i\}_{i=1}^{N_N^T}$ . The Galerkin projection of the problem (7.1) on these spaces, leads to the matricial equation: find  $(\mathbf{u}_N(\boldsymbol{\mu}), p_N(\boldsymbol{\mu}), T_N(\boldsymbol{\mu}))$  such that

$$\underbrace{\begin{bmatrix} \mathbf{D}_N^u(\boldsymbol{\mu}) + \mathbf{C}_N(\mathbf{u}_N; \boldsymbol{\mu}) & \mathbf{B}_N(\boldsymbol{\mu})^\top & -\mathbf{F}_N(\boldsymbol{\mu}) \\ -\mathbf{B}_N(\boldsymbol{\mu}) & 0 & 0 \\ 0 & 0 & \mathbf{A}_N(\mathbf{u}_N; \boldsymbol{\mu}) + \mathbf{D}_N^T(\boldsymbol{\mu}) \end{bmatrix}}_{M_N(\mathbf{u}_N; \boldsymbol{\mu})} \begin{bmatrix} \mathbf{u}_N(\boldsymbol{\mu}) \\ p_N(\boldsymbol{\mu}) \\ T_N(\boldsymbol{\mu}) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ G_N(\boldsymbol{\mu}) \end{bmatrix}, \quad (7.6)$$

where,  $\forall \boldsymbol{\mu} \in \mathcal{D}$  and for any  $1 \leq i, j \leq N_N^u$ ,  $1 \leq m \leq N_N^p$ ,  $1 \leq n, l \leq N_N^T$ ,

$$\begin{aligned} (\mathbf{D}_N^u(\boldsymbol{\mu}))_{ij} &= d_u(\zeta_j, \zeta_i; \boldsymbol{\mu}), & (\mathbf{C}_N(\mathbf{u}; \boldsymbol{\mu}))_{ij} &= c(\zeta_j, \mathbf{u}, \zeta_i; \boldsymbol{\mu}), & (\mathbf{F}_N(\boldsymbol{\mu}))_{i,l} &= f(\xi_l, \zeta_i; \boldsymbol{\mu}), \\ (\mathbf{B}_N(\boldsymbol{\mu}))_{mj} &= b(\zeta_j, \eta_m; \boldsymbol{\mu}) \\ (\mathbf{D}_N^T(\boldsymbol{\mu}))_{nl} &= d_T(\xi_l, \xi_n; \boldsymbol{\mu}), & (\mathbf{A}_N(\mathbf{u}; \boldsymbol{\mu}))_{nl} &= a(\xi_l, \mathbf{u}, \xi_n; \boldsymbol{\mu}), & (\mathbf{G}_N(\boldsymbol{\mu}))_n &= g(\xi_n; \boldsymbol{\mu}). \end{aligned} \quad (7.7)$$

The system (7.6) can be solved using a Newton method. For that, we have to evaluate,  $\forall (\mathbf{u}, p, T) \in V_N \times Q_N \times X_N$ ,  $\forall \boldsymbol{\mu} \in \mathcal{D}$  the Jacobian  $\mathbf{J}_N(\mathbf{u}, p, T, \boldsymbol{\mu})$  (or Fréchet Derivative) of the application  $\mathcal{F}(\boldsymbol{\mu}) : (\mathbf{u}, p, T) \rightarrow M(\mathbf{u}; \boldsymbol{\mu})(\mathbf{u}, p, T)^\top$ . We also need the residual  $\mathbf{R}_N(\mathbf{u}, p, T; \boldsymbol{\mu})$  of the system (7.6). The evaluation of these two objects during the online phase has to be independent of the finite element dimension  $\mathcal{N}$ .

Thanks to the quadratic form of our problem, we can easily write the Jacobian  $\mathbf{J}_N$ ,  $\forall (\mathbf{u}, p, T) \in V_N \times Q_N \times X_N$ ,  $\forall \boldsymbol{\mu} \in \mathcal{D}$

$$\mathbf{J}_N(\mathbf{u}, p, T; \boldsymbol{\mu}) = \begin{bmatrix} \mathbf{D}_N^u(\boldsymbol{\mu}) + \mathbf{C}_N(\mathbf{u}_N; \boldsymbol{\mu}) + \mathbf{C}_N^*(\mathbf{u}_N; \boldsymbol{\mu}) & \mathbf{B}_N(\boldsymbol{\mu})^\top & -\mathbf{F}_N(\boldsymbol{\mu}) \\ -\mathbf{B}_N(\boldsymbol{\mu}) & 0 & 0 \\ \mathbf{A}_N^*(T_N; \boldsymbol{\mu}) & 0 & \mathbf{A}_N(\mathbf{u}_N; \boldsymbol{\mu}) + \mathbf{D}_N^T(\boldsymbol{\mu}) \end{bmatrix} \quad (7.8)$$

where,  $\forall \boldsymbol{\mu} \in \mathcal{D}$  and for any  $1 \leq i, j \leq N_N^u$ ,  $1 \leq n, l \leq N_N^T$ ,

$$(\mathbf{C}_N^*(\mathbf{u}; \boldsymbol{\mu}))_{ij} = c(\mathbf{u}, \zeta_j, \zeta_i; \boldsymbol{\mu}), \quad (\mathbf{A}_N^*(\mathbf{u}; \boldsymbol{\mu}))_{nj} = a(T, \zeta_j, \xi_n; \boldsymbol{\mu}). \quad (7.9)$$

By using the affine decomposition assumption (2.2), we can rewrite the definitions of the matrices  $\mathbf{D}_N^u$ ,  $\mathbf{B}_N$ ,  $\mathbf{F}_N$  and  $\mathbf{D}_N^T$ ,  $\forall \boldsymbol{\mu} \in \mathcal{D}$  and for any  $1 \leq i, j \leq N_N^u$ ,  $1 \leq m \leq N_N^p$ ,  $1 \leq n, l \leq N_N^T$ ,

$$\begin{aligned} (\mathbf{D}_N^u(\boldsymbol{\mu}))_{i,j} &= \sum_{k=1}^{Q_{d^u}} \theta_k^{d^u}(\boldsymbol{\mu}) \underbrace{d_k^u(\zeta_j, \zeta_i)}_{\text{precomputed}}, & (\mathbf{D}_N^T(\boldsymbol{\mu}))_{n,l} &= \sum_{k=1}^{Q_{d^T}} \theta_k^{d^T}(\boldsymbol{\mu}) \underbrace{d_k^T(\xi_l, \xi_n)}_{\text{precomputed}} \\ (\mathbf{B}_N(\boldsymbol{\mu}))_{m,j} &= \sum_{k=1}^{Q_b} \theta_k^b(\boldsymbol{\mu}) \underbrace{b_k(\zeta_j, \eta_m)}_{\text{precomputed}}, & (\mathbf{F}_N(\boldsymbol{\mu}))_{i,l} &= \sum_{k=1}^{Q_f} \theta_k^f(\boldsymbol{\mu}) \underbrace{f_k(\xi_l, \zeta_i)}_{\text{precomputed}}. \end{aligned} \quad (7.10)$$

We use the same idea to precompute the quadratic contributions,  $\forall \boldsymbol{\mu} \in \mathcal{D}$ ,  $\forall (\mathbf{u}, p, T) \in V_N \times Q_N \times X_N$ ,

$$\begin{aligned}
(\mathbf{C}_N(\mathbf{u}; \boldsymbol{\mu}))_{ij} &= c(\boldsymbol{\zeta}_j, \mathbf{u}, \boldsymbol{\zeta}_i; \boldsymbol{\mu}) & (\mathbf{C}_N^*(\mathbf{u}; \boldsymbol{\mu}))_{ij} &= c(\mathbf{u}, \boldsymbol{\zeta}_j, \boldsymbol{\zeta}_i; \boldsymbol{\mu}) \\
&= \sum_{n=1}^{N_N^u} \mathbf{u}_n c(\boldsymbol{\zeta}_j, \boldsymbol{\zeta}_n, \boldsymbol{\zeta}_i; \boldsymbol{\mu}) & &= \sum_{n=1}^{N_N^u} \mathbf{u}_n c(\boldsymbol{\zeta}_n, \boldsymbol{\zeta}_j, \boldsymbol{\zeta}_i; \boldsymbol{\mu}) \\
&= \sum_{k=1}^{Q_c} \sum_{n=1}^{N_N^u} \theta_k^c(\boldsymbol{\mu}) \mathbf{u}_n \underbrace{c_k(\boldsymbol{\zeta}_j, \boldsymbol{\zeta}_n, \boldsymbol{\zeta}_i)}_{\text{precomputed}}, & &= \sum_{k=1}^{Q_c} \sum_{n=1}^{N_N^u} \theta_k^c(\boldsymbol{\mu}) \mathbf{u}_n \underbrace{c_k(\boldsymbol{\zeta}_n, \boldsymbol{\zeta}_j, \boldsymbol{\zeta}_i)}_{\text{precomputed}}, \\
(\mathbf{A}_N(\mathbf{u}; \boldsymbol{\mu}))_{ij} &= a(\boldsymbol{\xi}_j, \mathbf{u}, \boldsymbol{\xi}_i; \boldsymbol{\mu}) & (\mathbf{A}_N^*(T; \boldsymbol{\mu}))_{ij} &= a(T, \boldsymbol{\zeta}_j, \boldsymbol{\xi}_i; \boldsymbol{\mu}) \\
&= \sum_{n=1}^{N_N^u} \mathbf{u}_n a(\boldsymbol{\xi}_j, \boldsymbol{\zeta}_n, \boldsymbol{\xi}_i; \boldsymbol{\mu}) & &= \sum_{n=1}^{N_N^T} T_n a(\boldsymbol{\xi}_n, \boldsymbol{\zeta}_j, \boldsymbol{\xi}_i; \boldsymbol{\mu}) \\
&= \sum_{k=1}^{Q_a} \sum_{n=1}^{N_N^u} \theta_k^a(\boldsymbol{\mu}) \mathbf{u}_n \underbrace{a_k(\boldsymbol{\xi}_j, \boldsymbol{\zeta}_n, \boldsymbol{\xi}_i)}_{\text{precomputed}}, & &= \sum_{k=1}^{Q_a} \sum_{n=1}^{N_N^T} \theta_k^a(\boldsymbol{\mu}) T_n \underbrace{a_k(\boldsymbol{\xi}_n, \boldsymbol{\zeta}_j, \boldsymbol{\xi}_i)}_{\text{precomputed}}.
\end{aligned} \tag{7.11}$$

Then we have to store  $Q_a + Q_c$  order-3 tensors of size  $N_N^u \times N_N^u \times N_N^u$  and  $N_N^u \times N_N^T \times N_N^T$ . It remains reasonable (in term of memory cost) and allows to totally ignore the finite element dimension during the assembly of the Jacobian  $\mathbf{J}_N$ .

The residual  $\mathbf{R}_N(\mathbf{u}, p, T; \boldsymbol{\mu})$  can be expressed using the same submatrices,  $\forall \boldsymbol{\mu} \in \mathcal{D}$

$$\mathbf{R}_N(\mathbf{u}, p, T; \boldsymbol{\mu}) = \begin{bmatrix} \mathbf{D}_N^u(\boldsymbol{\mu}) + \mathbf{C}_N(\mathbf{u}; \boldsymbol{\mu}) & \mathbf{B}_N(\boldsymbol{\mu})^\top & -\mathbf{F}(\boldsymbol{\mu}) \\ -\mathbf{B}_N(\boldsymbol{\mu}) & & \mathbf{A}_N(\mathbf{u}; \boldsymbol{\mu}) + \mathbf{D}_N^T(\boldsymbol{\mu}) \end{bmatrix} \begin{bmatrix} \mathbf{u} \\ p \\ T \end{bmatrix} - \begin{bmatrix} 0 \\ 0 \\ \mathbf{G}_N(\boldsymbol{\mu}) \end{bmatrix}, \tag{7.12}$$

and so we can use the method described previously and the precomputed tensors to evaluate the residual  $\mathbf{R}_N$  independently of the finite element dimension. The complete online non-linear resolution is summarized in algorithm 9.

**Remark 11.** *During the Newton algorithm, the blocks coming from the bilinear contributions are only assembled once. The blocks coming from the trilinear contributions have to be reassembled for each new linear resolution.*

The convergence of the Newton algorithm during the online procedure was an issue for some problems. We had to consider some improvements to help the convergence. First, for any new parameter  $\boldsymbol{\mu}$ , we search the closest parameter (in the sense of the usual norm for vectors) in  $S_N$ . We note  $\bar{\boldsymbol{\mu}}$  this reference parameter. Then we initialize the Newton algorithm using the projection of the solution  $(\mathbf{u}(\bar{\boldsymbol{\mu}}), p(\bar{\boldsymbol{\mu}}), T(\bar{\boldsymbol{\mu}}))$ . This projection has been evaluated and stored during the offline phase. This initialization greatly improves the convergence of the solver. However, it might not be enough. When this procedure is not sufficient, we add a continuation on the parameter. The online problem is solved for a sequence  $(\boldsymbol{\mu}_k)_0^K$  such that

$\mu_0 = \bar{\mu}$  and  $\mu_K = \mu$ . This procedure is particularly costly, but it is only required for a few numbers of parameters and is not used if the direct resolution converges.

---

**Algorithm 9:** Newton Algorithm for the online non-linear resolution

---

**Input:**  $\mu$   
**Output:**  $(\mathbf{u}_N(\mu), p_N(\mu), T_N(\mu))$

```

1  Choose  $\delta_{tol}$ ; // Algorithm tolerance
2  Choose  $(\mathbf{u}^0, p^0, T^0)$ ; // Initial guess
3  Assemble linear contribution  $G_N(\mu)$ ;
4  Assemble bilinear contributions  $D_N^u, B_N, F_N, D_N^T$ ;
5  repeat
6  | Assemble trilinear contributions:  $C_N(\mathbf{u}^{k-1}, \mu), C_N^*(\mathbf{u}^{k-1}, \mu)$ ;
7  | Assemble trilinear contributions:  $A(\mathbf{u}^{k-1}; \mu), A^*(T^{k-1}; \mu)$ ;
8  | Assemble  $J_N^{k-1} = J_N(\mathbf{u}^{k-1}, p^{k-1}, T^{k-1}, \mu)$ , see (7.8);
9  | Assemble  $R_N^{k-1} = R(\mathbf{u}^{k-1}, p^{k-1}, T^{k-1}; \mu)$ , see (7.12);
10 | Solve  $J_N^{k-1} H^{k-1} = -R_N^{k-1}$ ;
11 |  $(\mathbf{u}^k, p^k, T^k) \leftarrow (\mathbf{u}^{k-1}, p^{k-1}, T^{k-1}) + H^{k-1}$ ;
12 until  $\|R_N(\mathbf{u}^k, p^k, T^k; \mu)\|_2 \leq \delta_{tol}$ ;
13 return  $(\mathbf{u}_N(\mu), p_N(\mu), T_N(\mu)) = (\mathbf{u}^k, p^k, T^k)$ ;
```

---

### 1.3 Treatment of the Stabilization Terms

We consider the system (2.2) and we add the stabilization terms. We obtain the new weak formulation: find  $(\mathbf{u}_h(\mu), p_h(\mu), T_h(\mu)) \in V_h \times Q_h \times X_h$  such that

$$\begin{aligned} a_u^*(\mathbf{u}_h(\mu), p_h(\mu), \mathbf{v}_h; \mu) - b(\mathbf{v}_h, p_h; \mu) + b(\mathbf{u}_h, q_h; \mu) &= f_u^*(T_h, \mathbf{v}_h; \mu), \\ a_T^*(T_h, \mathbf{u}_h, S_h; \mu) &= f_T^*(S_h; \mu), \\ \forall (\mathbf{v}_h, q_h, S_h) &\in V_h \times Q_h \times X_h, \end{aligned} \quad (7.13)$$

whith

$$\begin{aligned} a_u^*(\mathbf{u}_h, p_h, \mathbf{v}_h; \mu) &= c(\mathbf{u}_h, \mathbf{u}_h, \mathbf{v}_h; \mu) + d_u(\mathbf{u}_h, \mathbf{v}_h; \mu) + a_u^{\text{SDM}}(\mathbf{u}_h, p_h, \mathbf{v}_h; \mathbf{u}_h, \mu) \\ a_T^*(T_h, \mathbf{u}_h, S_h; \mu) &= a(T_h, \mathbf{u}_h, S_h; \mu) + d_T(T_h, S_h; \mu) + a_T^{\text{SDM}}(T_h, S_h; \mathbf{u}_h, \mu) \\ f_u^*(T_h, \mathbf{v}_h; \mu) &= f(T_h, \mathbf{v}_h; \mu) + f_u^{\text{SDM}}(T_h, \mathbf{v}_h; \mathbf{u}_h, \mu) \\ f_T^*(S; \mu) &= g(S; \mu) + f_T^{\text{SDM}}(T_h, S_h; \mathbf{u}_h, \mu), \end{aligned} \quad (7.14)$$

where the forms  $c, d_u, a, d_T, f$  and  $g$  are defined in (7.3). The operators  $a_u^{\text{SDM}}, a_T^{\text{SDM}}, f_u^{\text{SDM}}$  and  $f_T^{\text{SDM}}$  are the stabilization terms detailed in the chapter 3.

Those stabilization operators are non-linear and non-affine in parameters. By introducing these terms in the variational formulation, we can rewrite the Jacobian  $J_N$  and the residual  $R_N$  as

$$\begin{aligned} J_N(\mathbf{u}, p, T; \mu) &= J_N^{\text{tri}}(\mathbf{u}, p, T; \mu) + J_N^{\text{SDM}}(\mathbf{u}, p, T; \mu), \\ R_N(\mathbf{u}, p, T; \mu) &= R_N^{\text{tri}}(\mathbf{u}, p, T; \mu) + R_N^{\text{SDM}}(\mathbf{u}, p, T; \mu), \end{aligned} \quad (7.15)$$

where  $\mathbf{J}_N^{\text{tri}}$  and  $\mathbf{R}_N^{\text{tri}}$  are the contributions detailed respectively in (7.8) and (7.12).  $\mathbf{J}_N^{\text{SDM}}$  and  $\mathbf{R}_N^{\text{SDM}}$  are the contributions from the stabilization operators. Using the EIM for discrete operators introduced in the previous chapter, we can compute an affine approximation for these contributions,

$$\begin{aligned}\mathbf{J}_N^{\text{SDM}}(\mathbf{u}, p, T; \boldsymbol{\mu}) &\approx \sum_{m=1}^{M^J} \theta_m^J(\mathbf{u}, p, T, \boldsymbol{\mu}) \mathbf{J}_{m,N} \\ \mathbf{R}_N^{\text{SDM}}(\mathbf{u}, p, T; \boldsymbol{\mu}) &\approx \sum_{m=1}^{M^R} \theta_m^R(\mathbf{u}, p, T, \boldsymbol{\mu}) \mathbf{R}_{m,N}\end{aligned}\tag{7.16}$$

where the matrices  $\mathbf{J}_{m,N}$  and the vectors  $\mathbf{R}_{m,N}$  are pre-computed during the offline procedure. The online phase then consists in the evaluation of the  $\theta_m^J(\mathbf{u}, p, T, \boldsymbol{\mu})$  and  $\theta_m^R(\mathbf{u}, p, T, \boldsymbol{\mu})$  using the EIM algorithm. This allows to compute the Jacobian matrix and the Residual vector independently of the FE dimension.

**Remark 12.** *With very basic stabilization parameters we might retrieve a quadratic non-linearity for the stabilization terms. However, it is a particular case, and we will not detail this point. In such a situation, the stabilization terms can be integrated into the trilinear contribution of the aerothermal system.*

## 1.4 Reduction of a Turbulent Model

We had some difficulties with turbulent simulations and especially with the stabilization of the turbulent models. The final development of the full turbulent model is very recent, and then we had no time to start the implementation of the reduced version. However, we will present in this section how we planned to apply the RBM to RANS simulations. The presented method is only a theory and was not implemented.

We assume here that we have a fully operational turbulence model. This model allows to compute the turbulent viscosity  $\nu_t(\mathbf{u}, \boldsymbol{\mu})$ . In RANS models, this turbulent viscosity is usually computed from scalar quantities, solutions of one (or several) PDE(s). For instance, the Spalart Allmaras model and the  $k - \omega$  SST models presented in the chapter 4 require the resolution of one and two equations respectively.

A natural idea is to use the EIM to build an affine approximation of the turbulent viscosity  $\nu_t(\mathbf{u}, \boldsymbol{\mu})$ . However, the EIM requires to know for any new  $\mathbf{u}$  and  $\boldsymbol{\mu}$ , the exact value of  $\nu_t(\mathbf{u}, \boldsymbol{\mu})$  on the interpolation points. To obtain these values, we have to build reduced basis spaces for the scalar quantities of the turbulence model ( $\bar{\nu}(\mathbf{u}, \boldsymbol{\mu})$  for Spalart Allmaras or  $(k(\mathbf{u}, \boldsymbol{\mu}), \omega(\mathbf{u}, \boldsymbol{\mu}))$  for  $k - \omega$  SST). To simplify the explanations in the following paragraphs, we suppose now that the viscosity  $\nu_t$  is computed from one generic scalar quantity  $\bar{\nu}(\mathbf{u}, \boldsymbol{\mu})$ , solution of one PDE. The following method is easily adaptable in the case of a model with multiple PDE.

In most of the RANS models, the scalar field  $\bar{v}(\mathbf{u}, \boldsymbol{\mu})$  is expressed as the solution of a non-linear transport equation on the domain  $\Omega$ . We will solve this PDE in a adapted FE space  $P_h(\Omega)$ . The generic weak formulation for this problem can be written as find  $\bar{v}(\mathbf{u}, \boldsymbol{\mu}) \in P_h, \forall \mathbf{u}_h \in V_h, \forall \boldsymbol{\mu} \in \mathcal{D}$ , such that

$$c_t(\bar{v}_h, \gamma_h, \mathbf{u}_h; \boldsymbol{\mu}) + a_t(\bar{v}_h, \gamma_h; \boldsymbol{\mu}, \mathbf{u}_h) = f_t(\gamma_h), \quad \forall \gamma_h \in P_h, \quad (7.17)$$

where,  $\forall \boldsymbol{\mu} \in \mathcal{D}$

- $c_t : P_h \times P_h \times V_h \rightarrow \mathbb{R}$  is a continuous trilinear form,
- $a_t : P_h \times P_h \rightarrow \mathbb{R}$  is a continuous bilinear form,
- $f_t : P_h \rightarrow \mathbb{R}$  is a continuous linear form.

In practice,  $a_t$  is usually not bilinear and we suppose here that we use the EIM to obtain an affine approximation,  $\forall \boldsymbol{\mu} \in \mathcal{D}, \forall \mathbf{u}_h \in V_h$

$$\begin{aligned} a_t(\bar{v}_h, \gamma_h; \boldsymbol{\mu}, \mathbf{u}_h, \bar{v}_h) &= \sum_{k=1}^{Q_{a_t}} \theta_k^{a_t}(\boldsymbol{\mu}, \mathbf{u}_h, \bar{v}_h) a_{t,k}(\bar{v}_h, \gamma_h) \\ f_t(\gamma_h; \boldsymbol{\mu}) &= \sum_{k=1}^{Q_{f_t}} \theta_k^{f_t}(\boldsymbol{\mu}) f_{t,k}(\gamma_h). \end{aligned} \quad (7.18)$$

We also suppose that the convection term  $c_t$  admits an affine decomposition

$$c_t(\bar{v}_h, \gamma_h, \mathbf{u}_h; \boldsymbol{\mu}) = \sum_{k=1}^{Q_{c_t}} \theta_k^{c_t}(\boldsymbol{\mu}) a_{t,k}(\bar{v}_h, \gamma_h, \mathbf{u}_h). \quad (7.19)$$

During the offline procedure, the reduced spaces  $V_N, Q_N, X_N$  are built as described in the previous section. We also construct the space  $P_N$ ,

$$P_N = \text{span}\{\bar{v}_h(\mathbf{u}_h(\boldsymbol{\mu}_n), \boldsymbol{\mu}_n), 1 \leq n \leq N\} = \text{span}\{\psi_n, 1 \leq n \leq N\}, \quad (7.20)$$

where  $(\psi_n)_{n=1}^N$  is an orthonormalized basis of  $P_N$ . During this phase, we will also compute the quantities

$$\begin{aligned} c_{t,k}(\psi_i, \psi_j, \boldsymbol{\zeta}_l), \quad &\forall 1 \leq i, j \leq N, \quad \forall 1 \leq l \leq N^u, \\ a_{t,k}(\psi_i, \psi_j), \quad &\forall 1 \leq i, j \leq N, \\ f_{t,k}(\psi_i), \quad &\forall 1 \leq i \leq N. \end{aligned} \quad (7.21)$$

Then, during the online phase, we can evaluate the RB approximation  $\bar{v}_N(\mathbf{u}_N, \boldsymbol{\mu})$ ,  $\forall \mathbf{u}_N \in \mathbb{R}^{N^u}, \forall \boldsymbol{\mu} \in \mathcal{D}$  as the solution of the reduced non-linear system

$$(\mathbf{C}_t(\mathbf{u}_N, \boldsymbol{\mu}) + \mathbf{A}(\boldsymbol{\mu}, \mathbf{u}_h, \bar{v}_N)) \bar{v}_N = \mathbf{F}_t(\boldsymbol{\mu}), \quad (7.22)$$

with

$$\begin{aligned}
 (\mathbf{C}_t(\mathbf{u}_N, \boldsymbol{\mu}))_{i,j} &= \sum_{k=1}^{Q_{c_t}} \sum_{l=1}^{N^u} \theta_k^{c_t}(\boldsymbol{\mu}) c_{t,k}(\psi_i, \psi_j, \boldsymbol{\zeta}_l), \\
 (\mathbf{A}_t(\boldsymbol{\mu}, \mathbf{u}_N, \bar{\nu}_N))_{i,j} &= \sum_{k=1}^{Q_{a_t}} \theta_k^{a_t}(\boldsymbol{\mu}, \mathbf{u}_N, \bar{\nu}_N) a_{t,k}(\psi_i, \psi_j) \\
 (\mathbf{F}_t(\boldsymbol{\mu}))_i &= \sum_{k=1}^{Q_{f_t}} \theta_k^{f_t}(\boldsymbol{\mu}) f_{t,k}(\psi_i).
 \end{aligned} \tag{7.23}$$

Once this system is solved, the RB approximation  $\bar{\nu}_N(\mathbf{u}_N, \boldsymbol{\mu})$  can be used to evaluate the turbulent viscosity  $\nu_t(\mathbf{u}_N, \boldsymbol{\mu})$ . It is then possible to iterate between the two reduced models with classical fixed point iterations.

The method described might theoretically work. However, it implies the use of multiple EIMs. The resolution of the turbulent PDE is usually not trivial and very sensitive. There is no guaranty that the many approximations made through the EIMs provide a stable solution.

## 2 OPUS Test-Case: Application to the Cooling of Electronic Components

We present in this section our numerical results on a test case proposed in 2009 by Michel Fouquembergh and Annabelle Le-Hyariac, both from Airbus Group.

### 2.1 Problem Description

We consider a 2D model representative of the neighboring of an electronic component submitted to a cooling air flow. It is described by four geometrical domains in  $\mathbb{R}^2$  named  $\Omega_i$ ,  $i = 1, 2, 3, 4$ , see figure 7.1. We then want to compute the temperature field  $T$ , in the whole domain  $\Omega = \cup_{i=1}^4 \Omega_i$ , and the velocity/pressure  $(\mathbf{u}, p)$  in the sub-domain  $\Omega_4$ , as the solution of the aerothermal system of equations

$$\begin{aligned}
 \rho \mathbf{u} \cdot \nabla \mathbf{u} - 2\nabla \cdot (\mu \mathbf{S}(\mathbf{u})) + \nabla p &= 0, & \text{on } \Omega_4 \\
 \nabla \cdot \mathbf{u} &= 0, & \text{on } \Omega_4 \\
 \rho_i C_i \mathbf{v} \cdot \nabla T - \nabla \cdot (k_i \nabla T) &= Q_i, & \text{on } \Omega_i, \quad i = 1, 2, 3, 4
 \end{aligned} \tag{7.24}$$

where, in each subdomain  $\Omega_i$ ,  $k_i$  is the thermal diffusivity and  $Q_i$  is a volumic heat dissipated.  $\rho$  and  $\mu$  are respectively the density and the viscosity of the air. We do not consider natural convection in this case and so the coupling between the fluid and the thermal equations is only one way.

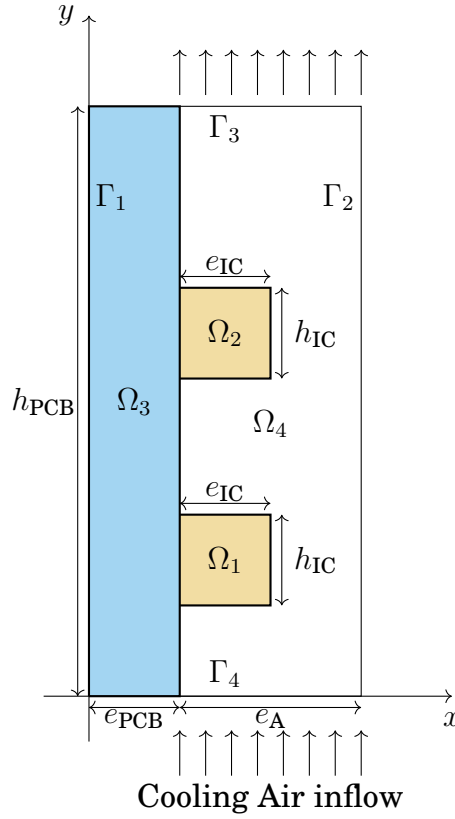


Figure 7.1 – Geometry of the problem

The subdomains  $\Omega_1$  and  $\Omega_2$  define two Integrated Circuits, respectively IC1 and IC2, simplified as two rectangles of height  $h_{IC}$  and width  $e_{IC}$ . These IC are soldered to the Printed Circuit Board (PCB) which is a rectangle of height  $h_{PCB}$  and width  $e_{PCB}$ . The Air is flowing along the PCB in the domain  $\Omega_4$ , of width  $e_A$ .

The notation for the external boundaries of the domain  $\Omega$ :

$$\begin{aligned}
 \Gamma_1 &= \{ x = 0, \quad 0 \leq y \leq h_{PCB} \}, \\
 \Gamma_2 &= \{ x = e_{PCB} + e_A, \quad 0 \leq y \leq h_{PCB} \}, \\
 \Gamma_3 &= \{ 0 \leq x \leq e_{PCB} + e_A, \quad y = h_{PCB} \}, \\
 \Gamma_4 &= \{ 0 \leq x \leq e_{PCB} + e_A, \quad y = 0 \},
 \end{aligned} \tag{7.25}$$

we also define the convenient notations:  $\Gamma_I = (\Omega_3 \cup \Omega_2 \cup \Omega_1) \cap \Omega_4$  for the internal boundary between fluid and solid domains and  $\Gamma_{ij} = \Gamma_i \cup \Gamma_j$  for the sub-boundaries.

### Boundary Conditions on the Temperature

- on  $\Gamma_1 \cup \Gamma_2 \cup \Gamma_3$ , a zero flux condition (Neumann)

$$-k \nabla T \cdot \mathbf{n} = 0, \tag{7.26}$$

- on  $\Gamma_4$ , imposed temperature (Dirichlet)

$$T = T_0; \quad (7.27)$$

- on all the internal boundaries, the continuity of the heat flux and temperature

$$T_i = T_j, \quad k_i \nabla T \cdot \mathbf{n} = -k_j \nabla T \cdot \mathbf{n} \quad (7.28)$$

### Boundary Conditions on the Fluid

- on  $\Gamma_3$ , a free stream condition

$$\mu \nabla \mathbf{u} \cdot \mathbf{n} - p \mathbf{n} = 0 \quad (7.29)$$

- on  $\Gamma_2 \cup \Gamma_I$ , a no-slip wall condition (Homogeneous Dirichlet)

$$\mathbf{u} = 0 \quad (7.30)$$

- on  $\Gamma_4$ , a Poiseuille Profile (Dirichlet)

$$\mathbf{u} = \mathbf{u}_P = \left( 0, \frac{3}{2(e_A - e_{IC})} D \left( 1 - \left( \frac{x - \left( \frac{e_A + e_{IC}}{2} + e_{PCB} \right)}{\frac{e_A - e_{IC}}{2}} \right)^2 \right) \right) \quad (7.31)$$

### Modifications of the Original Model

The boundary conditions presented above are slightly simplified compared to the conditions proposed in the original test case.

- The conditions of the temperature on  $\Gamma_1$  and  $\Gamma_2$  are supposed to be periodic conditions

$$\begin{aligned} T|_{x=0} &= T|_{x=e_{PCB}+e_A} \\ k_3 \nabla T|_{x=0} \cdot \mathbf{n}_3 &= -k_4 \nabla T|_{x=e_{PCB}+e_A} \cdot \mathbf{n}_4 \end{aligned} \quad (7.32)$$

- The conditions at the interfaces between  $IC_i$  and the PCB are supposed to be a thermal contact conductance

$$-k_i \nabla T \cdot \mathbf{n}_i - k_3 \nabla T \cdot \mathbf{n}_3 = r_{i3} (T_{\partial\Omega_i} - T_{\partial\Omega_3}) \quad (7.33)$$

where  $r_{i3}$  is a positive coefficient.

These conditions have been simplified into the respective conditions (96) and (98) because these kinds of conditions are not supported in Feel++.



### Inputs of the Model

We present in table 7.1 the different parameters of this case. We indicate a nominal value and eventually a range. The nominal value is the default. The variables with a range are the parameters of the model.

Name	Description	Nominal Value	Range	Units
IC				
$\rho_1 C_1 = \rho_2 C_2$	Heat Capacity	$1.4 \cdot 10^6$		$\text{J.m}^{-1}.\text{K}^{-1}$
$Q_1, Q_2$	Heat source	$10^6$	$[0, 10^6]$	$\text{W.m}^{-3}$
$k_1, k_2$	Thermal conductivity	2	$[0.2, 150]$	$\text{W.m}^{-1}.\text{K}^{-1}$
$e_{\text{IC}}$	Thickness	$2 \cdot 10^{-3}$		m
$h_{\text{IC}}$	Height	$70 \cdot 10^{-3}$		m
PCB				
$\rho_3 C_3$	Heat Capacity	$2 \cdot 10^6$		$\text{J.m}^{-1}.\text{K}^{-1}$
$Q_3$	Heat source	0		$\text{W.m}^{-3}$
$k_3$	Thermal conductivity	0.2		$\text{W.m}^{-1}.\text{K}^{-1}$
$e_{\text{PCB}}$	Thickness	$2 \cdot 10^{-3}$		m
$h_{\text{PCB}}$	Height	$130 \cdot 10^{-3}$		m
Air				
$\rho_4 C_4$	Heat Capacity	1100		$\text{J.m}^{-1}.\text{K}^{-1}$
$Q_4$	Heat source	0		$\text{W.m}^{-3}$
$k_4$	Thermal conductivity	$3 \cdot 10^{-2}$		$\text{W.m}^{-1}.\text{K}^{-1}$
$\rho$	Density	1.0		$\text{kg.m}^{-3}$
$\mu$	Viscosity	$1.8 \cdot 10^{-5}$		$\text{kg.m}^{-1}.\text{s}^{-1}$
$e_{\text{A}}$	Thickness	$14 \cdot 10^{-3}$		m
$T_0$	Inflow temperature	300		K
$D$	Inflow rate	7	$[0.5, 50]$	$10^{-3} \text{ m}^2.\text{s}^{-1}$

Table 7.1 – Opus Application: Parameters of the model

**Remark 13.** A geometric parameter was initially introduced in the benchmark configuration. The thickness  $e_{\text{A}}$  was supposed to vary. It appeared that some shapes produced recirculations at the outlet and then the Newton solver was not able to converge anymore. We are still working on this issue, and we hope to propose the results with geometric parameter very soon.

### Output of the Model

We compute,  $\forall \mu \in \mathcal{D}$  the mean temperature on the second IC,

$$s(\mu) = \int_{\Omega_2} T(\mu) d\Omega. \quad (7.34)$$

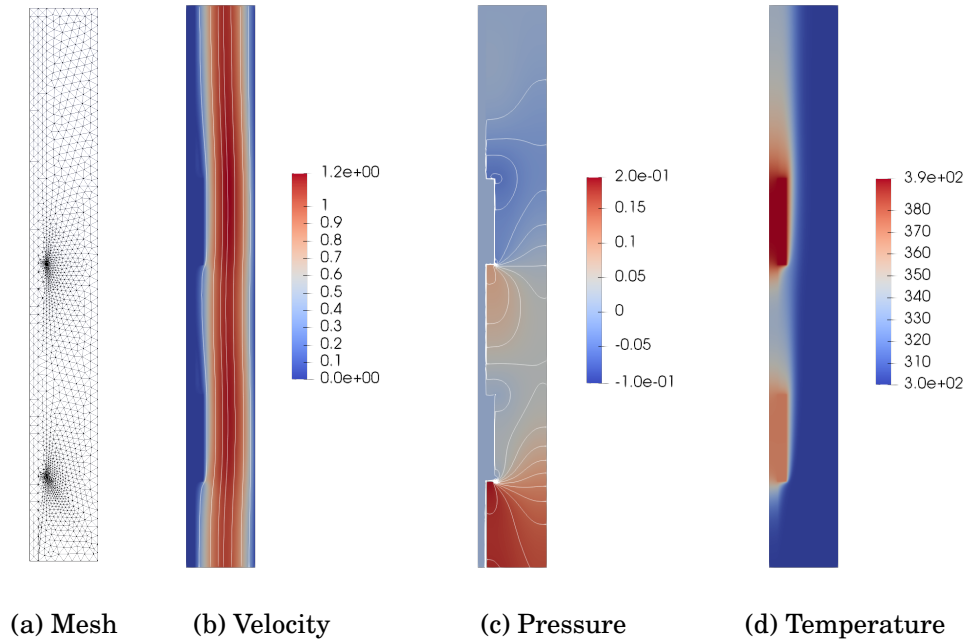


Figure 7.2 – Opus Application: mesh and example of solution profiles (RB) for parameter  $\mu = (150, 150, 10^6, 10^6, 10^{-2})$ . (a) coarser mesh used for the stabilized simulations, (b) velocity profile, in  $\text{m}\cdot\text{s}^{-1}$ , with streamlines, (c) pressure field, in Pa, with contour lines, (d) temperature field, in K

## 2.2 Numerical Results

The following results have been produced on a non-structured simplex mesh. The mesh is presented in figure 7.2 with an example of profiles for the temperature, velocity and the pressure ( $\mu = (150, 150, 10^6, 10^6, 10^{-2})$ ). The FE space is  $\mathbb{T}\mathbb{H}_{ch}^1(\Omega) \times \mathbb{L}_{ch}^2(\Omega)$ .

### No Stabilization

We first present some results with no stabilization. In order to keep a stable solution, we do two small modifications. First, the mesh has been refined by splitting all the cells in the mesh presented in figure 7.2. Then, the range for the inflow rate has been reduced, and we consider for now  $D \in [0.5 \cdot 10^{-3}, 10 \cdot 10^{-3}] \text{m}^2 \cdot \text{s}^{-1}$ .

We built the reduced basis using  $\delta_{tol} = 10^{-12}$  in the selective Gram-Schmidt algorithm 7, with 50 parameters in  $S_N$ . At the end of the offline procedure, the size of the spaces  $V_N \times Q_N \times X_N$  is  $89 \times 38 \times 50$ . We present the convergence of the errors for  $s$ ,  $u$ ,  $p$  and  $T$  on the figure 7.3. These results are very satisfying, considering the dimension of the parameter space. Of course, the convergence of the errors can still be improved by using an adapted error estimator in the greedy algorithm.

As expected, the time for the resolution of the reduced problem increases with the size of the reduced spaces. We present in figure 7.4 the mean runtime for

the reduced problem vs. the size  $N$  of  $S_N$ . To explain these runtimes, we also plotted (figure 7.5) the execution times for the main components of the resolution process. We observe that the assemblies of the reduced Jacobian and Residual become very expensive when  $N$  is growing. However, we figured out that the time for these operations depends on the dimension  $N$  but also on the max dimension  $N_{\max}$ . Indeed, if we compare the assembly times for different values of  $N_{\max}$ , see figure 7.6, we see a significant difference. This gap is probably the consequence of the block extraction methods used to assemble the reduced operators. The reduced matrices and vectors are stored using data structures from the C++ library Eigen. This library is very convenient for the extraction and the operations on submatrices. However, these extraction methods seem to depend on the size of the global matrix. This issue should be investigated in future works.

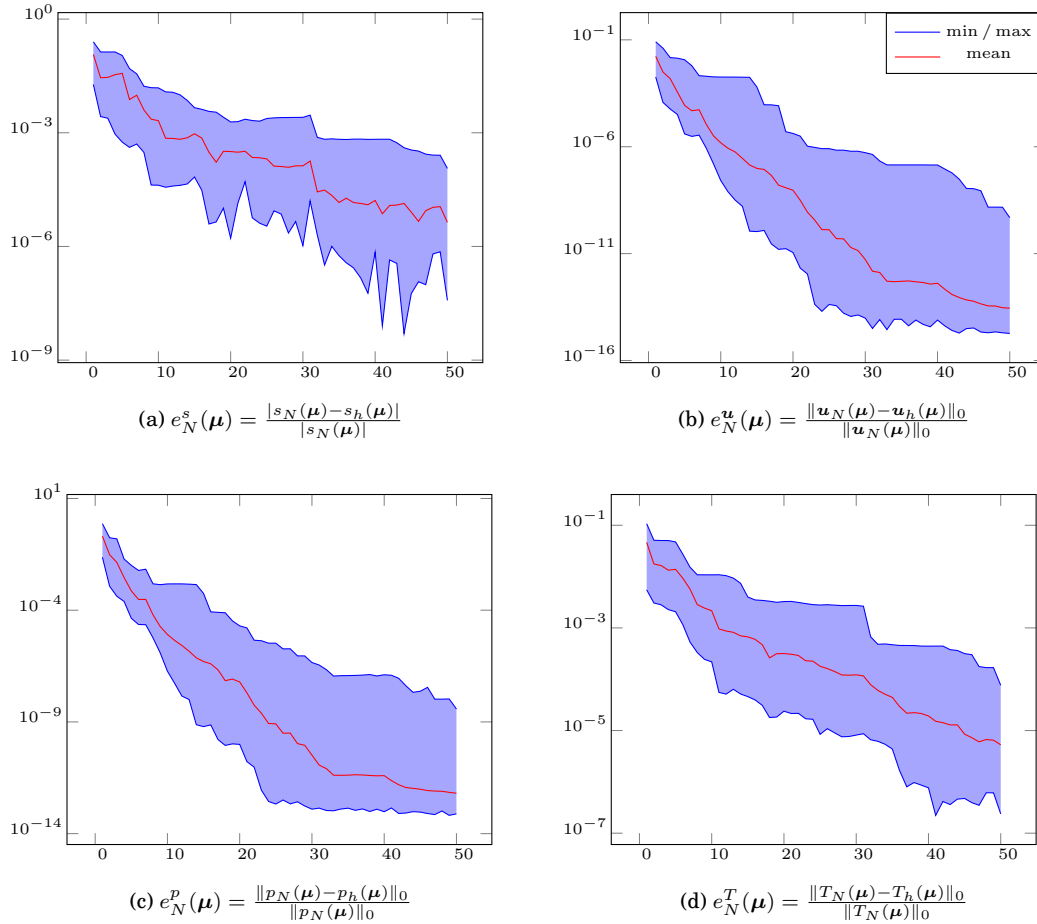


Figure 7.3 – Opus Application without Stabilization: convergence of the errors vs the size of  $N$ . Maximum minimum and average errors evaluated from 50 approximations.  $\delta_{tol} = 10^{-10}$  in selective Gram-Schmidt algorithm 7, maximum size of the reduced spaces  $89 \times 38 \times 50$

These first results with our RB aerothermal solver are very encouraging. However, the geometry was very basic, and thus there were no significant variations in the fluid flow. We are now interested in the validation of our solver on more complex problems and also for simulations with stabilization methods.

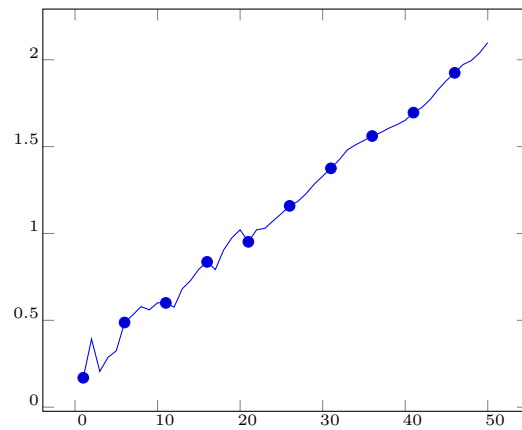


Figure 7.4 – Opus Application without Stabilization: average resolution time (in  $s$ ) for the reduced problem vs  $N$ .

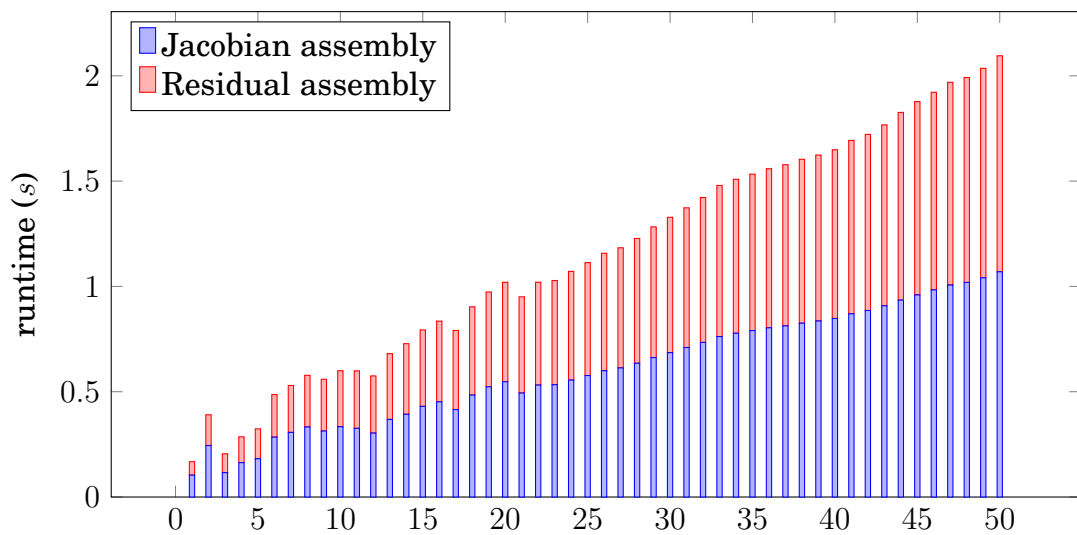


Figure 7.5 – Opus Application without Stabilization: Runtime of the main phases for the online resolution. Mean value for 50 runs. Other contributions are negligible

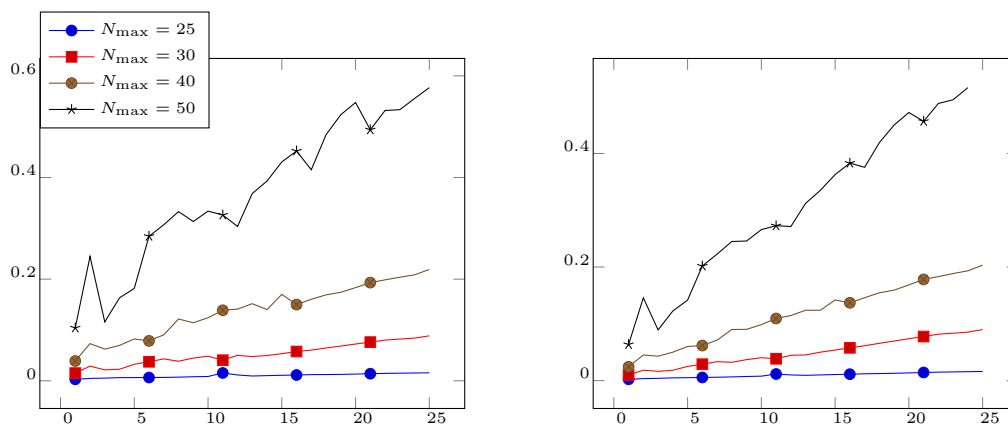


Figure 7.6 – Opus Application without Stabilization: assembly time for the reduced Jacobian (a) and the reduced Residual (b) for different values of  $N_{\max}$ . Mean values evaluated over 50 reduced resolutions.

## Results with Stabilization

We now introduce stabilization operators. For this test case we choose the formulation  $GLS_1$  for both the fluid and the temperature, see chapter 3. We use the non-refined mesh presented in figure 7.2 and the full range for the inflow rate,  $D \in [0.5 \cdot 10^{-3}, 50 \cdot 10^{-3}] \text{ m}^2 \cdot \text{s}^{-1}$ . The tolerance in the selective Gram-Schmidt algorithm is set to  $\delta_{tol} = 10^{-10}$ .

There are two possibilities for the treatment of the stabilization. We can either use an EIM for the stabilization of the whole system or two EIMs, one for each physic. We initially wanted to compare the results of both versions. However, it appears that the results with the monolithic construction of the EIM approximation was not convincing and we finally abandoned this idea. It seems necessary to build one EIM for each physic.

We present on figure 7.7 a comparison to the errors for classical EIM and SER construction, both with the same offline sampling and the same EIM size  $M = 25$ . The convergence of SER is very satisfying in comparison of the classic EIM. It confirms the preliminary results of the chapter 6. We also propose a plot of the minimal, maximal and averaged values of the errors for the construction using SER on figure 7.8. Regarding the convergence of the errors for the velocity and the pressure fields, we observe a deterioration of the results after  $N = 16$ . For now, we have no explanation for this phenomena, and we are still investigating this issue.

With these simulations, we highlighted another problem with our implementation of the discrete EIM. The computation of the EIM coefficients is relatively expensive and is limiting the performance of the reduced model. Even if the assembly is made on a minimal mesh, it requires between 0.1 and 0.3s per assembly. With the four EIM of this model and the non-linear iterations, the time spent in computing the non-linear coefficients becomes particularly expensive. To illustrate this issue, we present on figure 7.9 the runtimes of the principal components of the online resolution: the assembly of the Jacobian, the assembly of the Residual and the evaluation of the nonlinear coefficients. This excessive computation time, in addition to the runtime issue already exposed with the previous test-case, became a problem for the study of simulations with stabilization. We started thinking about some optimizations to fix this issue but it requires some implementation efforts and the modification could not be made in the context of this thesis.

Even if the numerical results are encouraging for this test-case, it remains some work to do in order to be completely satisfied. We have to identify the cause of the bad accuracy for the pressure and the velocity fields and we still have to optimize the computation of the non-linear coefficients using EIM. We will now test our framework on more complex geometries with the following problem.

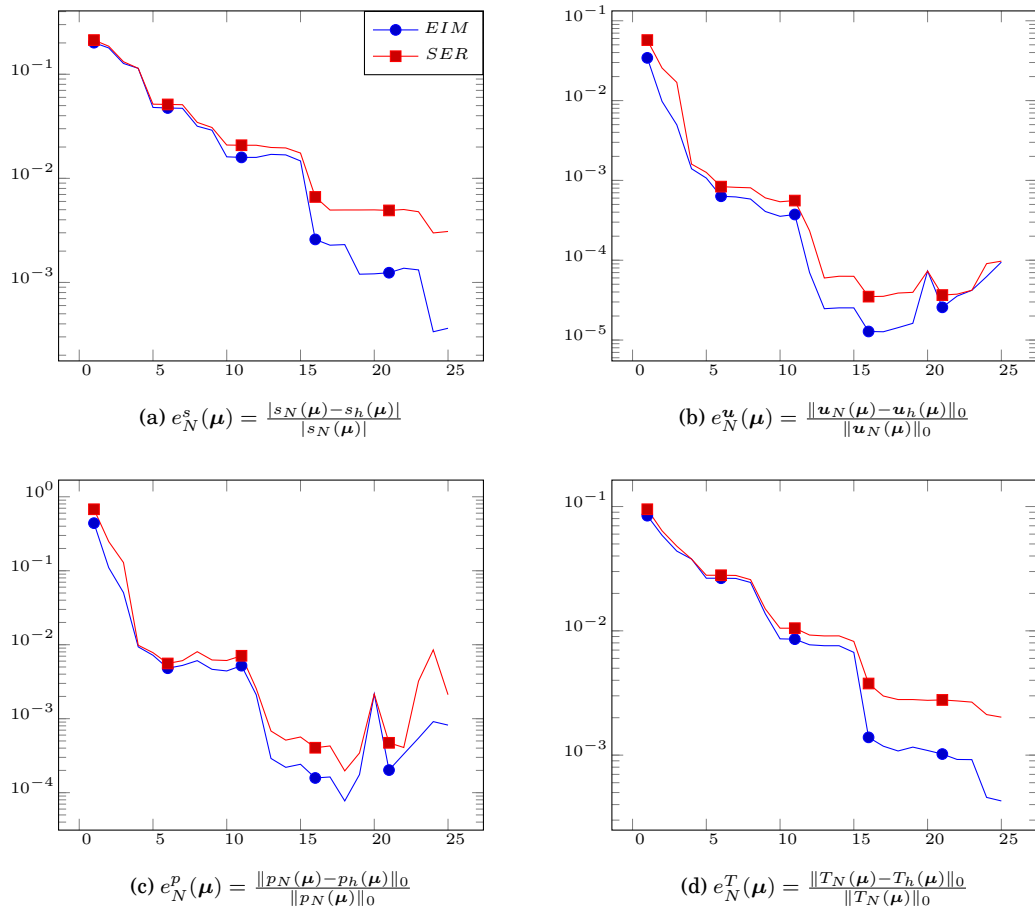


Figure 7.7 – Opus Application with Stabilization: convergence of the errors vs the size of  $N$ . Comparison of the maximum values for classical EIM and SER algorithm. Maximum size of the reduced spaces  $50 \times 25 \times 25$

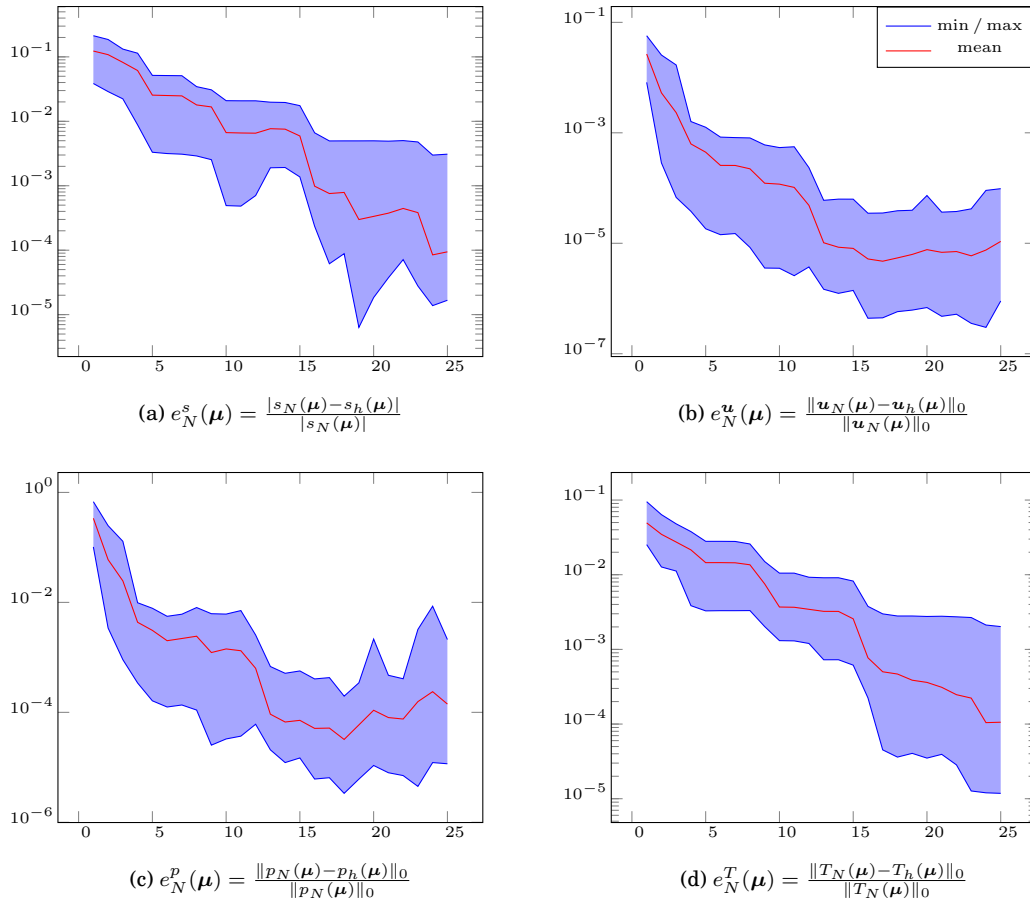


Figure 7.8 – Opus Application with Stabilization + SER: convergence of the errors vs  $N$ . Maximum minimum and average errors evaluated from 50 approximations. Maximum size of the reduced spaces:  $50 \times 25 \times 25$ . Number of basis vectors in EIM approximations  $M = 25$

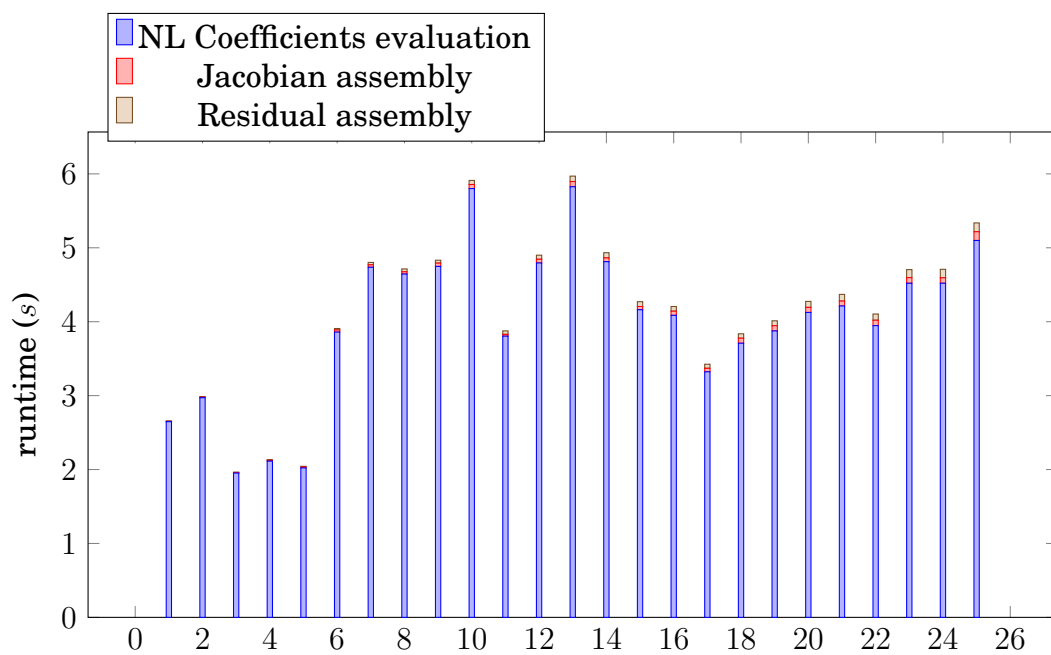


Figure 7.9 – Opus Application without Stabilization: Runtime of the main phases for the online resolution. Mean value for 50 runs. Other contributions are negligible



### 3 CHORUS Test-Case: Application to Aerothermal Simulations in an Airplane

We present in the test case proposed by Airbus Group in the context of the CHORUS project.

#### 3.1 Problem Description

The geometry of the problem is a 2D slice of an aircraft, see 7.10. The top of the geometry is the Cabin  $\Omega_C$ . In this part are six passengers  $P_i$ ,  $1 \leq i \leq 6$ , one inlet  $I_C$  for the cooled air and two outlets  $O_{C1}$  and  $O_{C2}$ . The bottom part is the bay of the aircraft. It contains two electrical equipment  $E_1$  and  $E_2$ , three bleed air ducts  $D_1$ ,  $D_2$  and  $D_3$  and three electrical harnesses  $H_1$ ,  $H_2$  and  $H_3$ . There are also two inlets  $I_{B1}$  and  $I_{B2}$  and two outlets  $O_{B1}$  and  $O_{B2}$ . For any of these items  $X$  we denote by  $\Gamma_X$  its boundary. All the remaining boundaries are denoted by  $\Gamma_W$ . We finally note  $\Omega = \Omega_C \cup \Omega_B$ .

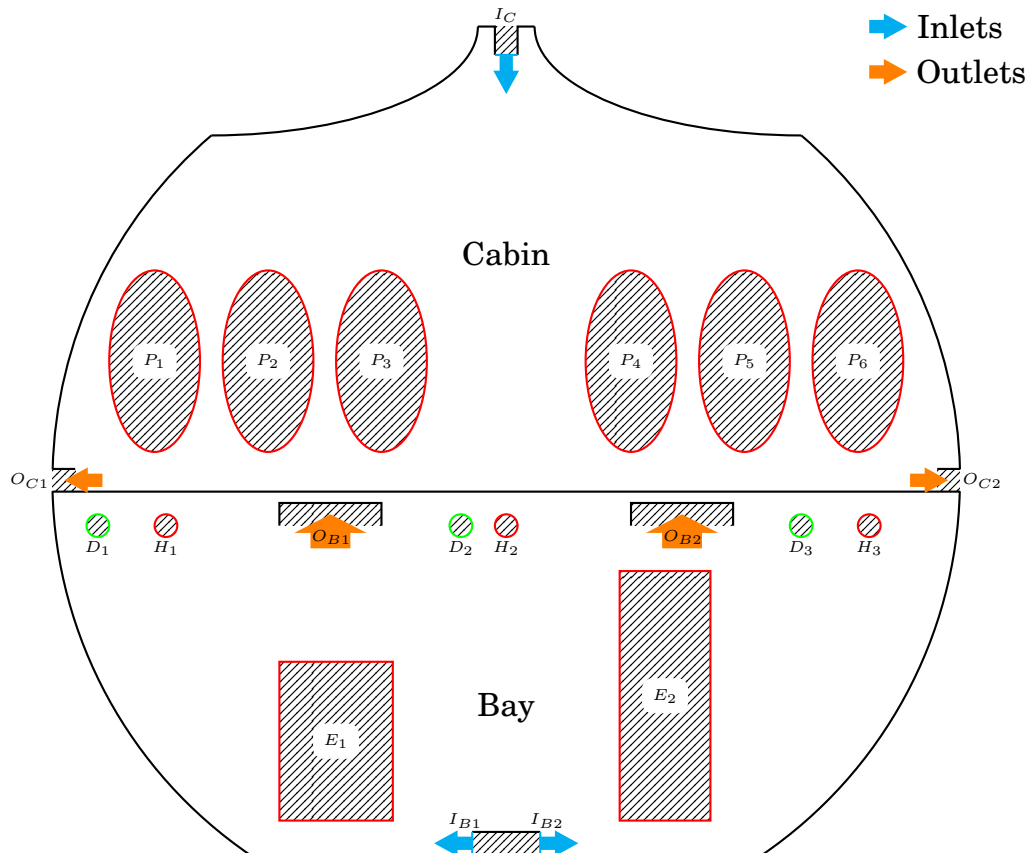


Figure 7.10 – Isometric 2D slice of the avionic bay

The aim of this test-case is the computation of the velocity field  $u$ , the pressure

field  $p$  and the temperature field  $T$ , as solutions of the strong problem

$$\begin{aligned} \rho \mathbf{u} \cdot \nabla \mathbf{u} - 2\nabla \cdot (\mu \mathbf{S}(\mathbf{u})) + \nabla p &= -\rho\beta(T - T_0)\mathbf{g}, \\ \nabla \cdot \mathbf{u} &= 0, \\ \mathbf{v} \cdot \nabla T - \nabla \cdot (\kappa \nabla T) &= 0, \end{aligned} \quad \text{on } \Omega. \quad (7.35)$$

The closing boundary conditions and the parameters of this system are detailed in the following paragraphs.

### Boundary Conditions on the Temperature

- **Passengers:** on  $\Gamma_P = \cup_{i=1}^6 \Gamma_{P_i}$ , a constant heat flux (Neumann condition)

$$-\kappa \nabla T \cdot \mathbf{n} = Q_{P_i}, \quad (7.36)$$

- **Inlets:** on  $\Gamma_{I_k}, k = I_C, I_{B1}, I_{B2}$ , a constant temperature (Dirichlet condition)

$$T = T_{I_k} \quad (7.37)$$

- **Ducts:** on  $\Gamma_{D_i}, 1 \leq i \leq 3$ , a constant temperature (Dirichlet condition)

$$T = T_{D_i} \quad (7.38)$$

- **Harnesses:** on  $\Gamma_{H_i}, 1 \leq i \leq 3$ , a constant flux (Neumann condition)

$$-\kappa \nabla T \cdot \mathbf{n} = Q_{H_i}, \quad (7.39)$$

- **Equipment:** on  $\Gamma_{E_i}, 1 \leq i \leq 2$ , a constant flux (Neumann condition), the same for both equipment

$$-\kappa \nabla T \cdot \mathbf{n} = Q_E, \quad (7.40)$$

- **All the other solid walls are supposed to be insulated (Homogeneous Neumann condition)**

$$\nabla T \cdot \mathbf{n} = 0 \quad (7.41)$$

### Boundary Conditions on the Fluid

- **Cabin Inlet:** on  $\Gamma_{I_C}$ , a Poiseuille profile (Dirichlet)

$$\mathbf{u} = \mathbf{u}_{I_C} = (0, -D_{I_C} 600(0.05 - x)(x + 0.05)) \quad (7.42)$$

- **Bay Inlets:** on  $\Gamma_{I_{B1}}$  and  $\Gamma_{I_{B2}}$ , a Poiseuille profile (Dirichlet)

$$\begin{aligned} \mathbf{u} = \mathbf{u}_{I_{B1}} &= -(D_{I_B} 600(y + 1.65)(-1.55 - y), 0), & \text{on } \Gamma_{I_{B1}} \\ \mathbf{u} = \mathbf{u}_{I_{B2}} &= (D_{I_B} 600(y + 1.65)(-1.55 - y), 0), & \text{on } \Gamma_{I_{B2}} \end{aligned} \quad (7.43)$$

- Outlets: on  $\Gamma_{Ok}$ ,  $k = C1, C2, B1, B2$ , a free stream condition

$$\mu \nabla \mathbf{u} \cdot \mathbf{n} - p \mathbf{n} = 0 \quad (7.44)$$

- All the other solid walls have no-slip conditions (Homogeneous Dirichlet)

$$\mathbf{u} = (0, 0) \quad (7.45)$$

### Modifications of the Original Model

The boundary conditions presented above are slightly simplified compared to the conditions proposed in the original test case.

- The floor between the two parts is not supposed to be insulated. We should model it as a solid part with a thickness of  $0.05m$ . This simplification allows decoupling the two resolutions. The modeling of the floor might be added in a future version of this application.
- The external walls of both the bay and the cabin are not supposed to be insulated. The walls should be modeled as solid parts with a thickness of  $0.1m$ . The outside temperature of the air is then a new parameter. This part might be added in a future version of the application.
- In the following numerical experiments, we will not consider the physical parameters of the air ( $\mu$  and  $\kappa$ ). Such parameters do not allow to reach a steady state without turbulence models. We will use larger viscosity and diffusivity to ensure the convergence of the solver. As soon as the turbulence models are fully operational, we will reconsider this simplification.
- The heat sources of the different passengers, harnesses and ducts are originally supposed to be independent. To keep the dimension of the parameter space reasonable, we suppose here that

$$\begin{aligned} T_{D_i} &= T_{D_j} = T_D, & \forall 1 \leq i, j \leq 3, \\ Q_{H_i} &= Q_{H_j} = Q_H, & \forall 1 \leq i, j \leq 3, \\ Q_{P_i} &= Q_{P_j} = Q_P, & \forall 1 \leq i, j \leq 6. \end{aligned} \quad (7.46)$$

### Inputs of the Model

We present in table 7.2 the different parameters of this case. We indicate a nominal value and eventually a range. The nominal value is the default. The variables with a range are the parameters of the model.

Name	Description	Nominal Value	Range	Units
<b>Cabin</b>				
$Q_P$	Passengers heat source	0.016	[0, 0.016]	K.s <sup>-1</sup>
$T_{IC}$	Inflow temperature	288	[273, 303]	K
$D_{IC}$	Inflow rate	0.5	[0.001, 1]	m <sup>2</sup> .s <sup>-1</sup>
<b>Bay</b>				
$Q_H$	Harnesses heat source	0.004	[0, 0.0045]	K.s <sup>-1</sup>
$Q_E$	Equipment heat source	0.05	[0, 0.05]	K.s <sup>-1</sup>
$T_D$	Ducts temperature	293	[293, 373]	K
$T_{IB}$	Inflow temperature	288	[273, 303]	K
$D_{IB}$	Inflow rate	0.5	[0.001, 1]	m <sup>2</sup> .s <sup>-1</sup>
<b>Air</b>				
$\kappa$	Thermal conductivity	$2.7 \cdot 10^{-3}$		m <sup>2</sup> .s <sup>-1</sup>
$\rho$	Density	1.0		kg.m <sup>-3</sup>
$\mu$	Viscosity	$1.8 \cdot 10^{-3}$		kg.m <sup>-1</sup> .s <sup>-1</sup>

Table 7.2 – Chorus Application: Parameters of the model

### Output of the Model

The original model was described with 15 quantities of interest. We chose here to focus on two of them: the mean temperature of the cabin and the mean temperature of the walls of the bay,

$$\begin{aligned}
 s_C(\boldsymbol{\mu}) &= \int_{\Omega_C} T(\boldsymbol{\mu}) \, d\Omega \\
 s_B(\boldsymbol{\mu}) &= \int_{\Gamma_{WB}} T(\boldsymbol{\mu}) \, d\Gamma
 \end{aligned}
 \tag{7.47}$$

## 3.2 Numerical Results

The following results have been produced on non-structured simplex meshes. The FE space is the same for both the cabin and the bay,  $\mathbb{TH}_{ch}^1(\Omega) \times \mathbb{L}_{ch}^2(\Omega)$ . The physical parameters presented previously guarantee to find a stationary state for the problem. It also ensures the stability of the numerical solution with a reasonable mesh. For that reason, we do not need to consider the stabilization operators for this test-case.

### Results on the Cabin

The mesh used for the simulation on the cabin part, and an example of solution profiles are proposed in figure 7.11. We built the reduced basis using  $\delta_{tol} = 10^{-10}$  in the selective Gram-Schmidt algorithm 7, with 50 parameters in  $S_N$ . At the end of

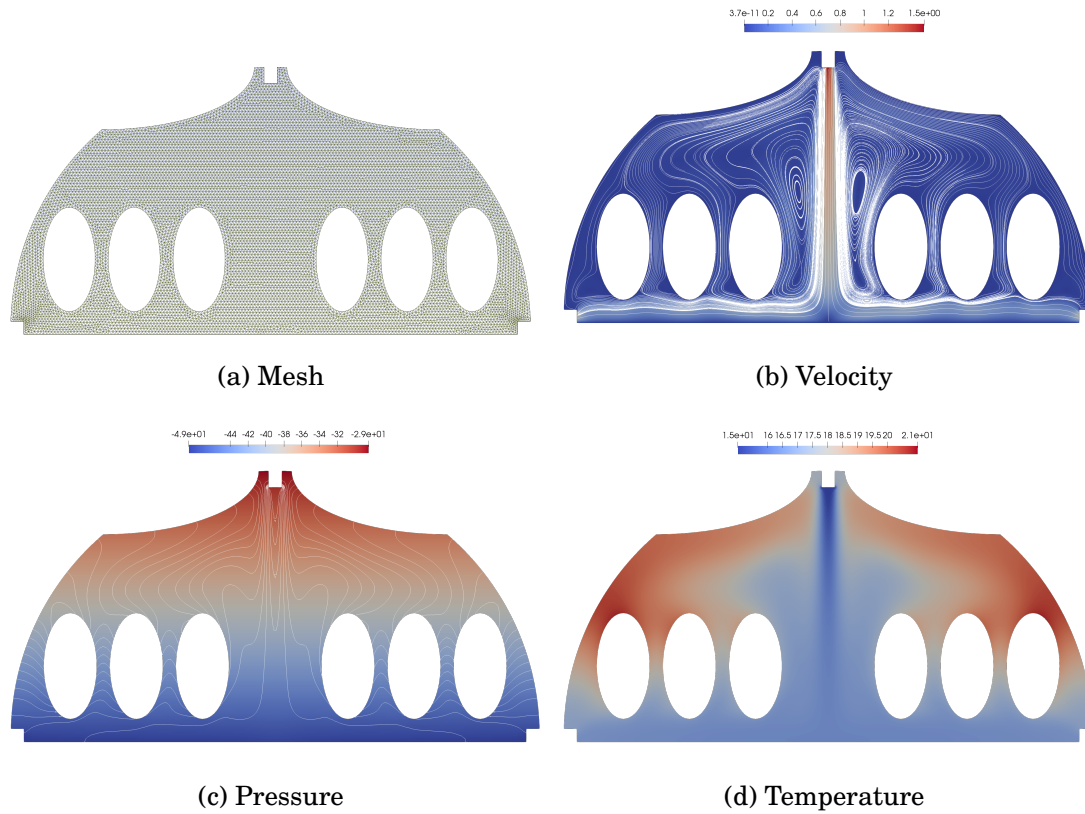


Figure 7.11 – Chorus Application, Cabin Part: (a) meshes and example of solution profiles ( $\mu = (1.6 \cdot 10^{-2}, 288, 1)$ ) (b) velocity profile, in  $\text{m.s}^{-1}$ , with streamlines, (c) pressure field, in Pa, with contour lines, (d) temperature field, in  $^{\circ}\text{C}$

the offline procedure, the size of the spaces  $V_N \times Q_N \times X_N$  is  $97 \times 47 \times 50$ . We present the convergence of the errors for  $s$ ,  $u$ ,  $p$  and  $T$  on the figure 7.12.

We observe a small deterioration of the results (for the max value) after  $N = 35$ . We have no explanation for these error pikes however these inaccurate results only happened for one parameter value. This kind of wrong approximation would probably disappear with a suitable error estimator and a thinner super sampling  $\Xi$ . Regarding the mean value of the errors, the results are very satisfying.

Again the resolution time for the online resolution grows very fast with  $N$ . For this test case, we also observed that the number of iterations in the non-linear solver is more consequent and fluctuating. We present the mean resolution times and number of iterations in figure 7.13 for  $N_{\max} = 50$  and  $N_{\max} = 25$  as a comparison. We observe the same gap between the runtime for  $N_{\max} = 25$  and  $N_{\max} = 50$ , but the number of iterations is broadly the same. The evolution of the time required for the assemblies of the reduced Jacobian and Residual is also presented on figure 7.14. Considering the fluctuating number of iterations in the non-linear solver, we decided to plot the mean value per iteration to have a better idea of the actual evolution.

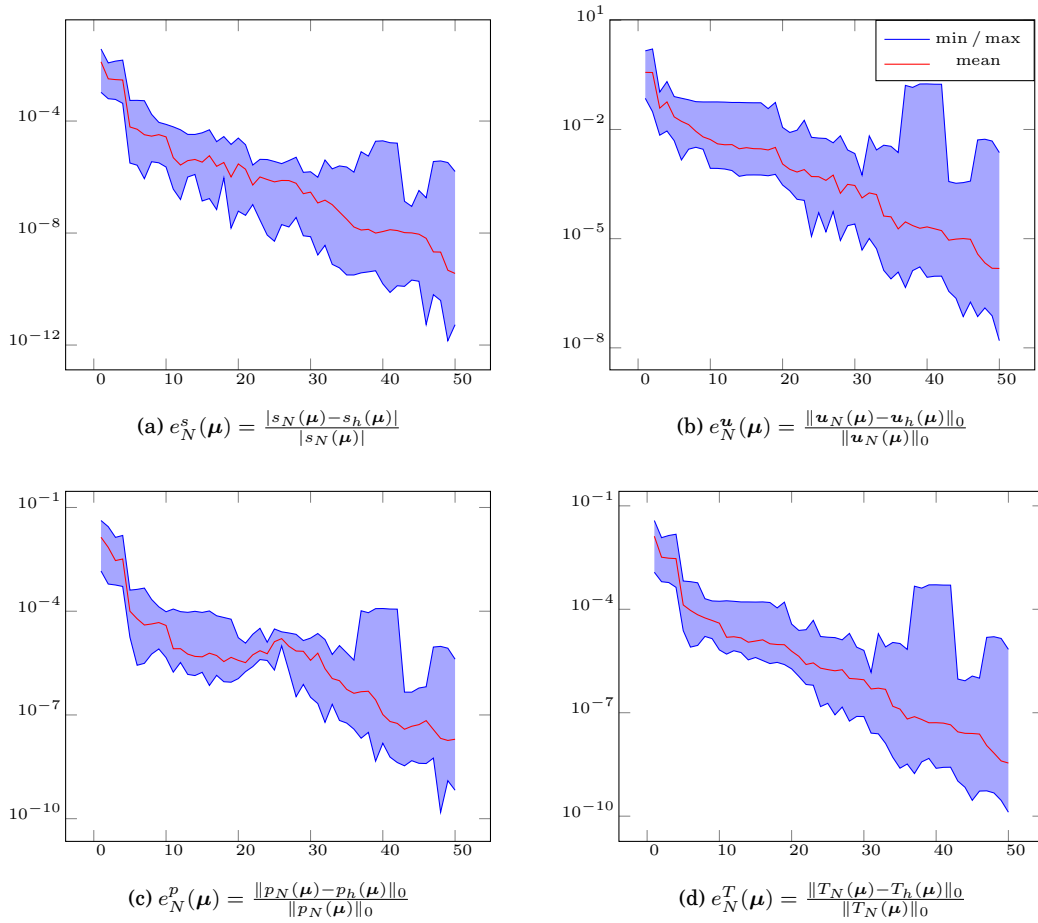


Figure 7.12 – Chorus Application, Cabin Part: convergence of the errors vs the size of  $N$ . Maximum minimum and average errors evaluated from 50 approximations.  $\delta_{tol} = 10^{-10}$  in selective Gram-Schmidt algorithm 7, maximum size of the reduced spaces  $97 \times 47 \times 50$

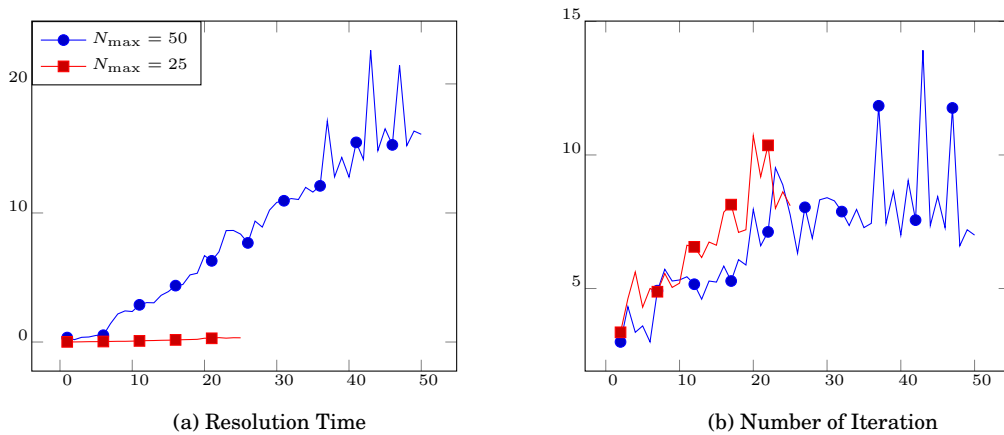


Figure 7.13 – Chorus Application, Cabin Part: average resolution time (in s) and number of iterations in the non-linear solver for the reduced problem vs  $N$ .

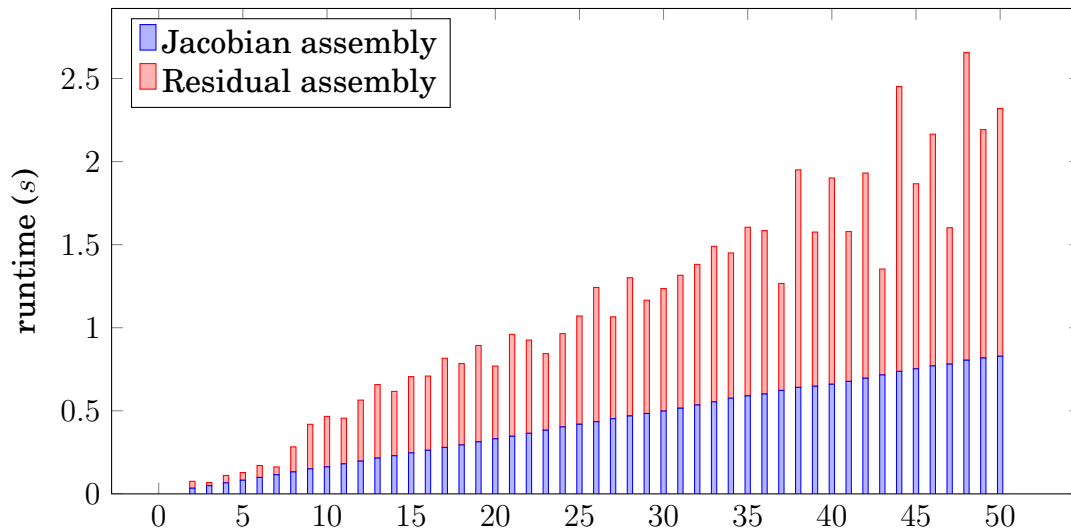


Figure 7.14 – Chorus Application, Cabin Part: Runtime of the main phases for the online resolution. Mean value for 50 runs. Other contributions are negligible. Jacobian and Residual assembly times per iteration.

## Results on the Bay

This part of the problem is more complex than the cabin. First, the dimension of the parameter space is larger, with five independent inputs. Then geometry is more complex and so is the fluid flow. The convergence of our solver was complicated even for the FEM resolution, and we had to use both physical and  $\Psi_{tc}$  continuation, see chapter 2. We present the mesh of simulation and an example of solution profiles on the figure 7.15

Regarding the convergence of the errors in figure 7.16 we observe that the reduced model is not able to correctly approximate the truth solution, even for large values of  $N$ . This bad accuracy can be explained by the complexity of the geometry and the larger dimension of the parameter space. We hope that an appropriate error estimator would improve these convergence results.

We did not plot any runtime information for this test-case, because these results are entirely similar to those of the cabin case.

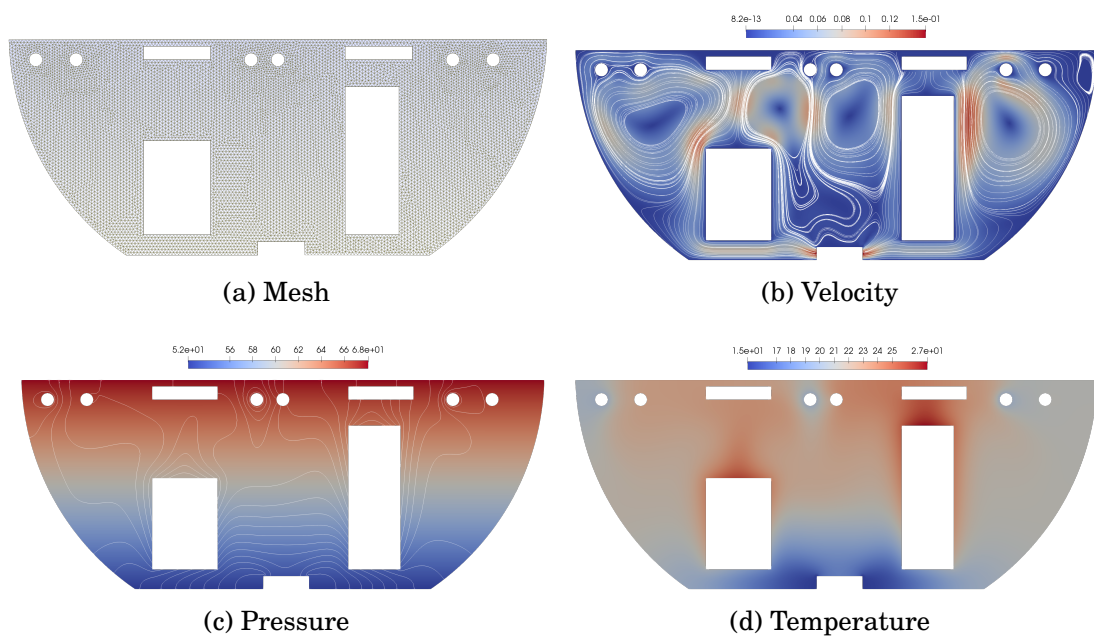


Figure 7.15 – Chorus Application, Bay Part: (a) meshes and example of solution profiles ( $\mu = (0.004, 0.05, 293, 288, 0.5)$ ) (b) velocity profile, in  $\text{m}\cdot\text{s}^{-1}$ , with streamlines, (c) pressure field, in Pa, with contour lines, (d) temperature field, in  $^{\circ}\text{C}$



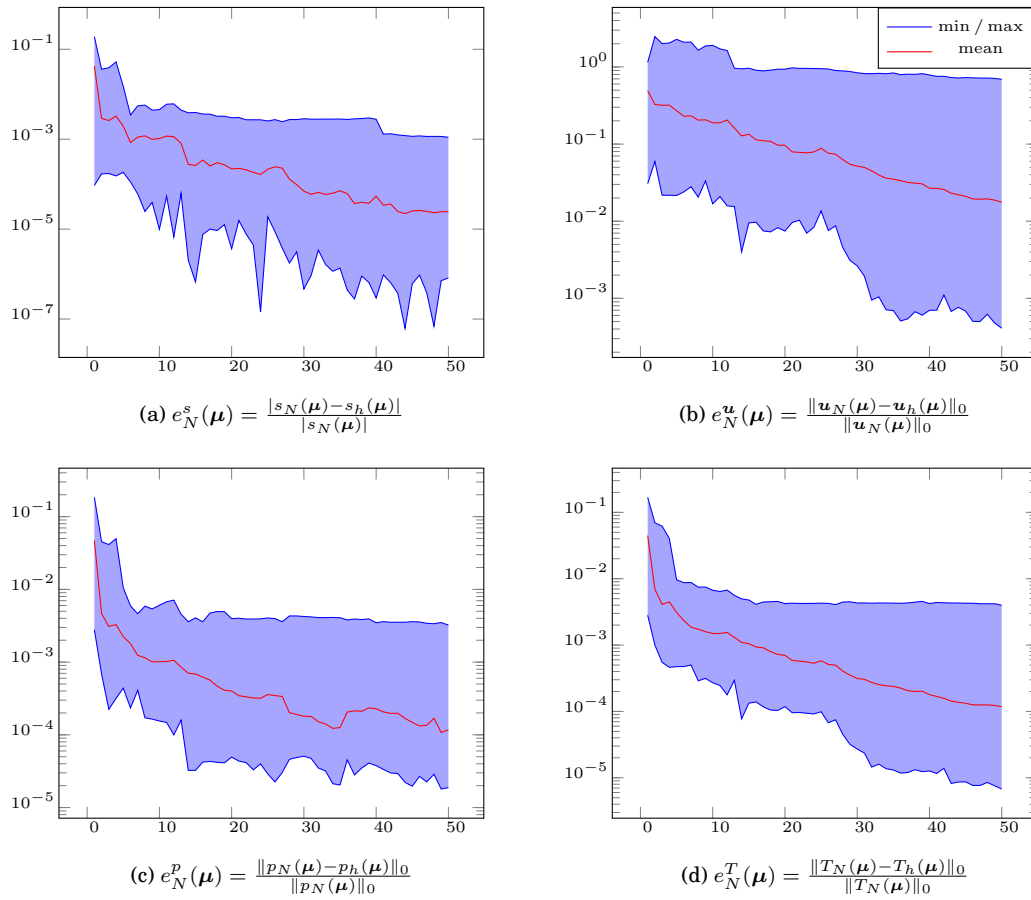


Figure 7.16 – Chorus Application, Bay Part: convergence of the errors vs the size of  $N$ . Maximum minimum and average errors evaluated from 50 approximations.  $\delta_{tol} = 10^{-10}$  in selective Gram-Schmidt algorithm 7, maximum size of the reduced spaces  $100 \times 50 \times 50$

## Conclusion

We proposed in this chapter a method for the reduction of an aerothermal model, using the Reduced Basis Method. The reduced spaces are built for each field,  $u$ ,  $p$ , and  $T$ . The non-linearity is solved using a Newton algorithm and standard reduction methods for quadratic problems. We also described a method for the reduction of stabilized operators, based on the EIM applied to discrete operators. Finally, we proposed some ideas for the reduction of turbulent problems, this last point remaining theoretical.

The results presented in this chapter are the outcome of this thesis. The resolutions of the FEM problems required most of the techniques developed in the part I (non-linear resolution, stabilization, continuations). The reduction procedure also used the different advanced methods presented in previous chapters (RB for multi-physics and saddle-point problems, EIM and SER for stabilization operators). We obtained excellent results on the opus test case and the cabin. The results for the stabilized simulation and the bay are less accurate and probably require some optimizations.

We propose here some points to study in future works to improve the reduced model:

- The error estimators are missing. We saw in the applications that the convergence of the errors might be difficult, especially when the dimension of the parameter space is growing. The development of an accurate error estimator would significantly improve this convergence. In a first time, we could imagine implementing the dual norm of the residual for aerothermal problems, as an error indicator.
- The theoretical study of the problem is missing too. For now we validated our framework with numerical experiments. It would be interesting to have also a theoretical study of the well-posedness of the reduced problem.
- The reduction of the turbulent models presented in the chapter 4.

We also highlighted two technical issues in the current implementation:

- The data structure used to store the reduced operators is not optimized. The operations on the reduced matrices and vectors become very expensive when the dimension of the reduced space is growing.
- The online computation of the EIM coefficients is still too slow. This part has already been optimized by assembling the online tensors only on the required degrees of freedom, see section 2. However, the cost of these evaluations is a limiting factor, especially when we use multiple EIM. It remains some

optimization to do, for instance, it is possible to store the evaluations of the FE basis functions for all the quadrature points of the interpolation mesh. This kind of technical improvements requires consequent implementation works and could not be achieved in the context of this thesis.

## **Part III**

# **Contribution to Feel++ Library**



---

## Aerothermal Framework

---

### Contents

---

<b>1</b>	<b>The Navier Stokes Solver . . . . .</b>	<b>164</b>
<b>2</b>	<b>The Advection Diffusion Solver . . . . .</b>	<b>165</b>
<b>3</b>	<b>Turbulence Models . . . . .</b>	<b>169</b>
3.1	The Stabilization Parameters . . . . .	170

---

The implementation of the aerothermal framework evolved a lot during the thesis. As explained in the first part of this report, we tried many strategies for the resolution of the coupled problems. When we eventually fixed the global idea of our solver, we decided to design a very generic solver to use it for many different applications. That kind of “black box” solvers already exists in Feel++ library and are called *toolboxes*. An aerothermal toolbox has been implemented since the beginning of this thesis, however, it does not have all the features required for our simulations. Thus, we started to implement our own toolbox with the advanced features presented in part I. These specific tools are now progressively adapted, optimized and integrated into the Feel++ aerothermal toolbox.

In this chapter, we present the main components of these aerothermal library. The models corresponding to these developments are detailed in the chapters 2, 3 and 4.

# 1 The Navier Stokes Solver

The resolution of the Navier-Stokes equations is managed by the class NS

---

```
template <typename meshType, int Order, int LMOOrder, int TOrder>
class NS
```

---

where

- `meshType` is the type of mesh. We can change the shape (simplex or hypercube), the dimension and the geometrical order. Eventually, we used `Mesh<Simplex<2, 1>>` in almost all our applications.
- `Order` is the polynomial order  $k$  of the Taylor Hood FE space used for the resolution :  $V_h \times Q_h = \text{TH}_{ch}^k$ , see chapter 2 for the definition.
- `LMOOrder` is the polynomial order of the Lagrange Multiplier space used to impose the slip condition when needed. This template argument was used to test the implementation, but the results are only concluding with order 0.
- `TOrder` is the polynomial order  $k$  used for the temperature field when the aerothermal option is activated,  $X_h = \mathbb{L}_{ch}^k$

This class has been designed to be very generic. Once the polynomial orders are chosen through the template parameters, we can set

- the viscosity  $\mu$  of the fluid ( as a constant or a scalar field),
- the density  $\rho$  of the fluid,
- the thermal diffusivity  $\kappa$  of the fluid,
- the thermal expansion coefficient  $\beta$ ,
- the source term of the NS equation.

Then we can choose some of the characteristics of the solver through the command line options:

- the solver type: Newton, Picard or Stokes resolution,
- the formulation (symmetric or not),
- the stabilization method used (if needed),
- the continuation method used (if needed),
- the turbulence method used (if needed).

Listing 8.1 – Minimal application solving the Turek 2D benchmark, using the NS class

---

```

1 int main( int argc, char** argv )
2 {
3     typedef Mesh<Simplex<2,1> mesh_t;           // define the mesh type
4
5     Environment env( _argc=argc, _argv=argv ); // create Feel environment
6
7     auto mesh = loadMesh( _mesh=new mesh_t );
8
9     auto ns_solver = NS<mesh_t>( mesh );       // instantiate the solver
10    ns_solver.setMu( 0.001 );                 // set the viscosity
11
12    auto U = ns_solver.solve();               // solve the system
13 }

```

---

We propose an elementary example of code in the snippet 8.1. This piece of code is a minimal application to solve the Turek benchmark 4.1. Once the `solve()` function is called (l.12), the solution  $U$  can be exported or used to compute quantities of interest.

The `solve()` method is the main function of the class. It may be called with the fixed point (Picard) solver or the Newton solver. However, the most recent features are not implemented in the Picard version and we will focus on the Newton algorithm. The global resolution routine is detailed in the algorithm 10. We also detail the crucial steps of the update of the Jacobian  $J$  and the Residual  $R$  in the respective algorithms 11 and 12. Note that most of the operators are updated with the Jacobian and reused in the Residual.

In the NS class, the number of blocks in the matrices/vectors can be set dynamically. By default, only the block velocity and pressure are built. If the problem presents slip conditions, we add blocks for the Lagrange multiplier space. Also, if we want to solve the temperature field, we can add the blocks dynamically for the temperature. This feature is particularly convenient and allows working with a very generic solver. This dynamic construction for the matrix is presented in the snippet 8.2, where  $X_h$ ,  $L_h$  and  $T_h$  are respectively the fluid, the Lagrange multiplier and the temperature FE spaces.

## 2 The Advection Diffusion Solver

In addition to the Navier-Stokes solver presented in the previous section, we needed a similar tool for the advection-diffusion equations. This class was initially implemented to test the different versions of the stabilization method. Later it has been used for the resolution of the turbulence model. To this purpose, we eventually added the possibility to have a reaction term in the equation.

The design of this class is very similar to the Navier-Stokes solver. We will only present the main features in this section. The prototype of the class is

---

```

template <typename MeshType, int Order=1>
class AD

```

---



Listing 8.2 – Example of Dynamic Block Construction for Matrices in Feel++

---

```

1 void initMatrix( sparse_matrix_ptr& M )
2 {
3     if ( M_slip_bc || M_aero )
4     {
5         BlocksBaseGraphCSR myblockGraph(M_nblocks,M_nblocks);
6         myblockGraph(0,0) = stencil(_test=Xh,_trial=Xh)->graph();
7
8         if ( M_slip_bc )
9         {
10            myblockGraph(0,M_nblocks-1) =
11                stencil(_test=Xh,_trial=Lh)->graph();
12            myblockGraph(M_nblocks-1,0) =
13                stencil(_test=Lh,_trial=Xh)->graph();
14            myblockGraph(M_nblocks-1,M_nblocks-1) =
15                stencil(_test=Lh,_trial=Lh)->graph();
16        }
17
18        if ( M_aero )
19        {
20            myblockGraph(0,1) = stencil(_test=Xh,_trial=Th)->graph();
21            myblockGraph(1,0) = stencil(_test=Th,_trial=Xh)->graph();
22            myblockGraph(1,1) = stencil(_test=Th,_trial=Th)->graph();
23        }
24        M = backend()->newBlockMatrix(_block=myblockGraph);
25    }
26    else
27        M = backend()->newMatrix(Xh,Xh);
28 }

```

---

**Algorithm 10:** Resolution algorithm in the NS class using Newton method:SOLVE ( $\mu, \kappa$ )

---

```

1  Assemble diffusion operators  $M_D^u = \mu \nabla \mathbf{u}_h : \nabla \mathbf{v}_h$ ;
2  Assemble null divergence operators  $M_B - p_h(\nabla \cdot \mathbf{v}_h) + (\nabla \cdot \mathbf{u}_h)q_h$ ;
3  if Temperature Field then
4      | Assemble diffusion operators  $M_D^T = \kappa \nabla T_h \cdot \nabla S_h$ ;
5  end
6  if Initialize with Stokes then
7      | Solve Stokes problem;
8  end
9  repeat
10     |  $\mathbf{J} \leftarrow \text{UPDATEJACOBIAN}(U)$ ; // see algorithm 11
11     |  $\mathbf{R} \leftarrow \text{UPDATERESIDUAL}(U)$ ; // see algorithm 12
12     | Solve  $\mathbf{J}\Delta = -\mathbf{R}$ ;
13     |  $U \leftarrow U + \Delta$ ;
14 until  $\frac{\|\mathbf{R}(U^k)\|}{\|\mathbf{R}(U^0)\|} \leq \delta_{Rtol}$  or  $\|\mathbf{R}(U^k)\| \leq \delta_{Atol}$ ;
15 return  $U$ ;

```

---

---

**Algorithm 11:** Updating step for the Jacobian, see resolution algorithm 10:

UPDATEJACOBIAN (U)

---

```

1  J ← 0;
2  Assemble trilinear contribution of J;
3  J ← J +  $M_B$ ; // null divergence operators
4  if Turbulence model then
5      Update  $\mu_t$  with the model; // Spalart Almaras or  $k-\omega$ 
6      Re-assemble Diffusion operator  $M_D^u(\mu + \mu_t)\nabla\mathbf{u}_h : \nabla\mathbf{v}_h$ ;
7      if Temperature Field then
8          Re-assemble Diffusion operator  $M_D^T = (\kappa + \frac{\mu_t}{Pr})\nabla T_h \cdot \nabla S_h$ ;
9      end
10 end
11 J ← J +  $M_D^u$ ;
12 if Temperature Field then
13     J ← J +  $M_D^T$ ;
14 end
15 if Use  $\Psi tc$  then
16     Update pseudo time step;
17     Update continuation matrix  $M_M$ ; // See chapter 2
18     J ← J +  $M_M$ ;
19 end
20 if Use Stabilization then // Detailed in chapter 3
21     Update Stabilization coefficient;
22     Assemble Stabilization operators  $M_S^u$  and  $M_S^T$ ;
23     J ← J +  $M_S^u + M_S^T$ ;
24 end
25 Edit J to set boundary conditions;
26 return J;

```

---

---

**Algorithm 12:** Updating step for the Residual, see resolution algorithm 10:

UPDATERESIDUAL (U)

---

```

1  R ← 0;
2  Assemble trilinear contribution of R;
3  Update right hand side FR;                                // If needed
4  R ← R + FR;
5  R ← R + MBU;
6  R ← R + MDuU;
7  if Temperature Field then
8  |   R ← R + MDTU;
9  end
10 if Use Stabilization then                                // Detailed in chapter 3
11 |   Update Stabilization coefficient;
12 |   Assemble Stabilization operators FSu and FST;
13 |   R ← R + FSu + FST;
14 end
15 Edit R to set boundary conditions;
16 return R;

```

---

where

- `meshType` is the type of mesh.
- `Order` is the polynomial order  $k$  of the Lagrange FE space used for the resolution :  $X_h = \mathbb{L}_{ch}^k$ .

Once the solver is initialized, we can set:

- the advection field,
- the diffusion field or constant,
- the reaction field or constant,

and again we can customize the solver with command line options to choose

- the stabilization method (if needed)
- the shock capturing method (if needed)

We propose a short example in the snippet 8.3. This application was used to resolve the thermal layer test case 4.1. The `solveSc()` method allows solving the system using shock capturing methods. The non-linearity of the model is then treated with a fixed point iterative method.

Listing 8.3 – Minimal application solving the Thermal Layer test case, using the AD class

---

```

1 int main( int argc, char** argv )
2 {
3     typedef Mesh<Simplex<2,1> mesh_t;
4
5     Environment env( _argc=argc, _argv=argv );
6
7     double sigma = doption("sigma");           // read from command line options
8     double mu = doption("mu");                 // read from command line options
9
10    auto mesh = loadMesh( _mesh=new mesh_t );
11    auto ad_solver = AD<mesh_t,2>( mesh, "tlayer" );
12
13    ad_solver.setDiff( mu );                     // set the diffusivity
14
15    // set the advection field v=(2*y,0)
16    ad_solver.adcField().on( elements(mesh), vec(2*Py(),cst(0.) ) );
17
18    if ( sigma!=0 )                             // set the reaction/mass constant
19        ad_solver.setMass( sigma );
20
21    auto u = ad_solver.solveSc();                 // solve the system
22 }

```

---

Listing 8.4 – Spalart Allmaras model: resolution loop using the AD class

---

```

1 M_solver = boost::make_shared<solver_type>(mesh);
2 M_solver->adcField().on( elements(mesh), idv(u) );
3
4 do
5 {
6     // update RHS and diffusion coefficient
7     M_solver->diffField().on( elements(mesh), 1./sigma*(cst(mu)+idv(mut)) );
8     M_solver->rhsField().on( elements(mesh),
9                             cb2/sigma*gradv(mut)*trans(gradv(mut))
10                            + cb1*(cst(1.)-ft2)*S*idv(mut)
11                            - (cw1*fw - cb1/k2*ft2)*idv(mut)*idv(mut)/d2 );
12
13    // solve the equation
14    auto results = M_solver->solve( 3*mu);
15
16    // check the norm of the increment
17    double norm2 = normL2( elements(mesh), idv(mut) );
18    diff = normL2( elements(mesh), idv(mut)-idv(results) )/norm2;
19
20    // iterate
21    mut.on( elements(mesh), max(idv(results), cst(0.) ) );
22    i++;
23 } while ( i<ioption(_name="sa.maxit" ) && doption( _name="sa.tol">diff );

```

---

### 3 Turbulence Models

A direct application of the AD class is the resolution of the equations of the turbulence models presented in the chapter 4. We implemented two class `TurbulenceSA` and `TurbulenceKOmega` both inherited from a base class

---

```

template <typename FluidSpaceType, int Order>
class TurbulenceBase

```

---

where `FluidSpaceType` is the type of the fluid FE element space and `Order` is the desired polynomial order for the FE space used to solve the turbulent equations.

Both models presented in this chapter involved the resolution of transport equations. Our advection diffusion solver is used within the turbulence class to solve these equations. The snippet 8.4 presents the resolution loop in our Spalart Allmaras class.

Besides the resolution of the transport equation, we also need to compute the

Listing 8.5 – Fast marching algorithm used to compute the field of the distance to wall

---

```

1 auto Yh = Pch<1>( mesh );
2 auto Yh0 = Pdh<0>(mesh);
3
4 auto thefms = fms( Yh );
5 auto phio = Yh->element();
6 auto mark = Yh0->element();
7
8 // the first cell of the mesh on walls is initialized with the distance hmin
9 for ( int i=0 ; i<wall_marker.size() ; i++ )
10     phio.on( _range=elementsWithMarkedFaces( mesh, wall_marker[i]),
11             _expr= hMin() );
12
13 // all the wall boundary are set to distance 0
14 for ( int i=0 ; i<wall_marker.size() ; i++ )
15     phio.on( markedfaces( mesh, wall_marker[i]), cst(0) );
16
17 // we mark the cells which are already initialized
18 for ( int i=0 ; i<wall_marker.size() ; i++ )
19     mark += project( _space=Yh0,
20                    _range= elementsWithMarkedFaces( mesh, wall_marker[i]),
21                    _expr=cst(1) );
22 mesh->updateMarker2( mark );
23
24 // the fast marching compute the distance field
25 M_distance = thefms->march( phio, true);

```

---

distance of each point of  $\Omega$  to the nearest point of the wall. The evaluation of this distance to the wall can be particularly challenging - or at least expensive - for complex meshes. We used the fast marching feature of the Feel++ library to compute this distance. This algorithm is entirely parallel and allows fast and accurate computation of the distance to the wall. More details on this feature and its implementation are proposed in [30]. Snippet 8.5 presents an example of using the fast marching to recover the distance to the wall.

### 3.1 The Stabilization Parameters

The final important tool of our framework is the class `StabCoeff`. It allows computing the stabilization parameter from the advection field and the diffusion field (or constant). The version of the stabilization parameter  $\tau$  and the version of the characteristic size  $h$  can be chosen with command line options. This class is used in both the NS and the AD class. Once the class is initialized, the coefficient  $\tau$  is computed by calling

---

```
auto tau = M_stab->update( idv(u), cst(mu) );
```

---

where `u` is the convection field and `mu` the diffusion (can be a field).

We detail some interesting points of this class in the next paragraphs.

#### Initialization: computation of $\lambda_K$

As explained in the chapter 3, the coefficient  $\lambda_K$  is an essential ingredient of efficient stabilization parameters. This coefficient is used to scale the local Reynolds number with respect to the polynomial order of the FE space. In practice,  $\lambda_K$  is

computed as the solution of the local eigenvalue problem

$$\lambda_K = \max_{v \in (\mathbb{P}^k(K) \setminus \mathbb{R})} \frac{\|\Delta v\|_K^2}{\|\nabla v\|_K^2}. \quad (8.1)$$

To solve this eigen value problem, we first assemble the laplacian and the gradient operators on each cell, in a global matrix. For that we use the Feel++ language,

---

```
auto fA = form2( _trial=Xh, _test=Xh, _matrix=matA );
fA = integrate( _range=elements(mesh),
               _expr=inner(laplacian(w), laplacian(w)) );
auto fB = form2( _trial=Xh, _test=Xh, _matrix=matB );
fB = integrate( _range=elements(mesh),
               _expr=inner(grad(w), grad(w)) );
```

---

As specified in (8.1) we have to exclude the constants from the solutions. We chose to add a penalisation term to take care of this constraint.

---

```
fB += integrate( _range=elements(mesh),
                _expr=penal*inner(id(w), id(w)) );
```

---

Once the operators are assembled, we can extract the local matrices associated to each cell.

---

```
for ( int i=0; i<nLocalIndices; ++i )
{
  for ( int j=0; j<nLocalIndices; ++j )
  {
    localMatA(i, j) = matA(extractIndices[i], extractIndices[j]);
    localMatB(i, j) = matB(extractIndices[i], extractIndices[j]);
  }
}
```

---

We use an Eigen solver to recover the eigen values of the generalized eigen values problem associated to these two matrices. The value of  $\lambda_K$  on each cell is the max of this eigen values.

---

```
GeneralizedSelfAdjointEigenSolver<MatrixXd> es( localMatA, localMatB );
VectorXd eigen_vls= es.eigenvalues();

double lambda = eigen_vls.maxCoeff();
lambdaK.on( idedelements( mesh, elt.id() ), cst( lambda) );
```

---

## Evaluation of the Characteristic Length $h$

We defined and studied three different definition of the characteristic length in the section 2.2.

- $h_m$  is the length of the smallest edge of the cell,
- $h_d = \frac{|\mathbf{a}|_p}{|\tilde{\mathbf{a}}|_p} \tilde{h}$  with  $\tilde{\mathbf{a}} = \mathbf{M}\mathbf{J}^{-1}\mathbf{a}$ , is the length of the cell in the direction of the current convection field  $\mathbf{a}$ .
- $h_h$  is a definition proposed by Harari and Hughes, evaluated from the measure of the cell and the vertex-to-centroid distance.

The evaluation of  $h_m$  is already implemented in Feel++ and did not required any development.

The evaluation of  $h_d$  was not very painful, thanks to Feel++ language facilities. As explain in the stabilization chapter, the evaluation of  $h_d$  is made by comparing the convection field  $a$  with the convection field  $\tilde{a}$  on a reference cell. To evaluate  $\tilde{a}$  we first define the matrix  $M$

---

```
auto M = mat<2,2>( cst(1.), cst(1./2), cst(0) , cst(std::sqrt(0.75)) );
```

---

then we can evaluate  $h_d$  using the Feel++ keyword `JinvT()` which compute the transposed inverse inverse of the jacobian of the geometric transformation associated to the current cell,

---

```
h_d.on( elements(mesh), 2*sqrt(inner(a)/inner(M*trans(JinvT())*a)) );
```

---

The evaluation of  $h_h$  requires some more work. The first issue is that the definition depends on the shape of the cells. We will use the 2D simplex as an illustration. For this particular case,  $h_h$  can be expressed as

$$h_h = \frac{4S}{\sqrt{3 \sum_{a=1}^3 |x_i - x_c|^2}}. \quad (8.2)$$

This quantity has to be evaluated on each cell. The following snippet shows how we compute  $h_h$  using Feel++ language and a loop over the elements of the mesh. Thanks to Feel++ parallel implementation, this evaluation is split using the domain decomposition on each process.

---

```
// loop over the elements of the mesh
for ( auto const& eltWrap : elements(mesh) )
{
    auto const& elt = unwrap_ref( eltWrap );

    auto b = elt.barycenter();
    double s = 0;

    // loop over the vertices of the cell
    for ( int k=0; k<n_vertices; k++ )
    {
        auto const& c = elt.point(k).node();
        // evaluate the vertex-to-centroid distance
        s += math::pow(c[0]-b[0],2) + math::pow(c[1]-b[1],2);
    }

    h_h.on( idedelements( mesh,elt.id() ), 4*meas()/math::sqrt( 3*s ) );
}
}
```

---

## Stabilization Parameter on a Reference Mesh

When using the Reduced Basis Method with geometric parameters, we rewrite the weak formulation on a reference domain, using the composition rules for differentiation operators. However, we also have to adapt the stabilization parameter. It has to be evaluated as if it was on the parametrized domain. In the following paragraphs, we denote by  $\tilde{J}$  the Jacobian of the mapping from the parametrized to the reference domain.

The evaluation of  $h_d$  on the reference domain is particularly convenient as long as we know  $\tilde{J}$ . We simply have to compose with  $\tilde{J}^{-1}$  in the classical formula.

---

```
h_d.on( elements(mesh), 2*sqrt(inner(a)/inner(M*trans(JinvT())*inv(M_J)*a)) );
```

---

where  $M_J$  represents  $\tilde{J}$ .

The computation of  $h_h$  on a reference domain was a bit more complex since we have to recompute the vertex-to-centroid distance for each new parametrized domain. To optimize this procedure, we first store discontinuous vector fields  $M_{di}$  containing for each cell the vector  $B\vec{V}_i$ , where  $B$  is the barycenter of the cell, and  $V_i$  is the  $i$ th vertices of the cell.

---

```
for ( auto const& eltWrap : elements(mesh) )
{
    auto const& elt = unwrap_ref( eltWrap );
    auto b = elt.barycenter();

    auto p0 = elt.point(0).node();
    auto p1 = elt.point(1).node();
    auto p2 = elt.point(2).node();

    M_d1.on( idedelements(mesh,elt.id()),
             (p0[0]-b[0])*oneX() + (p0[1]-b[1])*oneY() );
    M_d2.on( idedelements(mesh,elt.id()),
             (p1[0]-b[0])*oneX() + (p1[1]-b[1])*oneY() );
    M_d3.on( idedelements(mesh,elt.id()),
             (p2[0]-b[0])*oneX() + (p2[1]-b[1])*oneY() );
}
```

---

Then we can use these vector fields to compute the vertex-to-centroid distance, using the image of this vector by the jacobian matrix  $\tilde{J}$ .

---

```
h_h.on( elements(mesh),
        4*meas()*det(idv(M_J))/sqrt( 3*(inner(idv(M_J)*idv(M_d1))
        +inner(idv(M_J)*idv(M_d2))
        +inner(idv(M_J)*idv(M_d3))) ) );
```

---

We finally can evaluate the coefficient  $\lambda_K$  on the parametrized domain, using the relation  $\lambda_K h_h^2 = \text{cst}$ . Once we have the value of  $\tilde{h}_h$  on the parametrized domain, we can evaluate  $\lambda_K = \lambda_K \frac{h_h}{\tilde{h}_h}$ . With this last ingredient, we can compute any version of the stabilization parameter  $\tau$  on a reference domain, knowing only the Jacobian matrix  $\tilde{J}$ . We can set this matrix in the `StabCoeff` class, using the method `J()`

---

```
M_stab->J().on( elements(mesh), mat<2,2>(cst(1.),cst(0),cst(0),cst(1.)) );
M_stab->J().on( markedelements(mesh,"air4"),
               mat<2,2>(cst((ea-eic)/(ear-eic)),cst(0),cst(0),cst(1.)) );
auto tau = M_stab->update( rhoC_a*conv, kappa_a );
```

---





---

Reduced Basis Framework

---

**Contents**


---

<b>1</b>	<b>Reduced Basis by Block . . . . .</b>	<b>177</b>
1.1	Offline Procedure . . . . .	177
1.2	Online Procedure . . . . .	181
<b>2</b>	<b>Reduced Basis for Aerothermal Problems . . . . .</b>	<b>183</b>
2.1	Offline Procedure . . . . .	183
2.2	Online Procedure . . . . .	185

---

In the previous chapter, we detailed the developments of our aerothermal library. This library was not fully integrated into Feel++ yet. Hereafter, we will present implementation works made directly in the Feel++ library. We begin, in this chapter, with the development of the RB basis framework.

At the beginning of this thesis work, the RB framework was already well developed in Feel++ library, see [105]. The main features are:

- CRBM for elliptic and parabolic linear coercive problems
  - Primal/Dual resolution
  - Coercivity bounds evaluation through  $\min\text{-}\theta$  or SCM
  - Offline/Online residual evaluation
  - Efficient error bounds on the output

- CRBM approximation for non-affine linear coercive problems
  - EIM and error estimators
- RBM for non-linear and multiphysic problems
  - EIM + SER
  - RB for composite spaces (monolithic resolution)

The multiphysics problems were already treated in the framework, see [24], but the reduced basis spaces were always built monolithically. The implementation of composite reduced spaces was the first contribution we needed in Feel++ RB framework. This feature is essential, for instance, with saddle point Stokes problems when we need to build a velocity space with higher dimension than the pressure space, see chapter 5. We eventually implemented a *block* RB framework which allows building composite RB spaces with a block-structured resolution. In addition, we also implemented the tools necessary to deal with trilinear and especially aerothermal problems.

To keep the framework coherent and easy to maintain, it was essential to integrate the new features into the existing code. Thanks to the C++ inheritance paradigm, it was very convenient to derive new models from the classic CRB implementation. To have a better understanding of our choice, we first present an overview of this classic CRB framework on figure 9.1. The main classes are

- `ModelCrbBase` is used to implement a new RB model. Each new RB application is implemented in a new class inherited from `ModelCrbBase`. It provides the necessary interface with the RB framework. The user can then specify the characteristic of the model: geometry, FE space, parameter space, affine decomposition... `ModelCrbBase` also comes with some options used to set up the RB algorithm (transient or steady, linear or not, ...).
- `CRB` class executes the suitable offline algorithm and produces the RB space which is eventually stored in a database.
- `CRBModel` provides the FEM resolution algorithms specific to each problem.
- `EIM` is used to recover an affine approximation for non-affine or non-linear problems, see chapter 6.

The classes `ModelCrbBase` and `EIM` are not impacted by the implementation of our new features. We will essentially focus on the modifications in `CRB` and `CRBModel` which are derived respectively in `CRBBlock` and `CRBModelBlock` and later in `CRBAero` and `CRBModelAero`.

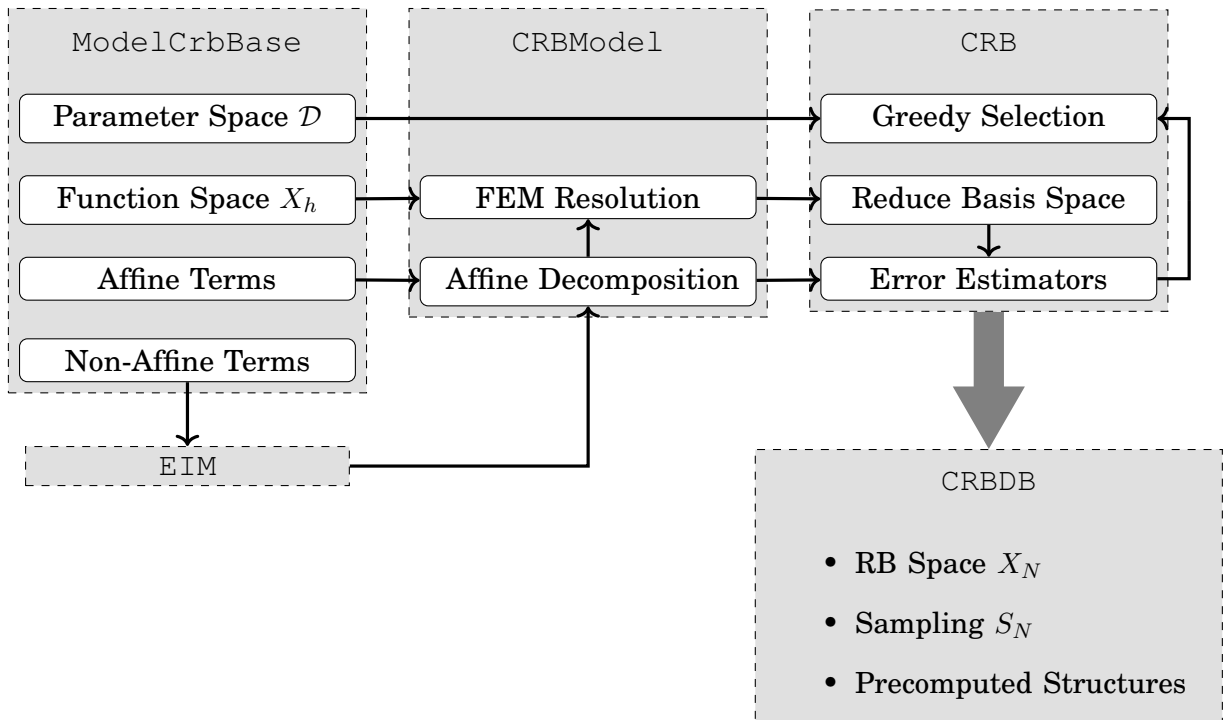


Figure 9.1 – Main classes in Feel++ RB framework

## 1 Reduced Basis by Block

Most of the modification needed for the block construction of the reduced spaces have been done in the new class

```

template <typename TruthModelType>
class CRBBlock : public CRB<TruthModelType>

```

Only a few functions have been rewritten in the inherited class. We detail the most important modifications in this section.

### 1.1 Offline Procedure

The offline construction of the reduced basis is entirely managed in the `offline()` method of `CRB`. We keep this algorithm unchanged, but we use virtual functions for the steps which are impacted by the block structure. In algorithm 13 we present a very generic version of the offline routine. The steps in blue have been rewritten in the `CRBBlock` class and will be detailed in further paragraphs.

**Remark 14.** We chose to implement a block version of the RB algorithm where the number of blocks is not fixed. It depends on the number of subspaces in the product spaces. To implement this class, we had to use meta-programming. In *Feel++*, a composite space is a class, mostly defined by template arguments. This construction is very powerful, but it might also provide some additional constraints.

---

**Algorithm 13:** Offline Construction of the Reduced Basis Space
 

---

```

1  Generate the super sampling  $\Xi \subset \mathcal{D}$ ;
2  if Use predefined sampling then
3    | Generate Sampling  $W$ ;
4  end
5  Chose  $\mu_1 \in \Xi$ ;
6  repeat
7    |  $N \leftarrow N + 1$ ;
8    |  $S_N \leftarrow S_{N-1} \cup \{\mu_N\}$ ;
9    |  $U(\mu_N) \leftarrow$  Solve the FE problem with  $\mu_N$ ;
10   | Add new vector in the Basis;
11   | Orthonormalize the basis;
12   | Precompute the RB matrix and vectors;
13   | if Use predefined sampling then
14     |  $\mu_{N+1} \leftarrow W[N + 1]$ ;
15   | else
16     | for  $\mu \in \Xi$  do
17       | Compute error indicator,  $\Delta_N(\mu)$  ;
18     | end
19     |  $\mu_{N+1} = \arg \max_{\mu \in \Xi} \Delta_N(\mu)$ ;
20   | end
21 until  $\Delta_N(\mu_{N+1}) < \delta_{tol}$  or  $N \geq N_{max}$ ;

```

---

Listing 9.1 – Minimal example of a loop over a template range, using meta-programming

---

```

1  static const int n_block = space_type::nSpaces;
2  typedef typename mpl::range_c< int, 0, n_block > rangespace_type;
3
4  template <typename T>
5  struct doSomethingOverEachSubSpacesByBlock
6  {
7      template <typename T>
8      void operator()( T const& t )const
9      {
10         // access a subspace using template value
11         auto subXh = Xh->template functionSpace<T::value>();
12         // now I can use this subspace...
13         ...
14     }
15 };
16
17 void doSomethingOverEachSubSpaces()
18 {
19     // instantiate the template struct
20     doSomethingOverEachSubSpacesByBlock my_struct();
21     // instantiate an object of type range
22     rangespace_type range;
23     // loop on the structure
24     boost::fusion::for_each( range, my_struct );
25 }

```

---

In the algorithms presented in this section, we regularly iterate over the subspaces of a composite space. This operation cannot be done in *Feel++* with a classical loop. One solution is to iterate, using tool like `boost::fusion::for_each`, over templated structures. We propose a minimal example of this construction in the snippet 9.1. Note that all the loops over the subspaces, in the following paragraphs, have been implemented using this kind of structures.

We first introduce some notations for the product space and the block structure. We denote by  $X_h = X_{h1} \times \dots \times X_{hR}$  a composite FE space.  $X_h$  can be naturally decomposed (in the sense of the physic) as the product of  $R$  subspaces noted  $X_{hr}$ ,  $1 \leq r \leq R$ . We also assume for each block the affine decomposition:  $\forall 1 \leq r, c \leq R$ ,

$$\begin{aligned}
 a^{rc}(u_r, u_c; \boldsymbol{\mu}) &= \sum_{q=1}^{Q^a} \theta_q^a(\boldsymbol{\mu}) a_q^{rc}(u_r, u_c) \\
 f^r(u_r; \boldsymbol{\mu}) &= \sum_{q=1}^{Q^f} \theta_q^f(\boldsymbol{\mu}) f_q^r(u_r)
 \end{aligned} \tag{9.1}$$

where  $a^{rc}(\cdot, \cdot; \boldsymbol{\mu}) : X_{hr} \times X_{hc} \times \mathcal{D} \rightarrow \mathbb{R}$ ,  $a^{rc}(\cdot, \cdot) : X_{hr} \times X_{hc} \rightarrow \mathbb{R}$  and  $f^r(\cdot; \boldsymbol{\mu}) : X_{hr} \times \mathcal{D} \rightarrow \mathbb{R}$ ,  $f^r(\cdot) : X_{hr} \rightarrow \mathbb{R}$ . Since we are detailing discrete algorithms, it is convenient to identify bilinear/linear forms and matrices/vectors. So we will transparently switch between  $a_q^{rc}$ ,  $f_q^r$  and  $A_q^{rc}$ ,  $F_q^r$ . The underlying block structure is such that

$$A_q = \begin{bmatrix} A_q^{11} & \dots & A_q^{1R} \\ \vdots & & \vdots \\ A_q^{R1} & \dots & A_q^{RR} \end{bmatrix}, \quad F_q = \begin{bmatrix} F_q^1 \\ \vdots \\ F_q^R \end{bmatrix}, \quad U = \begin{bmatrix} u_1 \\ \vdots \\ u_R \end{bmatrix} \tag{9.2}$$

We will not detail the steps [Add new vector in the Basis](#); and [Orthonormalize the basis](#) of the algorithm 13. In the first one, we simply add each subfield  $u_r$  in

the corresponding reduced subspace  $X_{N_r}$ . In the second one, we orthonormalize (if necessary) the basis vectors of each reduced sub-space  $X_{N_r}$ .

The assembly of the reduced matrices and vectors is a bit more technical than the previous steps. The aim is to construct,  $\forall 1 \leq r, c \leq R$  the matrices and vectors

$$\begin{aligned} A_{N_q}^{rc} &= a_q^{rc}(u_r, u_c) = u_r^\top A_q^{rc} u_c, \quad \forall 1 \leq q \leq Q^a, \\ F_{N_q}^r &= f_q^r(u_r) = u_r^\top F_q^r, \quad \forall 1 \leq q \leq Q^f. \end{aligned} \quad (9.3)$$

As presented in algorithm 14, it essentially consists in many nested loops. In this algorithm,  $n_r$  and  $n_c$  denote the number of basis vectors added respectively in  $X_{N_r}$  and  $X_{N_c}$ .

---

**Algorithm 14:** Block construction of the reduced matrices and vectors. cf 1.12 in algorithm 13

---

```

1  for 1 ≤ r ≤ R do                                     // Loop on the Row
2      for 1 ≤ c ≤ R do                                     // Loop on the Columns
3          for 1 ≤ q ≤ Qa do                               // Loop on the affine decomposition
4              /* Update the last rows of ANqrc. */
5              for Nr - nr ≤ i ≤ Nr do                 // Loop on the first basis
6                  for 1 ≤ j ≤ Nc do                       // Loop on the second basis
7                      ξi ← i-th basis in XNr;
8                      ξj ← j-th basis in XNc;
9                      (ANqrc)ij = ξi⊤ Aqrc xj;
10                 end
11            end
12            /* Update the last columns of ANqrc. */
13            for 1 ≤ i ≤ Nr do                             // Loop on the first basis
14                for Nc - nc ≤ j ≤ Nc do             // Loop on the second basis
15                    ξi ← i-th basis in XNr;
16                    ξj ← j-th basis in XNc;
17                    (ANqrc)ij = ξi⊤ Aqrc xj;
18                end
19            end
20        for 1 ≤ q ≤ Qf do                               // Loop on the affine decomposition
21            /* Update the last rows of FNqr. */
22            for Nr - nr ≤ i ≤ Nr do                 // Loop on the basis
23                ξi ← i-th basis in XNr;
24                (FNqr)i = ξi⊤ Fqr;
25            end
26    end

```

---

## Error Estimator

Using an efficient error bound requires to have two main components: bounds for the coercivity/inf-sup constants and dual norm of the residual. The first ingredient is problem dependent and cannot be implemented generically for the multiphysics problems. We already have solutions (min-theta algorithm, SCM) for the evaluation of the coercivity constant in the classic CRB framework but these methods have not been extended to the block structure due to lack of time.

The evaluation of the residual has been implemented but only in the particular case of 2 subspaces. The offline precomputation of the different structures can be very expensive since it requires the resolution of many linear systems to evaluate the Riesz representations, see appendix 3. In the original CRB class, these Riesz representations were reevaluated multiple times because it was inconceivable to store all these vectors in the memory. However, these multiple reevaluations are very costly and especially for block construction. We finally decided to store the Riesz representations on the disk and to reread them each time they are needed. The dual norm of the residual is then used as an error indicator for certain problems, see Stokes test case 3.

The implementation of the residual evaluation for an arbitrary number of subspaces is on the todo list.

## 1.2 Online Procedure

The online phase of the reduced basis methodology essentially consists in the resolution of the reduced problem. For now the CRBBlock online integrates a fixed point algorithm, but a Newton method has been implemented in CRBAero and will be detailed in a further section.

The fixed point procedure is detailed in the algorithm 15. We simplified the algorithm by assuming that all the reduced subspaces have the same dimension  $N$ . In practice, with the selective Gram-Schmidt orthonormalization<sup>7</sup> and the addition of the supremizer, the dimension of the sub-spaces may vary. The only remarkable thing is the assembly of the matrix  $A$  and the vector  $F$ . In practice, this step is realized using the block operators of the Eigen library which simplifies the treatment of the blocks. Here is an example of an addition is the first block of the matrix  $A$  of size  $N_0 \times N_0$

---

```
A.block(0,0,N0,N0) += betaAqm[q][m]*blockAqm[0][0][q][m].block(0,0,N0,N0);
```

---



---

**Algorithm 15:** Online fixed point procedure for block RB.
 

---

```

1  Input:  $N, \mu$ 
2  Create matrix  $A$  of size  $NR \times NR$ ;
3  Create vector  $F$  of size  $NR$ ;
4  Compute  $\theta_q^a(\mu)$  and  $\theta_q^f(\mu)$ ;
5  repeat
6    for  $1 \leq r \leq R$  do
7      for  $1 \leq c \leq R$  do
8        for  $1 \leq q \leq Q^a$  do
9          for  $(r-1)N \leq i < rN$  do
10           for  $(c-1)N \leq j < cN$  do
11              $(A)_{ij} \leftarrow (A)_{ij} + \theta_q^a(\mu)A_{Nq}^{rc}$ ;
12           end
13         end
14       end
15     for  $1 \leq q \leq Q^a$  do
16       for  $(r-1)N \leq i < rN$  do
17          $(F)_i \leftarrow (F)_i + \theta_q^f(\mu)F_{Nq}^r$ ;
18       end
19     end
20   end
21    $U^* \leftarrow U$ ;
22    $U \leftarrow$  Solution of  $AU = F$ ;
23   is_finished  $\leftarrow$  is_linear or  $\|U - U^*\|_2 \leq \delta_{\text{tol}}$ ;
24 until is_finished;
25 return  $U$ ;

```

---

## 2 Reduced Basis for Aerothermal Problems

The reduced basis construction for aerothermal problems, presented in the chapter 7, is made in a decoupled way. We construct three reduced sub-spaces: for the velocity, the pressure and the temperature. The `CRBBlock` class, introduced in the previous section, provides the tools to construct these reduced sub-spaces. In this section, we present some details of a new inherited class

---

```
template <typename TruthModelType>
class CRBAero : public CRBBlock<TruthModelType>
```

---

This class contains the algorithm necessary for the construction of the reduced basis and for the online evaluation of reduced solution.

### 2.1 Offline Procedure

We reuse the algorithm 13 of the previous section and we add in red, in algorithm 16, the steps which have to be modified for the RB aerothermal framework. The

---

#### Algorithm 16: Offline Construction of the Reduced Basis Space for Aerothermal RB

---

```

1  Generate the super sampling  $\Xi \subset \mathcal{D}$ ;
2  if Use predefined sampling then
3  |   Generate Sampling  $W$ ;
4  end
5  Chose  $\mu_1 \in \Xi$ ;
6  repeat
7  |    $N \leftarrow N + 1$ ;
8  |    $S_N \leftarrow S_{N-1} \cup \{\mu_N\}$ ;
9  |    $U(\mu_N) \leftarrow$  Solve the FE problem with  $\mu_N$ ;
10 |   Add new vector in the Basis;
11 |   Orthonormalize the basis;
12 |   Precompute the RB matrix and vectors;
13 |   Precompute the reduced trilinear contributions;
14 |   if Use predefined sampling then
15 |   |    $\mu_{N+1} \leftarrow W[N + 1]$ ;
16 |   else
17 |   |   for  $\mu \in \Xi$  do
18 |   |   |   Compute error indicator,  $\Delta_N(\mu)$ ;
19 |   |   end
20 |   |    $\mu_{N+1} = \arg \max_{\mu \in \Xi} \Delta_N(\mu)$ ;
21 |   end
22 until  $\Delta_N(\mu_{N+1}) < \delta_{tol}$  or  $N \geq N_{max}$ ;
```

---

resolution of the finite element problem is particular since it implies quadratic non-linearity. This part is modified in the class `CRBModelAero`. The resolution uses a Newton iterative method and the model provides the affine decomposition of the trilinear, bilinear and linear forms.

The precomputation of the reduced structures now integrates the evaluation of the order-3 tensors, as explained in the chapter 7, in equation (7.11). In the particular context of the aerothermal problems, we know that we only have quadratic contributions in certain blocks. We will only assemble the order-3 tensors  $T_q^1$  and  $T_q^2$  for the convection operators respectively of the fluid and temperature equations. In addition of the construction 14, the class `CRBAero` will also assemble these tensors, following the algorithm 17, where  $N_u$  and  $N_T$  are the current dimension of respectively the velocity and the temperature spaces,  $n_u$  and  $n_T$  are the number of basis vectors added in these spaces.

---

**Algorithm 17:** Precomputation of the reduced order-3 tensors

---

```

1  for  $1 \leq q \leq Q^c$  do
2    for  $1 \leq i \leq N_u$  do
3      for  $1 \leq j \leq N_u$  do
4        for  $1 \leq k \leq N_u$  do
5          if  $i \geq N_u - n_u$  or  $j \geq N_u - n_u$  or  $k \geq N_u - n_u$  then
6             $(T_q^1)_{ijk} \leftarrow c_q(\zeta_i, \zeta_j, \zeta_k)$ ;
7          end
8        end
9      end
10     end
11  end
12  for  $1 \leq q \leq Q^a$  do
13    for  $1 \leq i \leq N_T$  do
14      for  $1 \leq j \leq N_u$  do
15        for  $1 \leq k \leq N_T$  do
16          if  $i \geq N_T - n_T$  or  $j \geq N_u - n_u$  or  $k \geq N_T - n_T$  then
17             $(T_q^2)_{ijk} \leftarrow a_q(\xi_i, \zeta_j, \xi_k)$ ;
18          end
19        end
20      end
21    end
22  end

```

---

## 2.2 Online Procedure

The online procedure for the aerothermal framework is already well detailed in the chapter 7. We will not rewrite these algorithms here, but we decided to illustrate them with a detailed explanation of the code.

In the class `CRBAero`, the virtual method `onlineSolve` has been totally rewritten. Its prototype is

---

```
typename CRBAero<TruthModelType>::matrix_info_tuple
onlineSolve( size_type N, parameter_type const& mu, std::vector< vectorN_type > & uN,
             std::vector< vectorN_type > & uNdu, std::vector<vectorN_type> & uNold,
             std::vector<vectorN_type> & uNduold, std::vector< double > & output_vector,
             int K, bool print_rb_matrix, bool computeOutput ) const;
```

---

Once the different variables are initialized, we start with the assembly of the bilinear contributions for the Jacobian

---

```
for ( int q=0; q<model->sizeOfBilinearJ(); q++ )
{
  for ( int m=0; m<model->mMaxA(q); m++ )
  {
    // first row
    M_Jbil.block( 0, 0, N0, N0) += betaJqm[q][m]*blockAqm(0,0)[q][m].block(0,0,N0,N0);
    M_Jbil.block( 0, N0, N0, N1) += betaJqm[q][m]*blockAqm(0,1)[q][m].block(0,0,N0,N1);
    M_Jbil.block( 0, N0+N1, N0, N2) += betaJqm[q][m]*blockAqm(0,2)[q][m].block(0,0,N0,N2);
    // second row
    M_Jbil.block(N0, 0, N1, N0) += betaJqm[q][m]*blockAqm(1,0)[q][m].block(0,0,N1,N0);
    M_Jbil.block(N0, N0, N1, N1) += betaJqm[q][m]*blockAqm(1,1)[q][m].block(0,0,N1,N1);
    M_Jbil.block(N0, N0+N1, N1, N2) += betaJqm[q][m]*blockAqm(1,2)[q][m].block(0,0,N1,N2);
    // third row
    M_Jbil.block(N0+N1, 0, N2, N0) += betaJqm[q][m]*blockAqm(2,0)[q][m].block(0,0,N2,N0);
    M_Jbil.block(N0+N1, N0, N2, N1) += betaJqm[q][m]*blockAqm(2,1)[q][m].block(0,0,N2,N1);
    M_Jbil.block(N0+N1, N0+N1, N2, N2) += betaJqm[q][m]*blockAqm(2,2)[q][m].block(0,0,N2,N2);
  }
}
```

---

and for the Residual

---

```
for ( int q=0; q<model->sizeOfLinearR(); q++ )
{
  for ( int m=0; m<model->mMaxF(0,q); m++ )
  {
    M_Rli.segment( 0,N0) += betaRqm[0][q][m]*blockFqm(0)[q][m].head(N0);
    M_Rli.segment( N0,N1) += betaRqm[0][q][m]*blockFqm(1)[q][m].head(N1);
    M_Rli.segment( N0+N1,N2) += betaRqm[0][q][m]*blockFqm(2)[q][m].head(N2);
  }
}
```

---

then, the non-linear solver is called and the rest of the work will be done in the functions

---

```
void
updateJacobianOnline( const map_dense_vector_type& X, map_dense_matrix_type& J,
                    parameter_type const& mu , int N ) const
```

---

and

---

```
void
updateResidualOnline( const map_dense_vector_type& X, map_dense_vector_type& R,
                    parameter_type const& mu , int N ) const
```

---

The Jacobian is updated in three steps. First we add the bilinear contributions, already assembled before the resolution loop

---

```
J += M_Jbil;
```

---

then we assemble the trilinear contributions

---

```

auto betaTri = this->M_model->computeBetaTri( mu );
for ( int q=0; q<model->QTri(); q++ )
{
    for ( int k=0; k<N0; k++ )
    {
        for ( int i=0; i<N0; i++ )
        {
            J(k,i) += betaTri[q]*(blockTriqm(0,0)[q][k].row(i).head(N)).dot(X.segment(0,N0));
            J(k,i) += betaTri[q]*(blockTriqm(0,0)[q][k].col(i).head(N)).dot(X.segment(0,N0));
        }
    }
    for ( int k=0; k<N2; k++ )
    {
        for( int i=0; i<N2; i++ )
        {
            J(N0+N1+k,N0+N1+i) += betaTri[q]*(blockTriqm(2,0)[q][k].row(i)).dot(X.segment(0,N0));
        }

        for( int j=0; j<N0; j++ )
        {
            J(N0+N1+k,j) += betaTri[q]*(blockTriqm(2,0)[q][k].col(j)).dot(X.segment(N0+N1,N2));
        }
    }
}

```

---

and finally we assemble the non linear contributions (if needed), with coefficients given by the EIM

---

```

for ( int q=model->sizeOfBilinearJ(); q< model->Qa(); q++ )
{
    for ( int m=0; m<model->mMaxA(q); m++ )
    {
        // first row
        J.block( 0, 0, N0, N0) += betaJqm[q][m]*blockAqm(0,0)[q][m].block(0,0,N0,N0);
        J.block( 0, N0, N0, N1) += betaJqm[q][m]*blockAqm(0,1)[q][m].block(0,0,N0,N1);
        J.block( 0, N0+N1, N0, N2) += betaJqm[q][m]*blockAqm(0,2)[q][m].block(0,0,N0,N2);
        // second row
        J.block( N0, 0, N1, N0) += betaJqm[q][m]*blockAqm(1,0)[q][m].block(0,0,N1,N0);
        J.block( N0, N0, N1, N1) += betaJqm[q][m]*blockAqm(1,1)[q][m].block(0,0,N1,N1);
        J.block( N0, N0+N1, N1, N2) += betaJqm[q][m]*blockAqm(1,2)[q][m].block(0,0,N1,N2);
        // third row
        J.block( N0+N1, 0, N2, N0) += betaJqm[q][m]*blockAqm(2,0)[q][m].block(0,0,N2,N0);
        J.block( N0+N1, N0, N2, N1) += betaJqm[q][m]*blockAqm(2,1)[q][m].block(0,0,N2,N1);
        J.block( N0+N1, N0+N1, N2, N2) += betaJqm[q][m]*blockAqm(2,2)[q][m].block(0,0,N2,N2);
    }
}

```

---

We use the same idea for the update of the Residual. We start by adding the bilinear and linear contributions

---

```

R += M_Rli;
R += M_Jbil*X;

```

---

then the trilinear part

---

```

auto betaTri = this->M_model->computeBetaTri( mu );
for ( int q=0; q<model->QTri(); q++ )
{
    for ( int k=0; k<N0; k++ )
        for ( int i=0; i<N0; i++ )
            temp(k,i) = (blockTriqm(0,0)[q][k].row(i)).dot(X.segment(0,N0));
    for ( int k=0; k<N2; k++ )
        for ( int i=0; i<N2; i++ )
            temp(N0+N1+k,N0+N1+i) = (blockTriqm(2,0)[q][k].row(i)).dot(X.segment(0,N0));
    R += betaTri[q]*temp*X;
}

```

---

and finally the non-linear terms

---

```

for ( int q=model->sizeOfLinearR(); q<model->Ql(); q++ )
{
    for ( int m=0; m<model->mMaxF(0,q); m++ )
    {
        R.segment( 0,N0) += betaRqm[0][q][m]*blockFqm(0)[q][m].head(N0);
        R.segment( N0,N1) += betaRqm[0][q][m]*blockFqm(1)[q][m].head(N1);
        R.segment( N0+N1,N2) += betaRqm[0][q][m]*blockFqm(2)[q][m].head(N2);
    }
}

```

---

---

EIM Framework for Discrete Operators

---

**Contents**


---

<b>1</b>	<b>Design of the Implementation . . . . .</b>	<b>187</b>
<b>2</b>	<b>Construction of the Interpolation Space . . . . .</b>	<b>189</b>

---

The work on the EIM algorithm presented in the section 2 represent the most important implementation effort. This feature has been entirely coded and integrated into the Feel++ RB framework during this thesis. We present in this chapter some interesting details and advanced C++ techniques used in this implementation.

The Feel++ library already had an EIM framework when we started this implementation. This initial version uses analytic expressions and their evaluation on the quadrature points. Also, this version was not adapted to the treatment of discrete operators. To distinguish this version from the discrete one, we used the abusive notation DEIM for our new implementation. To remain consistent with the code we will keep using this denomination in this chapter.

## 1 Design of the Implementation

The DEIM framework have been split into four classes, with the inheritance relations detailed on the figure 10.1.

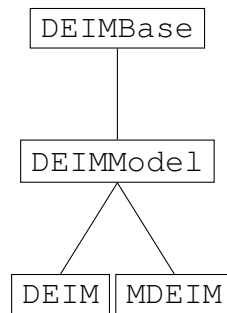


Figure 10.1 – Inheritance Diagram of the Classes in the DEIM Framework

### The class

---

```

template <typename ParameterSpaceType, typename SpaceType, typename TensorType>
class DEIMBase : public CRBDB

```

---

contains the essence of the DEIM algorithm. This class is managing all the algebraic operations which do not require an access to the RB model:

- The greedy algorithm
- The research of maximum entry in the tensors
- The evaluation of a specific entry in a tensor
- The computation of the  $\beta_m(u, \boldsymbol{\mu})$

The class is templated by the type of the tensor. It allowed to write only one class for both matrices and vectors. The class also inherits from the class CRBDB which manages to save the essential information in a database.

### The class

---

```

template <typename ModelType, typename TensorType>
class DEIMModel :
    public DEIMBase<typename ModelType::parameterspace_type,
                  typename ModelType::space_type, TensorType>

```

---

is the interface with the RB model. This class will call the assembly functions of the model and possibly the resolution functions for non-linear problems. It also contains the mapping for the degrees of freedom (DoFs) between the interpolation space  $X_{\mathcal{I}_M}$  and the FE space  $X_h$ . This point is detailed in the section 2.

### The class

---

```

template <typename ModelType>
class DEIM :
    public DEIMModel<ModelType, typename Backend<typename ModelType::value_type>::vector_type>

```

---

and

---

```

template <typename ModelType>
class MDEIM :
    public DEIMModel<ModelType,
                  typename Backend<typename ModelType::value_type>::sparse_matrix_type>

```

---

are specific to the type of tensors. `DEIM` is the instantiation for the vectors and `MDEIM` the instantiation for the matrices. These classes also provide methods for the construction of the interpolation mesh  $\mathcal{T}_{\mathcal{I}_M}$

The global structure of the DEIM algorithm is basic and did not present many difficulties for the implementation. The main challenge of the implementation was the creation of the interpolation mesh and space, required for the optimized online evaluation of the  $\beta_m(u, \boldsymbol{\mu})$ .

## 2 Construction of the Interpolation Space

At the end of the DEIM greedy algorithm, we obtain a set of indices  $\mathcal{I}_M$ . To optimized the online computation of the  $\beta_m(u, \boldsymbol{\mu})$ , we have to assemble reduced tensors  $T_{\mathcal{I}_M}(u, \boldsymbol{\mu})$  only on the associated DoFs. This procedure is detailed in the section 2. We will describe here some interesting details.

**Remark 15.** *In the following paragraphs, we only present the construction for a single FE space. The construction for product spaces was also implemented and is basically a loop of the following method on each space.*

### The Interpolation Mesh and Space

The first step is to extract the interpolation mesh  $T_{\mathcal{I}_M}$  from the full mesh  $\mathcal{T}_h$ . From the `index_list`, we have to create a list of the elements in  $\mathcal{T}_h$  that contains the indices in `index_list`

---

```
for ( auto index : index_list )
{
    auto searchGpDof = Xh->dof()->searchGlobalProcessDof( index );
    if ( boost::get<0>( searchGpDof ) )
    {
        size_type gpDof = boost::get<1>( searchGpDof );
        for ( auto const& dof : Xh->dof()->globalDof( gpDof ) )
        {
            size_type eltId = dof.second.elementId();
            if ( mesh->element( eltId ).isGhostCell() )
                continue;
            this->M_elts_ids.insert( eltId );
        }
    }
}
```

---

After this loop, each processor has a set of elements, `M_elts_ids`, containing the local elements to extract. The extraction function is then provided by the `Feel++` library

---

```
auto submesh = createSubmesh( mesh, idelements(mesh, this->M_elts_ids.begin(),
                                             this->M_elts_ids.end()) );
```

---

At this point, each processor contains its local part of the extracted `submesh`. To gather all the contributions, we will write them on the disk and then read the complete extracted mesh with each processor:



---

```

saveGMSHMesh( _mesh=submesh, _filename=this->name(true)+"-submesh.msh" );
Environment::worldComm().barrier();
auto seqmesh = loadMesh( _mesh=new mesh_type,
                        _filename=this->name(true)+"-submesh.msh",
                        _worldcomm= Environment::worldCommSeq() );

```

---

Here the `mpi barrier`, `Environment::worldComm().barrier();`, is essential to ensure that all the processors finished writing there contributions before the `seqmesh` is read. From this sequential mesh we can now build the interpolation space

---

```
Rh = this->newInterpolationSpace(seqmesh);
```

---

Once this space is generated, it will be used in the online model to assemble the reduced tensors  $T_{\mathcal{I}_M}(u, \boldsymbol{\mu})$ . However, we are not interested in all the value in this reduced tensor. We only need the entry corresponding to the index in  $\mathcal{I}_M$ . Here the issue is that the DoFs have entirely different indices in  $X_{\mathcal{I}_M}$  and  $X_h$ . To deal with that problem, we implemented a specific class to build the mapping between a parallel space and a sequential extracted space.

### Sequential/Parallel Mapping

In Feel++, once a submesh is extracted, all the elements are renumbered. The only data remaining unchanged is the numbering of the vertex of the mesh. We will use this property to build the mapping between the two spaces.

We denote by  $X_p$  and  $\mathcal{T}_p$  the parallel space and mesh. These data are split between all the processors. We also note  $X_s$  and  $\mathcal{T}_s$  the sequential space and mesh. All the processors have the entirety of these data. We assume that we have the relation  $\mathcal{T}_s \subset \mathcal{T}_p$ . The map between the two spaces has to contain the two-way correspondence  $\text{DoF}_s \longleftrightarrow (n_p, \text{DoF}_p)$  where  $n_p$  is number of the processor having the parallel DoF  $\text{DoF}_p$  corresponding to the sequential DoF  $\text{DoF}_s$ .

On each processor, we first store for each element in  $\mathcal{T}_s$  the numbering of its vertices

---

```

for ( auto const& eltWrap : elements(Xs->mesh()) )
{
    auto const& elt = unwrap_ref( eltWrap );
    std::set<int> pts_id;
    for ( int p=0; p<elt.nPoints(); p++ )
        pts_id.insert( elt.point(p).id() );
    elts_map_s[pts_id] = elt.id();
}

```

---

Then, each processor will look if its part of  $\mathcal{T}_p$  has elements in commun with  $\mathcal{T}_s$ . For each commun element, the map is then constructed by identifying the DoF associated in  $X_p$  and  $X_s$ . This procedure is detailed in the snippet 10.1.

Once the map is built, we have access to useful methods to do the correspondence between the two spaces:

Listing 10.1 – Construction of the map between the DoF of the sequential and the parallel space

---

```

1 // loop on the elements of the full mesh splitted between all procs
2 for ( auto const& eltWrap : elements(Xp->mesh()) )
3 {
4     auto const& elt = unwrap_ref( eltWrap );
5     std::set<int> pts_id;
6     for ( int p=0; p<elt.nPoints(); p++ )
7         pts_id.insert( elt.point(p).id() + M_shift );
8
9     // check if this elements is in the submesh of the sequential space
10    auto map_it = elts_map_s.find( pts_id );
11    if ( map_it!=elts_map_s.end() ) // the element exists in the seq mesh
12    {
13        int eid_s = map_it->second;
14
15        // loop on each ldof of the element : get the globaldof id associated
16        // and put it in the maps
17        for ( auto const& ldof : Xp->dof()->localDof(elt.id()) )
18        {
19            int gdof_s = Xs->dof()->localToGlobalId( eid_s, ldof.first.localDof() );
20            int gdof_p = ldof.second.index() ;
21
22            auto s_to_p_it = M_s_to_p.find( gdof_s );
23            if ( s_to_p_it==M_s_to_p.end() )
24            {
25                M_s_to_p[gdof_s] = std::make_pair( Xp->worldComm().globalRank(),
26                                                    gdof_p );
27                M_p_to_s[gdof_p] = gdof_s;
28            }
29        }
30    }
31 }

```

---

- `int parallelToSequential( size_type const& p_dof )`: returns the sequential DoF corresponding to the parallel Dof  $p\_dof$ .
- `std::pair<int,int> sequentialToParallel( size_type const& s_dof )`: returns the couple  $(n_p, \text{DoF}_p)$  corresponding to the sequential Dof  $s\_dof$ .
- `void project( element_type& us, element_type const& up )`: projects a parallel element from  $X_p$  on  $X_s$ . This function also gathers the contributions from all processors, this procedure is detailed in the snippet 10.2

This map and those methods are used in DEIM to extract the values of  $T_{I_M}(u, \mu)$  and to project the reduced basis on the interpolation space (for non-linear problems only).

**Listing 10.2** – Gathering of the contributions from all processors after the projection of a parallel element on the sequential space.

---

```
1 // loop on the elements of the full mesh splitted between all procs
2 template <typename ElementType>
3 void gather( ElementType& us )
4 {
5     int world_size = Environment::worldComm().globalSize();
6     auto Rh = us.functionSpace();
7     auto ut = Rh->element();
8     std::vector<ElementType> all_u( world_size, ut );
9
10    mpi::all_gather( Environment::worldComm().globalComm(),
11                   us,
12                   all_u );
13
14    for ( auto ut : all_u )
15    {
16        for ( int j=0; j<us.size(); j++ )
17        {
18            if ( ut(j)!=0 )
19                us(j) = ut(j);
20        }
21    }
22 }
```

---

### Conclusion

We started this thesis with the ambition to develop a reduced model for aerothermal simulations. We divided this work in two natural parts.

In the first part, we exposed the numerical methods implemented for modeling aerothermal flows. We considered the incompressible Navier-Stokes equations coupled with an energy equation in the Boussinesq approximation. We used a standard Finite Element discretization with a Galerkin projection to numerically solve this system of PDEs. We introduced a Newton method to solve the non-linear discrete system. To improve the convergence capability of this iterative method, we also introduced a physical and a pseudo-transient continuation. These tools allow to efficiently reach the stationary state of the system if it exists.

We also implemented different versions of the Streamline Diffusion Method (SDM) to stabilize the numerical resolution. These methods have two aims: the first one is to deal with non-physical oscillations appearing in the resolution of convection dominated problems. The second objective is the stabilization of the numerical solution in the regions of anisotropic meshing. This kind of mesh is used to model the boundary layers of turbulent flows. To fulfill these two goals, we compared the efficiency of two SDM formulations: the Streamline-Upwind Petrov Galerkin (SUPG) and the Galerkin Least Square (GLS). These methods have been studied with different expressions for the stabilization parameter. We paid a special attention to the evaluation of the characteristic size of the cells for anisotropic meshes and we proposed an efficient formulation for simplex with large aspect ratios. The implemented methods were validated on multiple numerical applications.

The objective of the second part of this thesis was the reduction of the model described previously.

The model order reduction method studied was the Reduced Basis Method (RBM). This method is particularly adapted for the reduction of problems discretized with the FEM. It provides an efficient offline/online decomposition and a rapid convergence when used with efficient error estimators. We presented in this report our implementation of the RBM for saddle-point and trilinear problems. The numerical results are very encouraging for the reduced aerothermal model. We obtained a nice convergence of the reduced solutions although we did not have any error estimator. However, we noticed the limit of our model for complex geometries with important variation of the flow. For these kind of problems, our reduced model was not able to provide accurate solutions.

In order to reduce the stabilization operators presented in the first part, we proposed to use the Empirical Interpolation Method (EIM). We implemented a discrete version of this algorithm, allowing to recover an affine decomposition of parametrized matrices and vectors. This method was successfully applied to obtain an affine approximation of the Jacobian and the Residual used for the resolution of stabilized aerothermal problems with a Newton algorithm. We also presented a discrete version of the Simultaneous EIM and RB (SER) algorithm using this new discrete version of EIM. Using SER for the reduction of non-linear problems greatly reduced the cost of the offline procedure by replacing FEM approximations with RB approximation in the EIM greedy algorithm.

The reduction of the turbulence models, described in the first part, was not tackled by lack of time. Nevertheless, we exposed some ideas for the development of reduced RANS model involving the construction of a reduced space for the turbulent quantities.

We eventually provided a reduced aerothermal model. However, this model still needs important improvements.

## Outlook

Apart from the issue with the turbulent model, the aerothermal framework is finally robust and efficient. Most of the future developments should concern the reduction method.

The first priority would concern the optimization of the online resolution. We highlighted two major issues with our numerical results. First, the assembly of the reduced operators becomes abnormally expensive when the dimension of the reduced spaces is growing. This issue leads to excessive online runtimes. This issue should be fixed by optimizing the data structure used for storing the reduced

operators. The second important issue concerns the assembly of the reduced tensors for the discrete EIM. Despite the use of a reduced interpolation space, this procedure is relatively expensive and is limiting the capabilities of the reduced models. Some optimizations are still possible. For instance, we could store the evaluation of the FE basis functions on all the quadrature points of the interpolation mesh.

Another important potential improvement is the implementation of error estimators for the aerothermal model. We observed some convergence issues in the numerical results which could be solved with a better selection of the parameters in the greedy algorithm. However, this point remains non-trivial and requires important theoretical and numerical efforts.

Finally, the ultimate outlook would be the complete development of a reduced turbulent model. This research subject is still very open in the community of the model order reduction and we hope that the contribution in this thesis can help somehow to advance on this problem.



# **Appendices**





# 1 Adaptive Time Stepping and Benchmarks

As an alternative to the standard BDF time discretization, we implemented an adaptive time-stepping algorithm proposed by Kay and Al. in [54, 38]. We recall in this section the main ingredients of the methods and we give some results.

Considering the results on the 2D Turek Benchmark and the difficulty to scale our PCD preconditioner with this formulation, the method does not actually seem to be adapted to our problem. We nevertheless dedicated a section to this algorithm which was very interesting to study and implement.

We consider the problem described in chapter 2. We will not detail the formulation for the whole system but only for the Navier-Stokes part. The method is similar for the energy equation and we give the final equations in section 1.1

We introduce a discretization of the time interval  $[0, t_f]$ ,  $\{t_n\}_{n=0,\dots,N}$  with  $t_0 = 0$  and  $t_N = t_f$ . We then define the time steps  $\{k_n\}_{n=1,\dots,N}$  by  $k_{n+1} = t_{n+1} - t_n$ . We denote by  $(\mathbf{u}^n, p^n, T^n)$  the velocity, pressure and temperature fields at iteration  $n$ .

## 1.1 Detailed formulation for Navier-Stokes

We consider the fluid equation in the system (63). We first decide to linearise the equation using an Oseen formulation with extrapolation of the convection velocity

$$\mathbf{w}^{n+1} = \left(1 + \frac{k_{n+1}}{k_n}\right) \mathbf{u}^n - \frac{k_{n+1}}{k_n} \mathbf{u}^{n-1}. \quad (1)$$

We choose a Crank Nicholson scheme for the time discretization. We find different propositions for the choice of the scheme in the literature but Crank Nicholson seems to produce the best results with this kind of time stepping algorithm. And so the weak formulation becomes

$$2\rho\left(\frac{\mathbf{u}^{n+1}}{k_{n+1}}, \mathbf{v}\right) + \rho(\mathbf{w}^{n+1} \cdot \nabla \mathbf{u}^{n+1}, \mathbf{v}) + \mu(\nabla \mathbf{u}^{n+1}, \nabla \mathbf{v}) - (p^{n+1}, \nabla \cdot \mathbf{v}) - (q, \nabla \cdot \mathbf{u}^{n+1}) = \quad (2)$$

$$2\rho\left(\frac{\mathbf{u}^n}{k_{n+1}}, \mathbf{v}\right) + \left(\frac{\partial \mathbf{u}^n}{\partial t}, \mathbf{v}\right) + (\mathbf{f}^n, \mathbf{v})$$

where  $f^n = \beta(T^n - T_0)\mathbf{g}$ .

We introduce now the discrete acceleration given by  $\mathbf{d}^n = \frac{\mathbf{u}^{n+1} - \mathbf{u}^n}{k_{n+1}}$ ,  $\mathbf{u}^{n+1} = \mathbf{u}^n + k_{n+1}\mathbf{d}^n$  and we rewrite the equation (2) with this new quantity,

$$\begin{aligned} \rho(\mathbf{d}^n, \mathbf{v}) + \rho k_{n+1}(\mathbf{w}^{n+1} \cdot \nabla \mathbf{d}^n, \mathbf{v}) + \mu k_{n+1}(\nabla \mathbf{d}^n, \nabla \mathbf{v}) - (p^{n+1}, \nabla \cdot \mathbf{v}) - (q, \nabla \cdot \mathbf{d}^n) = \\ \mu(\nabla u^n, \nabla \mathbf{v}) - \rho(\mathbf{w}^{n+1} \cdot \nabla u^n, \mathbf{v}) + \frac{1}{k_{n+1}}(q, \nabla \cdot \mathbf{u}^n) + \left(\frac{\partial \mathbf{u}^n}{\partial t}, \mathbf{v}\right) + (f, \mathbf{v}). \end{aligned} \quad (3)$$

We recognize a convection diffusion equation very close to the usual Navier-Stokes weak formulation. But now we want to solve this equation in order to find a couple  $(\mathbf{d}^n, p^{n+1}) \in V_h$ . The solution strategy and especially the preconditioner choice will be discussed later.

## Formulation for the Aerothermal System

With the same method we can rewrite our full system (63), introducing the rate of change of the temperature field  $D^n = \frac{T^{n+1} - T^n}{k_{n+1}}$ , the full aerothermal system now becomes, find  $(\mathbf{d}^n, p^{n+1}, D^n) \in X_h$  such that

$$\begin{aligned} \rho(\mathbf{d}^n, \mathbf{v}) + \rho k_{n+1}(\mathbf{w}^{n+1} \cdot \nabla \mathbf{d}^n, \mathbf{v}) + \mu k_{n+1}(\nabla \mathbf{d}^n, \nabla \mathbf{v}) - (p^{n+1}, \nabla \cdot \mathbf{v}) - (q, \nabla \cdot \mathbf{d}^n) = \\ (f^n, \mathbf{v}) - \mu(\nabla u^n, \nabla \mathbf{v}) - \rho(\mathbf{w}^{n+1} \cdot \nabla u^n, \mathbf{v}) + \frac{1}{k_{n+1}}(q, \nabla \cdot \mathbf{u}^n) + \left(\frac{\partial \mathbf{u}^n}{\partial t}, \mathbf{v}\right) \end{aligned}$$

$$(D^n, S) + \kappa k_{n+1}(\nabla D^n, \nabla S) + k_{n+1}(\mathbf{w}^{n+1} D^n, S) = -\kappa(\nabla T^n, \nabla S) - (\mathbf{w}^{n+1} T^n, S) + \left(\frac{\partial T^n}{\partial t}, S\right),$$

$$\forall (\mathbf{v}, q, S) \in X_h \quad (4)$$

Once we have solved the system (4) we can evaluate the different quantities required for the next step :

$$\mathbf{u}^{n+1} = \mathbf{u}^n + k_{n+1} \mathbf{d}^n \quad (5)$$

$$\frac{\partial \mathbf{u}^{n+1}}{\partial t} = 2\mathbf{d}^n - \frac{\partial \mathbf{u}^n}{\partial t} \quad (6)$$

$$T^{n+1} = T^n + k_{n+1} D^n \quad (7)$$

$$\frac{\partial T^{n+1}}{\partial t} = 2D^n - \frac{\partial T^n}{\partial t} \quad (8)$$

## 1.2 Time step computation

The choice of the time step is based on comparison of two solution : the first one given by our resolution of (4) after the update (5), the second solution is given by an explicit evaluation using a second order Adams-Baschforth scheme. This comparison is made for both velocity and temperature field.

First the AB solution are computed

$$\mathbf{u}_{AB2}^{n+1} = \mathbf{u}^n + \frac{k_{n+1}}{2} \left[ \left(1 + k_{n+1}/k_n\right) \frac{\partial \mathbf{u}^n}{\partial t} + \frac{k_{n+1}}{k_n} \frac{\partial \mathbf{u}^{n-1}}{\partial t} \right], \quad (9)$$

$$T_{AB2}^{n+1} = T^n + \frac{k_{n+1}}{2} \left[ \left(1 + k_{n+1}/k_n\right) \frac{\partial T^n}{\partial t} + \frac{k_{n+1}}{k_n} \frac{\partial T^{n-1}}{\partial t} \right]. \quad (10)$$

The errors between the solution and the AB extrapolation are computed using the formula

$$e_u^{n+1} = \frac{\|\mathbf{u}^{n+1} - \mathbf{u}_{AB2}^{n+1}\|}{3(1 + k_{n+1}/k_n)}, \quad e_T^{n+1} = \frac{\|T^{n+1} - T_{AB2}^{n+1}\|}{3(1 + k_{n+1}/k_n)} \quad (11)$$

and then used to evaluate the next time step, taking in account the error on both velocity and temperature fields.

$$k_{n+2} = k_{n+1} \left( \frac{\varepsilon}{\sqrt{(e_T^{n+1})^2 + (e_u^{n+1})^2}} \right)^{1/3}. \quad (12)$$

Here the parameter  $\varepsilon$  is a tolerance used to act on the time step evolution. If this tolerance is too big the time step will grow too quickly and one will lose in accuracy. But with too small tolerance, the time step will not take off and then one loses the interest of this algorithm. We typically use a tolerance  $\varepsilon$  between  $10^{-2}$  and  $10^{-6}$ , this choice greatly depends of the physics of your problem.

### Initialization of the algorithm

Since the algorithm is not self-starting, one has to compute acceptable initial state for the different fields of interest  $(\mathbf{u}, p, T)$ , for the discrete acceleration and for rate of change.

The first time step is computed by the resolution of at time  $t_1 = k_1$ . In practice we use very small value of  $k_1$  ( $10^{-8}$  for instance). We start by the computation of  $(\mathbf{u}^1, p^1, T^1)$  such that

$$\begin{aligned} \mu(\nabla \mathbf{u}^1, \nabla \mathbf{v}) - (p^1, \nabla \cdot \mathbf{v}) - (q, \nabla \cdot \mathbf{u}^1) &= 0, \\ \kappa(\nabla T^1, \nabla S) &= 0, \quad \forall (\mathbf{v}, q, T) \in V_h \times X_h. \end{aligned} \quad (13)$$

These initial solutions are then used to compute initial acceleration  $d^1$  and rate of change  $D^1$  such that

$$\begin{aligned} \rho(d^1, \mathbf{v}) - (p^1, \nabla \cdot \mathbf{v}) - (q, \nabla \cdot d^1) &= -\mu(\nabla \mathbf{u}^1, \nabla \mathbf{v}) - \rho(\mathbf{u}^1 \cdot \nabla \mathbf{u}^1, \mathbf{v}) \\ \left(\frac{\partial D^1}{\partial t}, S\right) &= -\kappa(\nabla T^1, S) - (\mathbf{u}^1 \cdot \nabla T^1, S), \quad \forall (\mathbf{v}, q, S) \in V_h \times X_h \end{aligned} \quad (14)$$

For the second time step we first have to find  $(u^2, p^2)$  and  $T^2$  such that

$$\begin{aligned} \frac{2}{k_2}(\mathbf{u}^2, \mathbf{v}) + \mu(\nabla \mathbf{u}^2, \nabla \mathbf{v}) + (\mathbf{w}^2 \cdot \nabla \mathbf{u}^2, \mathbf{v}) - (p^2, \nabla \cdot \mathbf{v}) - (q, \nabla \cdot \mathbf{u}^2) &= \frac{2}{k_2}(\mathbf{u}^1, \mathbf{v}) + (\mathbf{d}^1, \mathbf{v}), \\ \frac{2}{k_2}(T^2, S) + \kappa(\nabla T^2, \nabla S) + (\mathbf{w}^2 \cdot \nabla T^2, S) &= \frac{2}{k_2}(T^1, S) + (D^1, S), \\ \forall(\mathbf{v}, q, S) \in V_h \times X_h. & \end{aligned} \quad (15)$$

This allows to evaluate  $d^2$  and  $D^2$  with the formula

$$\begin{aligned} d^2 &= \frac{2}{k_2}(\mathbf{u}^2 - \mathbf{u}^1) - \mathbf{d}^1 \\ D^2 &= \frac{2}{k_2}(T^2 - T^1) - D^1. \end{aligned} \quad (16)$$

Since we now have the acceleration and the rate of change at time  $t^2 = 2k_1$ , we can start the algorithm. The value of  $\mathbf{u}^2$  and  $T^2$  will be replaced by using the update (4). And then the algorithm can begin with the resolution of (5) for the next time step.

### Averaging algorithm

The method introduced in this section is prone to “ringing” effects especially with large tolerance  $\varepsilon$  or when we are close to a steady state. More details about this behavior can be found in [54, 38].

To avoid these oscillations effects, we chose to implement an averaging algorithm as proposed in [54]. This algorithm is invoked periodically, every  $n^*$  steps. For such a step, we first compute the next time step  $k_{n+1}$ , using the method described in section 1.2. We also compute the current velocity  $\mathbf{u}^n$  and the current acceleration  $\frac{\partial \mathbf{u}^n}{\partial t}$  with the formula (4).

The time step is then modified as follows

$$t_n = \frac{t_{n-1} + t^*}{2}, \quad k_n = \frac{k^*}{2} \quad (17)$$

and we can compute the new value for the fields  $\mathbf{u}^n$  and  $\frac{\partial \mathbf{u}^n}{\partial t}$

$$\mathbf{u}^n = \frac{\mathbf{u}^{n-1} + \mathbf{u}^*}{2}, \quad \frac{\partial \mathbf{u}^n}{\partial t} = \frac{1}{2} \left( \frac{\partial \mathbf{u}^{n-1}}{\partial t} + \frac{\partial \mathbf{u}^*}{\partial t} \right). \quad (18)$$

with

$$t^* = t_n, \quad k^* = k_n, \quad \mathbf{u}^* = \mathbf{u}^n, \quad \frac{\partial \mathbf{u}^*}{\partial t} = \frac{\partial \mathbf{u}^n}{\partial t}. \quad (19)$$

The value of  $k_{n+1}$  choose previously keep unchanged and we can solve the next time step at  $t_{n+1} = t_n + k_{n+1}$ . We gave here the details of the averaging for the fluid part but the process is exactly the same for the thermal part of the system.

### 1.3 Turek Benchmarks

These benchmarks are proposed on <http://www.featflow.de>. Here we briefly recall the main parameters of the model. We then present our results obtained with our implementation of the framework presented in chapter 2.

The results on these two benchmarks were significant to evaluate the efficiency of our implementation.

#### Configuration

The parameters of the fluid are  $\rho = 1.0$  and  $\mu = 0.001$ . The geometry of the model is a pipe without a circular cylinder of radius  $r = 0.05$ , as described in figure 2

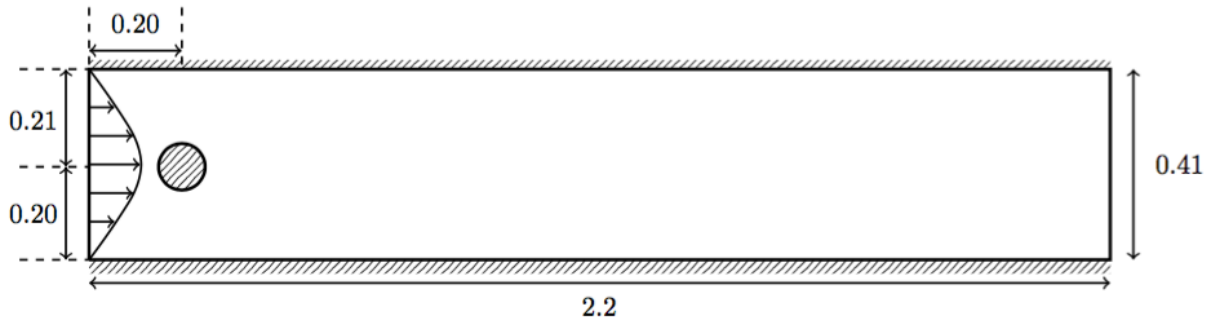


Figure 2 – Geometry of the Problem

No-slip boundary conditions are imposed on upper and lower walls and on the boundary of the cylinder. The left edge is the input of the pipe where a Poiseuille profile is prescribed,

$$u(0, y) = \left( \frac{4U_y(0.41 - y)}{0.41^2}, 0 \right) \quad (20)$$

with the maximum amplitude of the velocity  $U_y = 0.3$ . The right edge is the output of the pipe and with a do-nothing condition

$$\bar{\sigma} = 0. \quad (21)$$

With this configuration, the Reynolds number is  $Re = 20$  and the system turns into stationary state.

#### Quantities of interest

In this benchmark we are interested in three quantities. All this quantities are measured in the steady state.

First we want to evaluate the drag and lift coefficients, respectively  $C_D$  and  $C_L$ , defined by

$$C_D = \frac{2}{\bar{U}^2 L} F_D \quad C_L = \frac{2}{\bar{U}^2 L} F_L \quad (22)$$

with  $\bar{U} = 0.2$  the mean velocity at the input,  $L = 0.1$  the characteristic length of the flow configuration and  $(F_D, F_L)$  such that

$$(F_D, F_L) = \int_{\Gamma_C} \bar{\sigma} n d\Gamma \quad (23)$$

where  $\Gamma_C$  is the boundary of the cylinder and  $n$  is the outer normal vector of the circle.

The last quantity of interest is the pressure difference  $\Delta_p = p(a_1) - p(a_2)$  between the points  $a_1 = (0.15, 0.2)$  and  $a_2 = (0.25, 0.2)$ .

We will compare our results with the reference which gives  $C_D = 5.57953523384$ ,  $C_L = 0.010618948146$  and  $\Delta_p = 0.11752016697$ .

## Mesh Convergence

We first try to find the best compromise between accuracy and computation cost for our mesh. We chose to use a mesh with 2 characteristic sizes: a default size  $h$  for the whole domain and a smaller size  $h_c$  applied on the cylinder. It results in a mesh particularly refined around the cylinder. An example of such a mesh is proposed in figure 3.

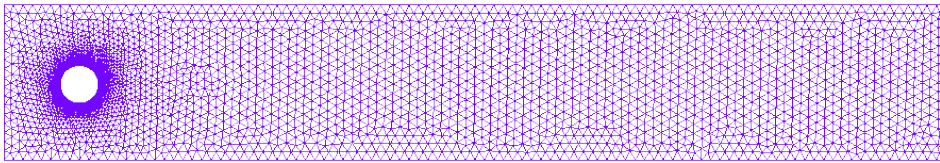


Figure 3 – Mesh used for the simulation

To find the best values for  $h$  and  $h_c$  we ran many simulation and compare the results for  $C_D$ ,  $C_L$  and  $\Delta_p$  with the reference values for different value of  $(h, h_c)$ . The errors in percent, with respect to the reference results, are reported in tables 4, 5 and 6.

## Reference Case

We are not only interested in the validation of our algorithms for the steady state but also for the transient state. For that we start by running a simulation with a

$h \backslash h_c$	0.1	0.05	0.025	0.01	0.005	0.00025	0.001
0.1	9.72	9.72	5.14	1.69	0.60	0.26	0.17
0.05		9.20	4.94	1.51	0.48	0.17	0.06
0.025			4.51	1.54	9.47	0.16	0.05

Figure 4 –  $C_D$  error for different values of  $h$  and  $h_c$  (in %)

$h \backslash h_c$	0.1	0.05	0.025	0.01	0.005	0.00025	0.001
0.1	407.17	407.17	71.95	149.20	63.77	11.58	77.98
0.05		2118.51	160.67	1.00	12.28	1.83	3.15
0.025			71.02	13.18	12.12	0.43	0.09

Figure 5 –  $C_L$  error for different values of  $h$  and  $h_c$  (in %)

very small time step. The result will be our reference to know if we are accurate enough in the transient state.

This reference simulation is ran on the mesh described in section 1.3 with a time step  $\delta t = 0.01$  and an order 2 BDF formulation.

The results over time for  $C_D$ ,  $C_L$  and  $\Delta_p$  are presented in the figure 7. We consider that the system is in a steady state once the difference  $\frac{\|u^{n+1}-u^n\|}{\|u^n\|}$  is smaller than a chosen tolerance  $\varepsilon_{\text{steady}}$ . For this simulation we choose  $\varepsilon_{\text{steady}} = 10^{-5}$  since we got very good results for all the quantities of interest with this tolerance.

Considering these results and the results from the mesh convergence study, it appears that the lift coefficient ( $C_L$ ) is the more sensitive and so we decided to focus on this quantity to find the optimal time-step (BDF) and tolerance (CNAB2) for this simulation.

### BDF : time step choice

We tried 3 time steps for BDF algorithm,  $\delta t = 0.1, 0.25, 0.5$ . The results are summarized in figure 8 . If both time step  $\delta t = 0.05$  and  $\delta t = 0.1$  give acceptable results, it is clear that the transient state is not well modeled with bigger time step. However we note that the system still converge to the steady state with  $\delta t = 0.25$  but it takes much more time.

### CNAB2 : configuration choice

For CNAB2 algorithm, we tried 3 different tolerances,  $\varepsilon = 10^{-5}, Re^{-5}, 10^{-4}$ . But we had some issue with the growing time step and the CN scheme. When time step becomes to big (typically when you are close to the steady state), the scheme loses in accuracy and generates some non-physics oscillations. To avoid theses oscillations



$h \setminus h_c$	0.1	0.05	0.025	0.01	0.005	0.00025	0.001
0.1	1.33	1.33	0.21	1.07	0.88	0.71	0.61
0.05		0.74	0.28	0.10	0.14	0.06	0.05
0.025			1.04	0.03	0.02	0.02	0.01

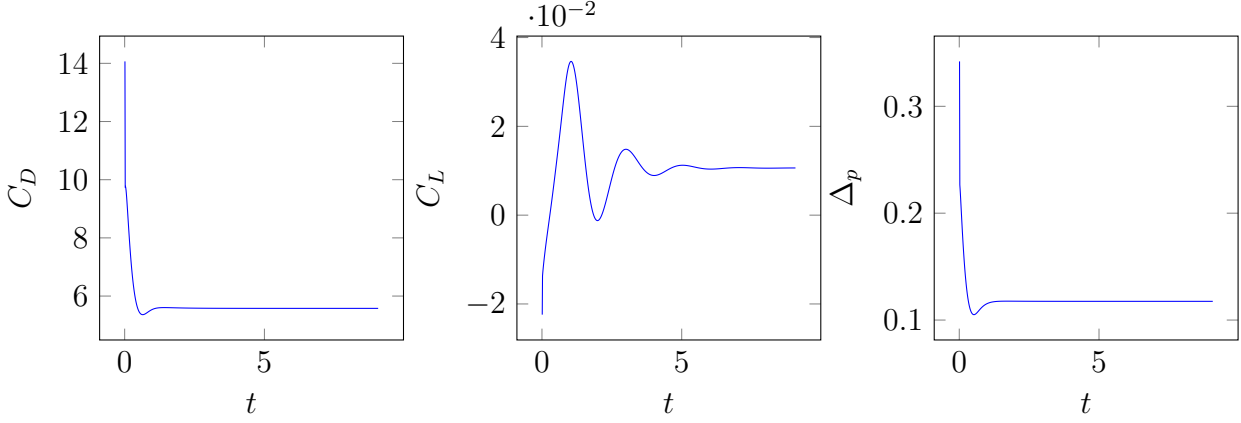
Figure 6 –  $\Delta_p$  error for different values of  $h$  and  $h_c$  (in %)

Figure 7 – Evolution of the 3 quantities of interest over time for the reference case

we had to use a fixed algorithm in the temporal iterations. We iterate until the quantity  $\frac{\|u^{n+1,i+1} - u^{n+1,i}\|}{\|u^{n+1,i}\|}$  is smaller than the tolerance  $\varepsilon_{\text{pic}}$  where  $u^{n+1,i}$  is the solution after the  $i$ -th fixed point iteration. We fix  $u^{n+1,0} = w^{n+1}$  as defined in (1).

We tried 2 values for this tolerance,  $\varepsilon_{\text{pic}} = 10^{-4}, 10^{-5}$ . The results are presented figures 9 and 10

We finally chose to keep the configuration with  $\varepsilon = 10^{-4}$  which was the most efficient in terms of simulation time and number of resolutions.

## Numerical Results

In the table 11 are our results and the error compared to the reference. We also registered the total time of the simulation  $t_{\text{simu}}$ , the number of time-steps  $n_{\text{steps}}$  and the number of resolutions of the Naviers-Stokes system  $n_{\text{NS}}$ .

We ran all theses simulations on 4 cores. The number of DoF is  $37797 \times 4862$ . The BDF case correspond to a usual second order backward differentiation formula with a constant time step chosen as described in section 1.3. The CNAB2 case corresponds to the adaptive time stepping algorithm presented in chapter 2 with the configuration chosen previously, in section 1.3, with  $n^* = 10$ .

In both cases we iterate in time until the quantity  $\Delta_u$  becomes smaller than a

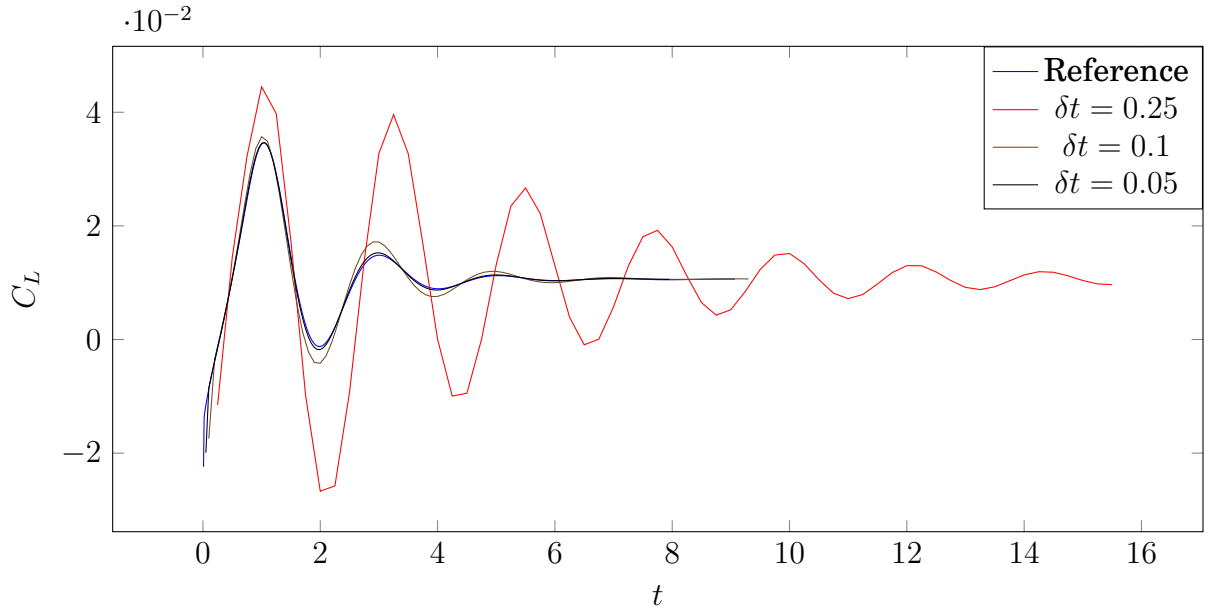


Figure 8 – Comparison of the lift coefficient for 3 different time steps

tolerance  $\varepsilon_{\text{steady}}$ . In practice we choose  $\varepsilon_{\text{steady}} = 10^{-4}$  and  $\Delta_u$  is defined as

$$\Delta_u = \frac{\|u^{n+1} - u^n\|}{\|u^n\|}. \quad (24)$$

When using iterative solver (PCD), we typically set the tolerance for the relative error between two KSP iterations to  $\text{rtol} = 10^{-9}$ .

### **Turek 2D**, $Re = 100$

The benchmark configuration is exactly the same than in section 1.3. We consider the same equations, fluid parameters and geometry, see figure 2.

The only parameter which is different in this configuration is the inflow velocity profile. We still consider a parabolic profile,

$$u(0, y) = \left( \frac{4U_y(0.41 - y)}{0.41^2}, 0 \right) \quad (25)$$

but impose this time a maximum velocity  $U_y = 1.5$ . and so the mean velocity is  $\bar{U} = 1$  and it leads to a Reynolds number  $Re = 100$ .

For this Reynolds number, the flow turns into a time-periodic behavior with a vortex shedding behind the cylinder. See <http://www.featflow.de> for more details.

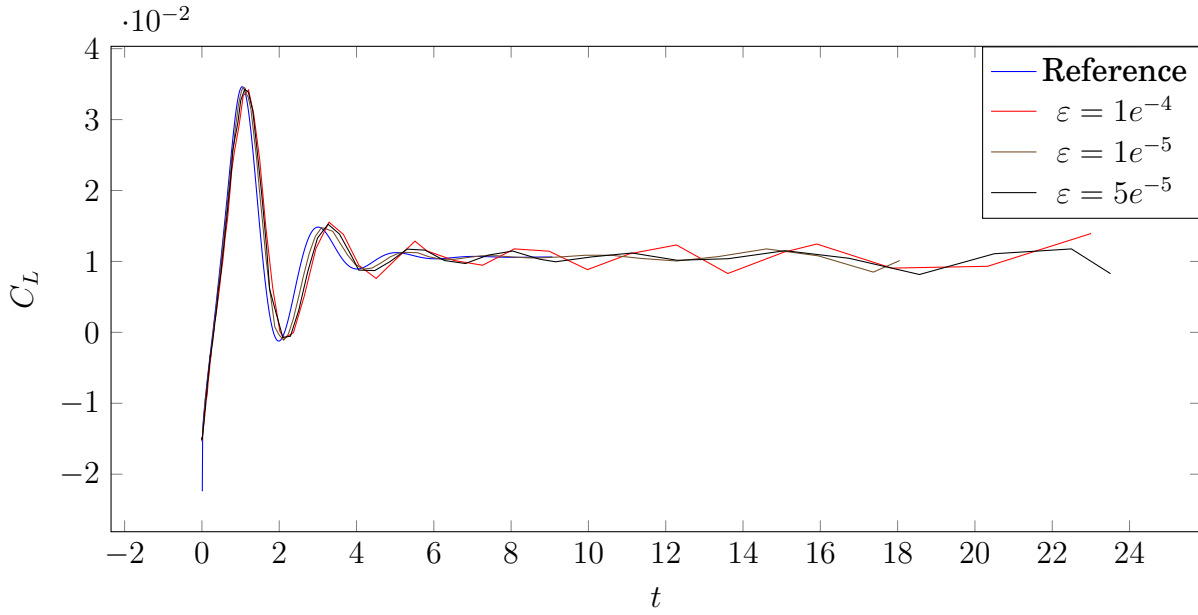


Figure 9 – Comparison of the lift coefficient for 3 different tolerances  $\varepsilon$  with a Picard tolerance  $\varepsilon = 10^{-4}$

## Data Measurement

The measure of the following data is made during a fully developed time-periodic solution. For that we wait for 25 seconds simulation, at this time the flow is fully developed. Then we take an arbitrary cycle between 25 and 30 seconds of simulation times to measure numerical data.

We are first interested in the value of  $C_L$ ,  $C_D$  and  $\Delta_p$  as defined in section 4.1. For these quantities we will compare the max, the min, the amplitude and the mean value against the reference results.

We also define  $f$  the frequency of the periodic flow. If a cycle starts at time  $t_0$  where the lift coefficient  $C_L$  is the smallest, it ends at time  $t_1 = t_0 + 1/f$  when  $C_L$  is smallest again.  $1/f$  is then the length of the cycle. With this definition of the frequency, we can define the Strouhal number

$$\text{St} = \frac{Lf}{\overline{U}} \quad (26)$$

Our measure for  $f$  and  $\text{St}$  will also be compared with the reference data.

For this benchmark we only used direct solver with LU resolution. PCD may be validated later on 3D benchmark. Our results are presented in table 12 for different meshes, time step or tolerance. In red are the non acceptable values of the measure with respect to the reference bounds recalled in table 13.

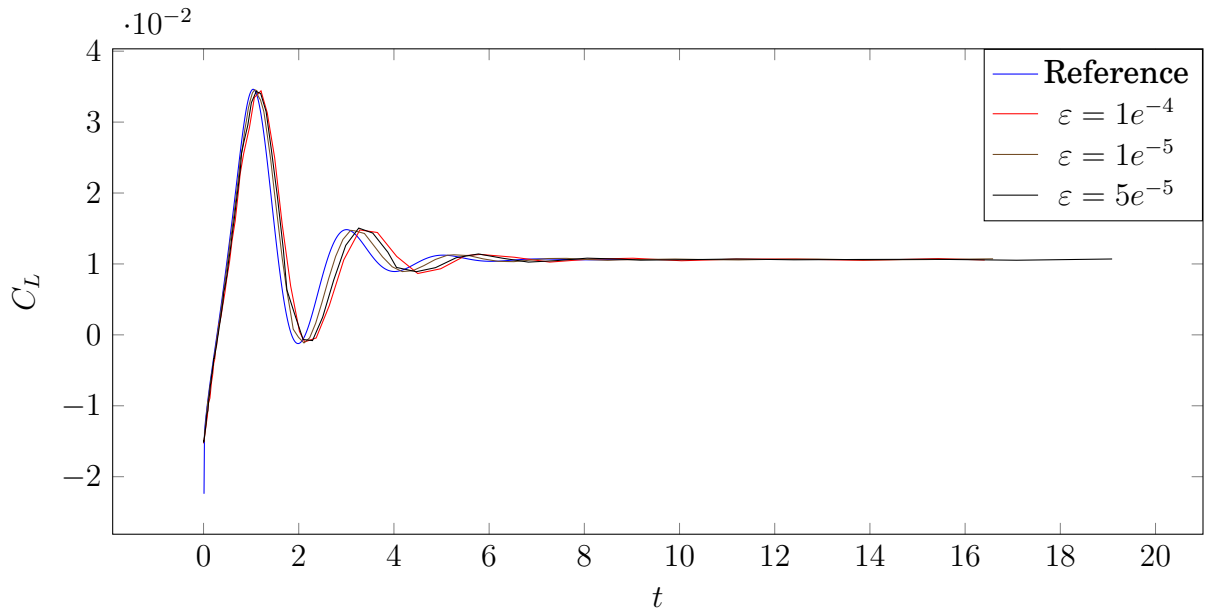


Figure 10 – Comparison of the lift coefficient for 3 different tolerances  $\varepsilon$  with a picard tolerance  $\varepsilon = 10^{-5}$

	BDF				CNAB2			
	LU		PCD		LU		PCD	
$C_D$	5.57694	[0.047]	5.57694	[0.047]	5.57694	[0.047]	5.56076	[0.33]
$C_L$	0.0106564	[0.35]	0.0106567	[0.35]	0.01056	[0.35]	0.0109929	[3.52]
$\Delta_p$	0.117512	[0.007]	0.117512	[0.007]	0.117512	[0.007]	0.11711	[0.34]
$t_{\text{simu}}$	124s		242s		184s		546s	
$n_{\text{step}}$	93		93		61		61	
$n_{\text{NS}}$	93		93		147		152	

Figure 11 – Summary of benchmark results for different configurations, [ error w.r.t. reference in %]

	BDF					CNAB2		
	$\delta t = 0.01$			0.005	0.0025	$\varepsilon = 10^{-4}$	$10^{-5}$	$10^{-6}$
$h$	0.025	0.01	0.005	0.01	0.01	0.01	0.01	0.01
$C_{D\text{max}}$	<b>3.2699</b>	<b>3.2722</b>	<b>3.2724</b>	3.2389	3.2302	3.2333	3.2273	3.2268
$C_{L\text{max}}$	<b>1.0655</b>	<b>1.0666</b>	<b>1.0666</b>	1.0071	0.99173	0.99391	0.98453	0.98473
$\Delta_p$	<b>2.50229</b>	<b>2.50718</b>	<b>2.5074</b>	2.49234	2.48609	2.43213	2.48715	2.48501
St	<b>0.31250</b>	0.30303	0.30303	0.30303	0.30303	0.30239	0.30120	0.30130
$n_{\text{iter}}$	3000	3000	3000	6000	12000	6504	14218	30925
$t_{\text{simu}}$		10100s		19800s		28000s		135000s

Figure 12 – Results for turek 2D benchmark,  $Re = 100$

$C_{D\text{max}}$	3.22 to 3.24
$C_{L\text{max}}$	0.99 to 1.01
$\Delta_p$	2.46 to 2.5
St	0.295 to 0.305

Figure 13 – Reference Bounds for the quantities of interest

## 2 Nitsche's Method for Slip Condition

We consider the system

$$\begin{cases} \mathbf{u} \cdot \nabla \mathbf{u} + \nabla p - \nabla \cdot (\mu \nabla \mathbf{u}) = \mathbf{f}, \\ \nabla \cdot \mathbf{u} = 0, \\ \mathbf{u} \cdot \mathbf{n} = g_n, \quad \text{on } \Gamma \\ (\mathbf{T}(p, \mathbf{u})\mathbf{n})_\tau = \mathbf{g}_\tau, \quad \text{on } \Gamma \end{cases} \quad (27)$$

For any vector  $\mathbf{u}$  we denote by  $(\mathbf{u})_n = (\mathbf{u} \cdot \mathbf{n})\mathbf{n}$  the component of  $\mathbf{u}$  orthogonal to  $\Gamma$  and  $(\mathbf{u})_\tau = \mathbf{u} - (\mathbf{u})_n$  the component of  $\mathbf{u}$  tangential to  $\Gamma$ .

We apply the usual method, detailed in section 2, to recover the following weak discrete formulation,

$$\begin{aligned} \int_{\Omega} (\mathbf{u}_h \cdot \nabla) \mathbf{u}_h \cdot \mathbf{v}_h \, d\Omega + \int_{\Omega} \mu \nabla \mathbf{u}_h : \nabla \mathbf{v}_h \, d\Omega - \int_{\Omega} p_h \nabla \cdot \mathbf{v}_h \, d\Omega - \alpha \int_{\Omega} q_h \nabla \cdot \mathbf{u}_h \, d\Omega \\ - \int_{\Gamma} \mathbf{T}(\mathbf{u}_h, p_h) \mathbf{n} \cdot \mathbf{v}_h \, d\Gamma = \int_{\Omega} \mathbf{f} \cdot \mathbf{v}_h, \quad \forall (q, \mathbf{v}_h) \in V_h \times Q_h. \end{aligned} \quad (28)$$

We rewrite the term on  $\Gamma$  using the normal/tangential decomposition of the vectors

$$\int_{\Gamma} \mathbf{T}(\mathbf{u}_h, p_h) \mathbf{n} \cdot \mathbf{v}_h \, d\Gamma = \int_{\Gamma} (\mathbf{T}(\mathbf{u}_h, p_h) \mathbf{n})_n \cdot (\mathbf{v}_h)_n \, d\Gamma + \int_{\Gamma} (\mathbf{T}(\mathbf{u}_h, p_h) \mathbf{n})_\tau \cdot (\mathbf{v}_h)_\tau \, d\Gamma, \quad (29)$$

and we can now apply the boundary condition on  $(\mathbf{T}(\mathbf{u}_h, p_h))_\tau$  given by (27). We also add the symmetric/antisymmetric term for the normal component (depending of the value of  $\delta$ ) and we finally add a penalization terms to control the value of  $\mathbf{u}_h \cdot \mathbf{n}$ ,

$$\begin{aligned} \int_{\Gamma} \mathbf{T}(\mathbf{u}_h, p_h) \mathbf{n} \cdot \mathbf{v}_h \, d\Gamma = \int_{\Gamma} (\mathbf{T}(\mathbf{u}_h, p_h) \mathbf{n})_n \cdot (\mathbf{v}_h)_n \, d\Gamma + \delta \int_{\Gamma} (\mathbf{T}(\mathbf{v}_h, q_h) \mathbf{n})_n \cdot (\mathbf{u}_h \cdot \mathbf{n} - g_n) \mathbf{n} \, d\Gamma \\ - \gamma \int_{\Gamma} \frac{1}{h_S} (\mathbf{u}_h \cdot \mathbf{n} - g_n) \mathbf{v}_h \cdot \mathbf{n} \, d\Gamma + \int_{\Gamma} \mathbf{g}_\tau \cdot (\mathbf{v}_h)_\tau \, d\Gamma, \end{aligned} \quad (30)$$

where  $\delta = \pm 1$ ,  $h_S$  is the measure of the surface of the current element and  $\gamma$  is a positive constant. The choice of the different constants  $(\alpha, \delta, \gamma)$  and its consequences will be discussed in the following subsection.

Introducing the expression (30) in equation (28), we obtain the new weak formulation

$$\begin{aligned} \int_{\Omega} (\mathbf{u}_h \cdot \nabla) \mathbf{u}_h \cdot \mathbf{v}_h \, d\Omega + \int_{\Omega} \mu \nabla \mathbf{u}_h : \nabla \mathbf{v}_h \, d\Omega - \int_{\Omega} p_h \nabla \cdot \mathbf{v}_h \, d\Omega - \alpha \int_{\Omega} q_h \nabla \cdot \mathbf{u}_h \, d\Omega \\ - \int_{\Gamma} (\mathbf{T}(\mathbf{u}_h, p_h) \mathbf{n})_n \cdot (\mathbf{v}_h)_n \, d\Gamma - \delta \int_{\Gamma} (\mathbf{T}(\mathbf{v}_h, q_h) \mathbf{n})_n \cdot (\mathbf{u}_h)_n \, d\Gamma + \gamma \int_{\Gamma} \frac{1}{h_S} (\mathbf{u}_h \cdot \mathbf{n}) (\mathbf{v}_h \cdot \mathbf{n}) \, d\Gamma \\ = -\delta \int_{\Gamma} (\mathbf{T}(\mathbf{v}_h, q_h) \mathbf{n})_n \cdot g_n \mathbf{n} \, d\Gamma + \gamma \int_{\Gamma} \frac{1}{h_S} g_n \mathbf{v}_h \cdot \mathbf{n} \, d\Gamma + \int_{\Gamma} \mathbf{g}_\tau \cdot (\mathbf{v}_h)_\tau \, d\Gamma + \int_{\Omega} \mathbf{f} \cdot \mathbf{v}_h, \\ \forall (q, \mathbf{v}_h) \in V_h \times Q_h. \end{aligned} \quad (31)$$

## 2.1 Validation

### Convergence Study with Stokes Equations

We first consider a square  $\Omega = [-1, 1] \times [-1, 1]$  with a non structured triangular mesh of characteristic size  $h$ . On this domain, we impose an exact solution to the Stokes equations with  $\mu = 1$ , on all the boundary  $\Gamma$ . As a comparison we first impose this condition with strong Dirichlet (DIR) condition and then using the penalization method for slip condition (SLIP), as written in (31). The exact solution imposed is  $\mathbf{u}_e = (2y(1 - x^2), -2x(1 - y^2))$ . With this solution we have

$$(\mathbf{T}(\mathbf{u}_e, p_e)\mathbf{n})_\tau = \mathbf{u}_e = (2y(1 - x^2), -2x(1 - y^2)) \quad (32)$$

on the boundary of the domain. We then use the formulation (31) with

$$g_n = 0, \quad \mathbf{g}_\tau = (2y(1 - x^2), -2x(1 - y^2)), \quad (33)$$

and the source term  $f$  is chosen in order to have  $p_2 = 0$ . The discrete solution is computed for different values of  $h = 0.2, 0.1, 0.05, 0.025$  and for the 4 different possibilities for the couple  $(\alpha, \delta)$ . We also tested different values of  $\gamma$ . These results are summarized in the figure 14. According to these results, only two configurations have the good convergence rates :  $\alpha = 1, \delta = 1$  and  $\alpha = -1, \delta = -1$ , with a slight accuracy advantage for this last case. In the following convergence studies we will only consider the configurations  $\alpha = \delta$ .

The second convergence test is made on a ring, the domain is defined by  $1 \leq x^2 + y^2 \leq 4$ . The exact solution  $\mathbf{u}_e = (-y\sqrt{x^2 + y^2}, x\sqrt{x^2 + y^2})$  is imposed strongly on the inside boundary ( $x^2 + y^2 = 1$ ) and with slip conditions on the outside boundary ( $x^2 + y^2 = 4$ ). With this exact solution we can evaluate

$$g_n = 0, \quad \mathbf{g}_\tau = (\mathbf{T}(\mathbf{u}_e, p_e)\mathbf{n})_\tau = (-2y, 2x). \quad (34)$$

Again we choose  $p = 0$  and the source term is computed to fit with these parameters.

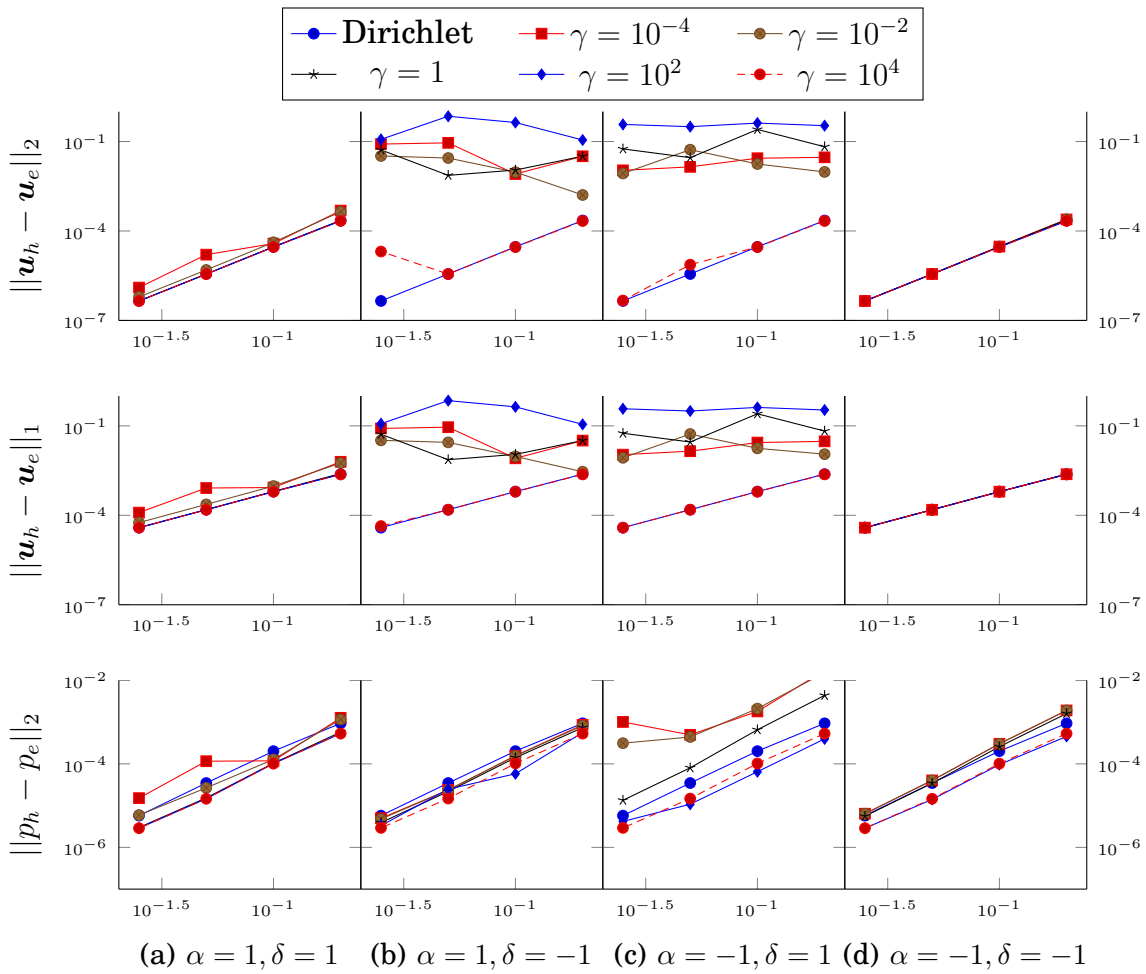


Figure 14 – Convergence study on square test case

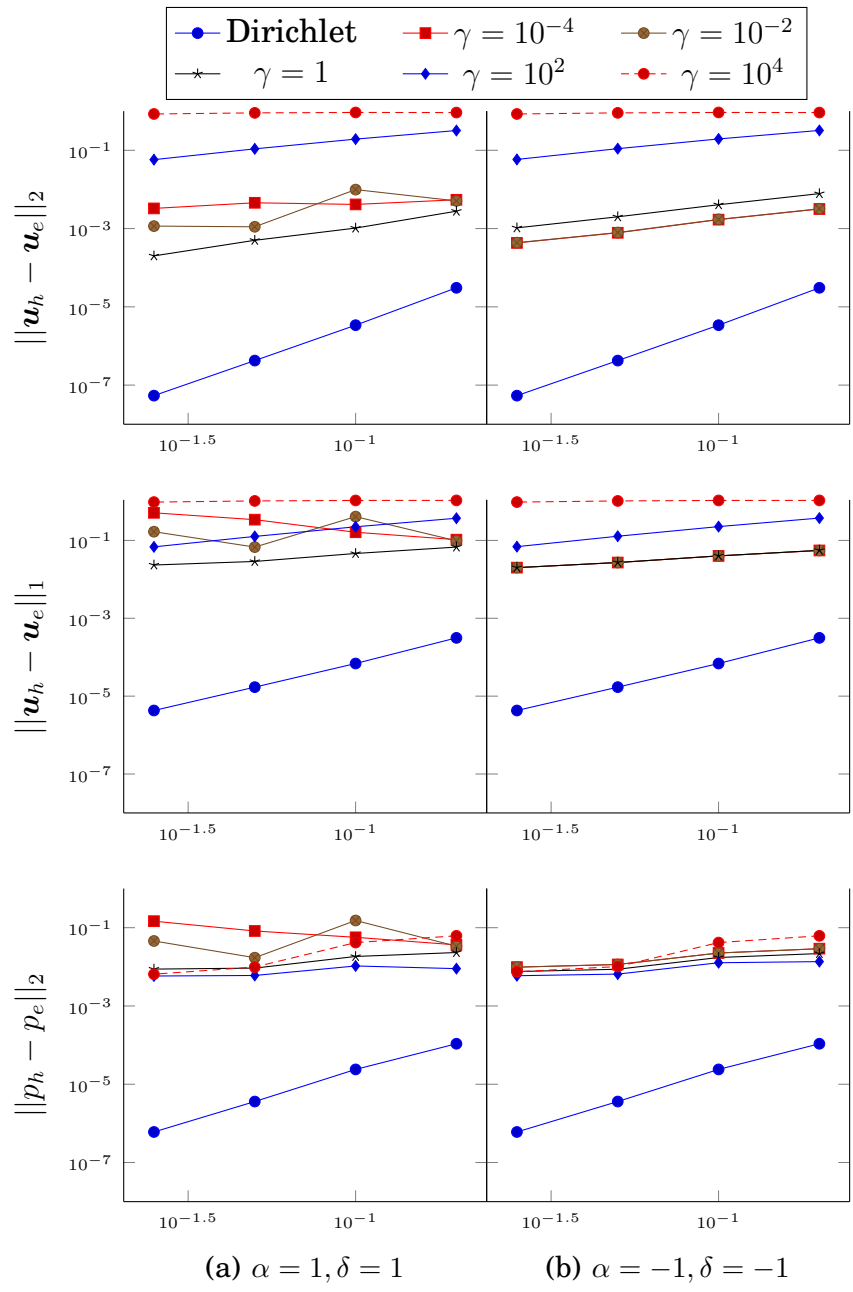


Figure 15 – Convergence study on ring test case



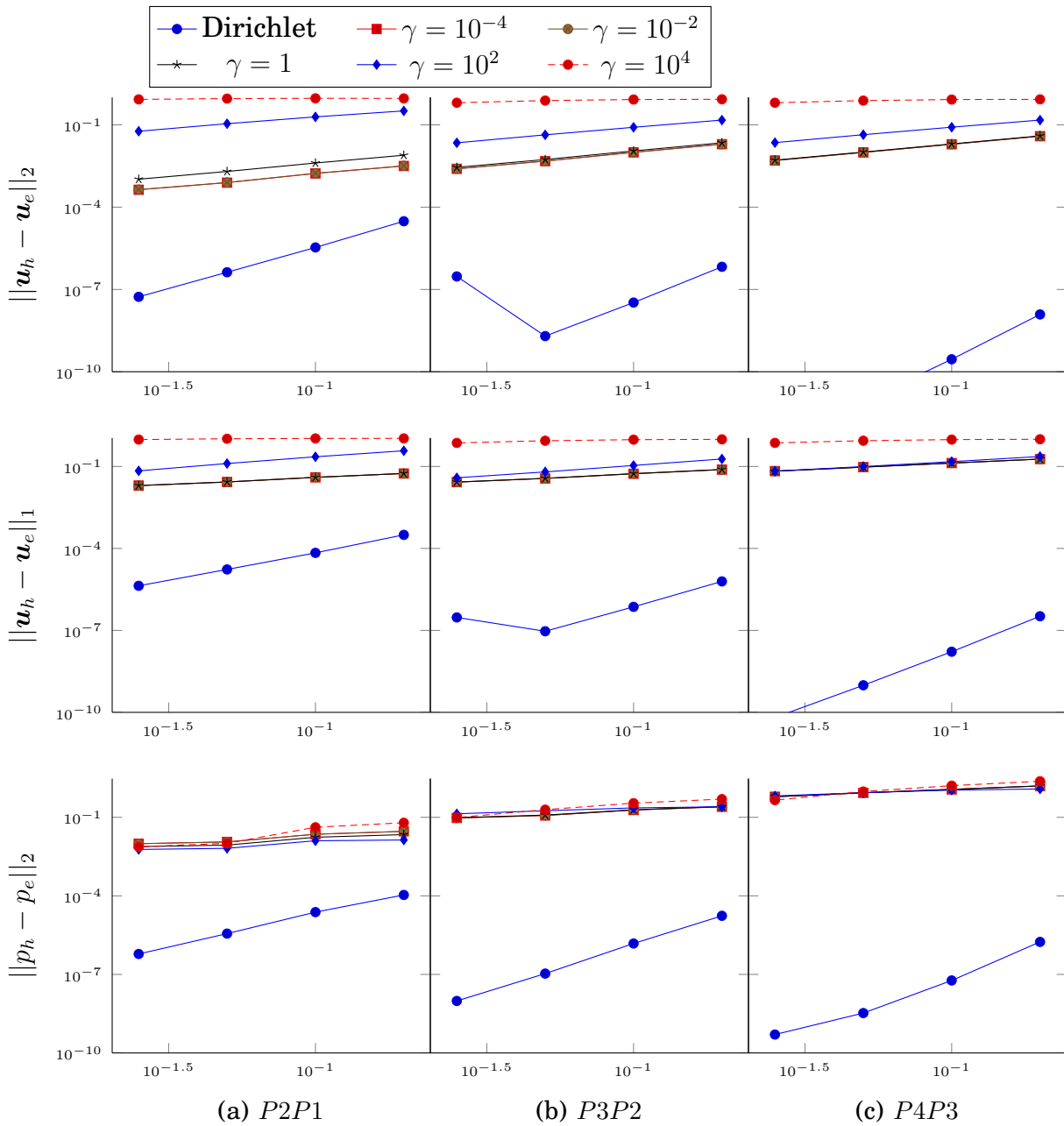


Figure 16 – Convergence study on ring test case G1

### 3 Offline/Online Strategy for Residual Evaluation in the CRB

#### 3.1 Coercive Problems

The residual defined in (5.22) can be rewritten using the affine decomposition of the problem

$$r_N(v; \boldsymbol{\mu}) = \sum_{q=1}^{Q_f} \theta_q^f(\boldsymbol{\mu}) f_q(v) - \sum_{q=1}^{Q_a} \theta_q^a a_q(u_N(\boldsymbol{\mu}); \boldsymbol{\mu}) \quad (35)$$

and then the Riesz representation satisfies

$$(\hat{e}_N(\boldsymbol{\mu}), v)_{\bar{\boldsymbol{\mu}}} = \sum_{q=1}^{Q_f} \theta_q^f(\boldsymbol{\mu}) f_q(v) - \sum_{q=1}^{Q_a} \theta_q^a a_q(u_N(\boldsymbol{\mu}); \boldsymbol{\mu}), \quad \forall v \in X_h. \quad (36)$$

Now using the projection of  $u_N(\boldsymbol{\mu})$  on the RB space, as in (5.12), we obtain

$$(\hat{e}_N(\boldsymbol{\mu}), v)_{\bar{\boldsymbol{\mu}}} = \sum_{q=1}^{Q_f} \theta_q^f(\boldsymbol{\mu}) f_q(v) - \sum_{q=1}^{Q_a} \sum_{i=1}^N \theta_q^a u_{Ni}(\boldsymbol{\mu}) a_q(\xi_i; \boldsymbol{\mu}), \quad \forall v \in X_h. \quad (37)$$

And then, using the linear superposition principle, we can write  $\hat{e}_N(\boldsymbol{\mu})$  as

$$\hat{e}_N(\boldsymbol{\mu}) = \sum_{q=1}^{Q_f} \theta_q^f(\boldsymbol{\mu}) \hat{f}_q - \sum_{q=1}^{Q_a} \sum_{i=1}^N \theta_q^a u_{Ni}(\boldsymbol{\mu}) \hat{a}_{qi} \quad (38)$$

where the Riesz representations  $\hat{f}_q$  and  $\hat{a}_{qi}$  satisfy

$$(\hat{f}_q, v)_{\bar{\boldsymbol{\mu}}} = f_q(v), \quad \forall v \in X_h, \quad 1 \leq q \leq Q_f, \quad (39)$$

$$(\hat{a}_{qi}, v)_{\bar{\boldsymbol{\mu}}} = a_q(\xi_i, v), \quad \forall v \in X_h, \quad 1 \leq i \leq N, \quad 1 \leq q \leq Q_a. \quad (40)$$

Finally we can write

$$\|\hat{e}_N(\boldsymbol{\mu})\|_{\bar{\boldsymbol{\mu}}}^2 = C^{ff}(\boldsymbol{\mu}) + \sum_{i=1}^N \sum_{j=1}^N (u_N(\boldsymbol{\mu}))_i (C_n^{aa}(\boldsymbol{\mu}))_{ij} (u_N(\boldsymbol{\mu}))_j - 2 \sum_{i=1}^N (u_N(\boldsymbol{\mu}))_i (C_N^{fa}(\boldsymbol{\mu}))_i \quad (41)$$

$$= C^{ff}(\boldsymbol{\mu}) + u_N(\boldsymbol{\mu})^\top C_N^{aa}(\boldsymbol{\mu}) u_N(\boldsymbol{\mu}) - 2 u_N(\boldsymbol{\mu})^\top C_N^{fa}(\boldsymbol{\mu}) \quad (42)$$

with,

$$C^{ff}(\boldsymbol{\mu}) = \sum_{q=1}^{Q_f} \sum_{q'=1}^{Q_f} \theta_q^f(\boldsymbol{\mu}) \theta_{q'}^f(\boldsymbol{\mu}) \underbrace{(\hat{f}_q, \hat{f}_{q'})_{\bar{\boldsymbol{\mu}}}}_{\text{precomputed}}, \quad (43)$$

$$C_N^{fa}(\boldsymbol{\mu}) = \sum_{q=1}^{Q_f} \sum_{q'=1}^{Q_a} \theta_q^f(\boldsymbol{\mu}) \theta_{q'}^a(\boldsymbol{\mu}) \underbrace{C_{qq',N}^{fa}}_{\text{precomputed}}, \quad (C_{qq',N}^{fa})_i = (\hat{f}_q, \hat{a}_{q'i})_{\bar{\boldsymbol{\mu}}}, \quad 1 \leq i \leq N, \quad (44)$$

$$C_N^{aa}(\boldsymbol{\mu}) = \sum_{q=1}^{Q_a} \sum_{q'=1}^{Q_a} \theta_q^a(\boldsymbol{\mu}) \theta_{q'}^a(\boldsymbol{\mu}) \underbrace{C_{qq',N}^{aa}}_{\text{precomputed}}, \quad (C_{qq',N}^{aa})_{ij} = (\hat{a}_{qi}, \hat{a}_{q'j})_{\bar{\boldsymbol{\mu}}}, \quad 1 \leq i, j \leq N. \quad (45)$$

The precomputed quantities are evaluated once during the online procedure and allows a efficient online evaluation, independent of the finite element dimension  $\mathcal{N}$ .

Using the same method, we can express Riesz representation norm of the residual for the dual problem as

$$\|\hat{e}_N^{du}(\boldsymbol{\mu})\|_{\bar{\boldsymbol{\mu}}}^2 = C^{ll}(\boldsymbol{\mu}) + \psi_N(\boldsymbol{\mu})^\top C_N^{aa}(\boldsymbol{\mu}) \psi_N(\boldsymbol{\mu}) + 2\psi_N(\boldsymbol{\mu})^\top C_N^{la}(\boldsymbol{\mu}) \quad (46)$$

### 3.2 Saddle-Point Problems

Just like in the coercive case, we start by introducing the Riesz representation of each residual which satisfy

$$\begin{aligned} (\hat{r}_N^u(\boldsymbol{\mu}), \mathbf{v})_V &= r_N^u(\mathbf{v}; \boldsymbol{\mu}), \quad \forall \mathbf{v} \in V_h \\ (\hat{r}_N^p(\boldsymbol{\mu}), q)_Q &= r_N^p(q; \boldsymbol{\mu}), \quad \forall q \in Q_h \end{aligned} \quad (47)$$

Using the affine decomposition (5.39) of the forms we can rewrite the relation (47) as,

$$\begin{aligned} (\hat{r}_N^u(\boldsymbol{\mu}), \mathbf{v})_V &= \sum_{k=1}^{Q_f} \theta_k^f(\boldsymbol{\mu}) f_k(\mathbf{v}) - \sum_{k=1}^{Q_a} \theta_k^a(\boldsymbol{\mu}) a_k(\mathbf{u}_N(\boldsymbol{\mu}), \mathbf{v}) - \sum_{k=1}^{Q_b} \theta_k^b(\boldsymbol{\mu}) b_k(\mathbf{v}, p_N(\boldsymbol{\mu})), \\ (\hat{r}_N^p(\boldsymbol{\mu}), q)_Q &= \sum_{k=1}^{Q_g} \theta_k^g(\boldsymbol{\mu}) g_k(q) - \sum_{k=1}^{Q_b} \theta_k^b(\boldsymbol{\mu}) b_k(\mathbf{u}_N(\boldsymbol{\mu}), q). \end{aligned} \quad (48)$$

Now with the projection on the RB spaces we can write

$$\begin{aligned} (\hat{r}_N^u(\boldsymbol{\mu}), \mathbf{v})_V &= \sum_{k=1}^{Q_f} \theta_k^f(\boldsymbol{\mu}) f_k(\mathbf{v}) - \sum_{k=1}^{Q_a} \sum_{i=1}^{2N} \theta_k^a(\boldsymbol{\mu}) u_{Ni}(\boldsymbol{\mu}) a_k(\boldsymbol{\zeta}_i, \mathbf{v}) - \sum_{k=1}^{Q_b} \sum_{i=1}^N \theta_k^b(\boldsymbol{\mu}) p_{Ni}(\boldsymbol{\mu}) b_k(\mathbf{v}, \eta_i), \\ (\hat{r}_N^p(\boldsymbol{\mu}), q)_Q &= \sum_{k=1}^{Q_g} \theta_k^g(\boldsymbol{\mu}) g_k(q) - \sum_{k=1}^{Q_b} \sum_{i=1}^{2N} \theta_k^b(\boldsymbol{\mu}) u_{Ni}(\boldsymbol{\mu}) b_k(\boldsymbol{\zeta}_i, q), \end{aligned} \quad (49)$$

and then using the linear superposition principle

$$\begin{aligned} \hat{r}_N^u(\boldsymbol{\mu}) &= \sum_{k=1}^{Q_f} \theta_k^f(\boldsymbol{\mu}) \hat{f}_k - \sum_{k=1}^{Q_a} \sum_{i=1}^{2N} \theta_k^a(\boldsymbol{\mu}) u_{Ni}(\boldsymbol{\mu}) \hat{a}_{ki} - \sum_{k=1}^{Q_b} \sum_{i=1}^N \theta_k^b(\boldsymbol{\mu}) p_{Ni}(\boldsymbol{\mu}) \hat{b}_{ki}^\top, \\ \hat{r}_N^p(\boldsymbol{\mu}) &= \sum_{k=1}^{Q_g} \theta_k^g(\boldsymbol{\mu}) \hat{g}_k - \sum_{k=1}^{Q_b} \sum_{i=1}^{2N} \theta_k^b(\boldsymbol{\mu}) u_{Ni}(\boldsymbol{\mu}) \hat{b}_{ki}, \end{aligned} \quad (50)$$

### 3. OFFLINE/ONLINE STRATEGY FOR RESIDUAL EVALUATION IN THE CRB217

where the Riesz representations  $\hat{a}_{ki}$ ,  $\hat{b}_{ki}$ ,  $\hat{b}_{ki}^\top$ ,  $\hat{f}_k$  and  $\hat{g}_k$  satisfy

$$\begin{aligned}
(\hat{a}_{ki}, \mathbf{v})_V &= a_k(\boldsymbol{\zeta}_i, \mathbf{v}), \quad \forall \mathbf{v} \in V, 1 \leq k \leq Q_a, 1 \leq i \leq 2N \\
(\hat{b}_{ki}^\top, \mathbf{v})_V &= b_k(\mathbf{v}, \eta_i), \quad \forall \mathbf{v} \in V, 1 \leq k \leq Q_b, 1 \leq i \leq N \\
(\hat{b}_{ki}, q)_Q &= b_k(\boldsymbol{\zeta}_i, q), \quad \forall q \in Q, 1 \leq k \leq Q_b, 1 \leq i \leq 2N \\
(\hat{f}_k, \mathbf{v})_V &= f_k(\mathbf{v}), \quad \forall \mathbf{v} \in V, 1 \leq k \leq Q_f \\
(\hat{g}_k, q)_Q &= g_k(q), \quad \forall q \in Q, 1 \leq k \leq Q_g.
\end{aligned} \tag{51}$$

And finally we obtain

$$\begin{aligned}
\|\hat{r}_N^{\mathbf{u}}(\boldsymbol{\mu})\|_V^2 &= C^{ff}(\boldsymbol{\mu}) + \mathbf{u}_N(\boldsymbol{\mu})^\top C_N^{aa}(\boldsymbol{\mu}) \mathbf{u}_N(\boldsymbol{\mu}) + p_N(\boldsymbol{\mu})^\top C_N^{b^\top b^\top}(\boldsymbol{\mu}) p_N(\boldsymbol{\mu}) \\
&\quad + 2\mathbf{u}_N(\boldsymbol{\mu})^\top C_N^{ab^\top}(\boldsymbol{\mu}) p_N(\boldsymbol{\mu}) - 2\mathbf{u}_N(\boldsymbol{\mu})^\top C_N^{af}(\boldsymbol{\mu}) - 2p_N(\boldsymbol{\mu})^\top C_N^{b^\top f}(\boldsymbol{\mu}) \\
\|\hat{r}_N^p(\boldsymbol{\mu})\|_Q^2 &= C^{gg}(\boldsymbol{\mu}) + \mathbf{u}_N(\boldsymbol{\mu})^\top C_N^{bb}(\boldsymbol{\mu}) \mathbf{u}_N(\boldsymbol{\mu}) - 2\mathbf{u}_N(\boldsymbol{\mu})^\top C_N^{bg}(\boldsymbol{\mu})
\end{aligned} \tag{52}$$

with

$$\begin{aligned}
C^{ff}(\boldsymbol{\mu}) &= \sum_{k=1}^{Q_f} \sum_{k'=1}^{Q_f} \theta_k^f(\boldsymbol{\mu}) \theta_{k'}^f(\boldsymbol{\mu}) \underbrace{(\hat{f}_k, \hat{f}_{k'})_V}_{\text{precomputed}}, & C^{gg}(\boldsymbol{\mu}) &= \sum_{k=1}^{Q_g} \sum_{k'=1}^{Q_g} \theta_k^g(\boldsymbol{\mu}) \theta_{k'}^g(\boldsymbol{\mu}) \underbrace{(\hat{g}_k, \hat{g}_{k'})_Q}_{\text{precomputed}}, \\
C_N^{af}(\boldsymbol{\mu}) &= \sum_{k=1}^{Q_f} \sum_{k'=1}^{Q_a} \theta_k^f(\boldsymbol{\mu}) \theta_{k'}^a(\boldsymbol{\mu}) \underbrace{C_{kk',N}^{af}}_{\text{precomputed}}, & (C_{kk',N}^{af})_i &= (\hat{f}_k, \hat{a}_{k'i})_V, \quad 1 \leq i \leq 2N, \\
C_N^{b^\top f}(\boldsymbol{\mu}) &= \sum_{k=1}^{Q_f} \sum_{k'=1}^{Q_b} \theta_k^f(\boldsymbol{\mu}) \theta_{k'}^b(\boldsymbol{\mu}) \underbrace{C_{kk',N}^{b^\top f}}_{\text{precomputed}}, & (C_{kk',N}^{b^\top f})_i &= (\hat{f}_k, \hat{b}_{k'i}^\top)_V, \quad 1 \leq i \leq N, \\
C_N^{bg}(\boldsymbol{\mu}) &= \sum_{k=1}^{Q_g} \sum_{k'=1}^{Q_b} \theta_k^g(\boldsymbol{\mu}) \theta_{k'}^b(\boldsymbol{\mu}) \underbrace{C_{kk',N}^{bg}}_{\text{precomputed}}, & (C_{kk',N}^{bg})_i &= (\hat{g}_k, \hat{b}_{k'i})_Q, \quad 1 \leq i \leq N, \\
C_N^{aa}(\boldsymbol{\mu}) &= \sum_{k=1}^{Q_a} \sum_{k'=1}^{Q_a} \theta_k^a(\boldsymbol{\mu}) \theta_{k'}^a(\boldsymbol{\mu}) \underbrace{C_{kk',N}^{aa}}_{\text{precomputed}}, & (C_{kk',N}^{aa})_{ij} &= (\hat{a}_{ki}, \hat{a}_{k'j})_V, \quad 1 \leq i, j \leq 2N, \\
C_N^{bb}(\boldsymbol{\mu}) &= \sum_{k=1}^{Q_b} \sum_{k'=1}^{Q_b} \theta_k^b(\boldsymbol{\mu}) \theta_{k'}^b(\boldsymbol{\mu}) \underbrace{C_{kk',N}^{bb}}_{\text{precomputed}}, & (C_{kk',N}^{bb})_{ij} &= (\hat{b}_{ki}, \hat{b}_{k'j})_Q, \quad 1 \leq i, j \leq N, \\
C_N^{b^\top b^\top}(\boldsymbol{\mu}) &= \sum_{k=1}^{Q_b} \sum_{k'=1}^{Q_b} \theta_k^b(\boldsymbol{\mu}) \theta_{k'}^b(\boldsymbol{\mu}) \underbrace{C_{kk',N}^{b^\top b^\top}}_{\text{precomputed}}, & (C_{kk',N}^{b^\top b^\top})_{ij} &= (\hat{b}_{ki}^\top, \hat{b}_{k'j}^\top)_Q, \quad 1 \leq i, j \leq N.
\end{aligned} \tag{53}$$

Here again, the precomputed quantities are evaluated only once during the offline phase and stored to be reused online. This way, during the online phase the computation of the residuals remain independent of the FE dimension  $\mathcal{N}$ . The same method is used to evaluate the dual norms of the dual residuals.

## 4 Detailed Affine Decomposition for RB with Stokes Flow

We present here the detailed formulation of the affine decomposition associated to the test case in section 3.

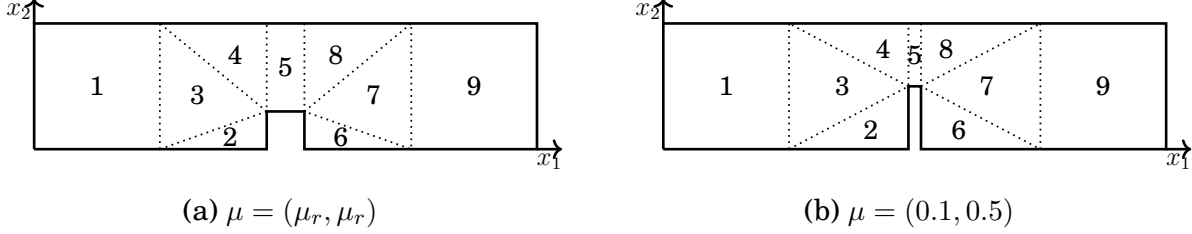


Figure 17 – Domain Decomposition used for the mapping

The domain decomposition is presented on figure 17. The associated mappings on each sub-domain is

$$\mathcal{G}^1(\tilde{\mathbf{x}}, \mu) = \mathcal{G}^7(\tilde{\mathbf{x}}, \mu) = \tilde{\mathbf{x}}$$

$$\mathcal{G}^2(\tilde{\mathbf{x}}, \mu) = \begin{pmatrix} \frac{2-\mu_r}{2-\mu_1} & 0 \\ 0 & \frac{\mu_r}{\mu_2} \end{pmatrix} \tilde{\mathbf{x}} + \begin{pmatrix} \frac{\mu_r - \mu_1}{2-\mu_1} \\ 0 \end{pmatrix}, \quad \mathcal{G}^6(\tilde{\mathbf{x}}, \mu) = \begin{pmatrix} \frac{2-\mu_r}{2-\mu_1} & 0 \\ 0 & \frac{\mu_r}{\mu_2} \end{pmatrix} \tilde{\mathbf{x}} + \begin{pmatrix} \frac{3(\mu_r - \mu_1)}{2-\mu_1} \\ 0 \end{pmatrix}$$

$$\mathcal{G}^3(\tilde{\mathbf{x}}, \mu) = \begin{pmatrix} \frac{2-\mu_r}{2-\mu_1} & 0 \\ \frac{2(\mu_r - \mu_2)}{2-\mu_1} & 1 \end{pmatrix} \tilde{\mathbf{x}} + \begin{pmatrix} \frac{\mu_r - \mu_1}{2-\mu_1} \\ \frac{-2(\mu_r - \mu_2)}{2-\mu_1} \end{pmatrix}, \quad \mathcal{G}^7(\tilde{\mathbf{x}}, \mu) = \begin{pmatrix} \frac{2-\mu_r}{2-\mu_1} & 0 \\ \frac{-2(\mu_r - \mu_2)}{2-\mu_1} & 1 \end{pmatrix} \tilde{\mathbf{x}} + \begin{pmatrix} \frac{3(\mu_r - \mu_1)}{2-\mu_1} \\ \frac{6(\mu_r - \mu_2)}{2-\mu_1} \end{pmatrix}$$

$$\mathcal{G}^4(\tilde{\mathbf{x}}, \mu) = \begin{pmatrix} \frac{2-\mu_r}{2-\mu_1} & 0 \\ 0 & \frac{1-\mu_r}{1-\mu_2} \end{pmatrix} \tilde{\mathbf{x}} + \begin{pmatrix} \frac{\mu_r - \mu_1}{2-\mu_1} \\ \frac{\mu_r - \mu_2}{1-\mu_2} \end{pmatrix}, \quad \mathcal{G}^8(\tilde{\mathbf{x}}, \mu) = \begin{pmatrix} \frac{2-\mu_r}{2-\mu_1} & 0 \\ 0 & \frac{1-\mu_r}{1-\mu_2} \end{pmatrix} \tilde{\mathbf{x}} + \begin{pmatrix} \frac{3(\mu_r - \mu_1)}{2-\mu_1} \\ \frac{\mu_r - \mu_2}{1-\mu_2} \end{pmatrix}$$

$$\mathcal{G}^5(\tilde{\mathbf{x}}, \mu) = \begin{pmatrix} \frac{\mu_r}{0} & 0 \\ 0 & \frac{1-\mu_r}{1-\mu_2} \end{pmatrix} \tilde{\mathbf{x}} + \begin{pmatrix} \frac{-2(\mu_r - \mu_1)}{\mu_1} \\ \frac{\mu_r - \mu_2}{1-\mu_2} \end{pmatrix}$$

It allows to write the affine decomposition of the problem, using the relations,

$$a(\mathbf{u}, \mathbf{v}; \mu) = \sum_{k=1}^P G_{mj}^k(\mu) G_{lj}^k J^k(\mu) \int_{\Omega^k} \frac{\partial u_i}{\partial x_m} \frac{\partial v_i}{\partial x_l} d\Omega, \quad f(\mathbf{v}) = 0 \quad (54)$$

$$b(\mathbf{u}, q) = \sum_{k=1}^P -G_{ji}^k J^k(\mu) \int_{\Omega^k} \frac{\partial u_i}{\partial x_j} q d\Omega, \quad g(q) = 0, \quad (55)$$

we have,

$$a(\mathbf{u}, \mathbf{v}; \mu) = \sum_{s=1}^{Q_a} \theta_a^s(\mu) a^s(\mathbf{u}, \mathbf{v}), \quad (56)$$

$$b(\mathbf{u}, q; \mu) = \sum_{s=1}^{Q_b} \theta_b^s(\mu) b^s(\mathbf{u}, q), \quad (57)$$

with

$$\begin{aligned}\theta_a^1(\mu) &= 1, & \theta_a^2(\mu) &= \mu_1, & \theta_a^3(\mu) &= \frac{1 - \mu_2}{\mu_1}, & \theta_a^4(\mu) &= \frac{1}{2 - \mu_1}, \\ \theta_a^5(\mu) &= \frac{\mu_2}{2 - \mu_1}, & \theta_a^6(\mu) &= \frac{\mu_2^2}{2 - \mu_1}, & \theta_a^7(\mu) &= \frac{2 - \mu_1}{\mu_2}, & \theta_a^8(\mu) &= \frac{\mu_1}{1 - \mu_2}, \\ \theta_a^9(\mu) &= \frac{2 - \mu_1}{1 - \mu_2},\end{aligned}$$

and

$$a^1(\mathbf{u}, \mathbf{v}) = \int_{\Omega^{1,9}} \nabla \mathbf{u} : \nabla \mathbf{v} \, d\Omega + \int_{\Omega^{3,7}} \frac{2}{2 - \mu_r} \frac{\partial \mathbf{u}}{\partial x_2} \cdot \frac{\partial \mathbf{v}}{\partial x_2} \, d\Omega$$

$$a^2(\mathbf{u}, \mathbf{v}) = \int_{\Omega^{3,7}} \frac{-1}{2 - \mu_r} \frac{\partial \mathbf{u}}{\partial x_2} \cdot \frac{\partial \mathbf{v}}{\partial x_2} \, d\Omega,$$

$$a^3(\mathbf{u}, \mathbf{v}) = \int_{\Omega^5} \frac{\mu_r}{1 - \mu_r} \frac{\partial \mathbf{u}}{\partial x_1} \cdot \frac{\partial \mathbf{v}}{\partial x_1} \, d\Omega$$

$$\begin{aligned}a^4(\mathbf{u}, \mathbf{v}) &= \int_{\Omega^{3,7}} (2 - \mu_r) \frac{\partial \mathbf{u}}{\partial x_2} \cdot \frac{\partial \mathbf{v}}{\partial x_2} + \frac{4\mu_r^2}{2 - \mu_r} \frac{\partial \mathbf{u}}{\partial x_2} \cdot \frac{\partial \mathbf{v}}{\partial x_2} \, d\Omega + \int_{\Omega^{4,8}} \frac{2 - \mu_r}{1 - \mu_r} \frac{\partial \mathbf{u}}{\partial x_1} \cdot \frac{\partial \mathbf{v}}{\partial x_1} \, d\Omega \\ &+ \int_{\Omega^3} 2\mu_r \left( \frac{\partial \mathbf{u}}{\partial x_1} \cdot \frac{\partial \mathbf{v}}{\partial x_2} + \frac{\partial \mathbf{u}}{\partial x_2} \cdot \frac{\partial \mathbf{v}}{\partial x_1} \right) \, d\Omega + \int_{\Omega^7} -2\mu_r \left( \frac{\partial \mathbf{u}}{\partial x_1} \cdot \frac{\partial \mathbf{v}}{\partial x_2} + \frac{\partial \mathbf{u}}{\partial x_2} \cdot \frac{\partial \mathbf{v}}{\partial x_1} \right) \, d\Omega \\ &, \end{aligned}$$

$$\begin{aligned}a^5(\mathbf{u}, \mathbf{v}) &= \int_{\Omega^{2,6}} \frac{2 - \mu_r}{\mu_r} \frac{\partial \mathbf{u}}{\partial x_1} \cdot \frac{\partial \mathbf{v}}{\partial x_1} \, d\Omega + \int_{\Omega^{3,7}} \frac{-8\mu_r}{2 - \mu_r} \frac{\partial \mathbf{u}}{\partial x_1} \cdot \frac{\partial \mathbf{v}}{\partial x_1} \, d\Omega \\ &+ \int_{\Omega^3} -2 \left( \frac{\partial \mathbf{u}}{\partial x_1} \cdot \frac{\partial \mathbf{v}}{\partial x_2} + \frac{\partial \mathbf{u}}{\partial x_2} \cdot \frac{\partial \mathbf{v}}{\partial x_1} \right) \, d\Omega + \int_{\Omega^7} 2 \left( \frac{\partial \mathbf{u}}{\partial x_1} \cdot \frac{\partial \mathbf{v}}{\partial x_2} + \frac{\partial \mathbf{u}}{\partial x_2} \cdot \frac{\partial \mathbf{v}}{\partial x_1} \right) \, d\Omega \\ &+ \int_{\Omega^{4,8}} -\frac{2 - \mu_r}{1 - \mu_r} \frac{\partial \mathbf{u}}{\partial x_1} \cdot \frac{\partial \mathbf{v}}{\partial x_1} \, d\Omega,\end{aligned}$$

$$a^6(\mathbf{u}, \mathbf{v}) = \int_{\Omega^{3,7}} \frac{4}{2 - \mu_r} \frac{\partial \mathbf{u}}{\partial x_2} \cdot \frac{\partial \mathbf{v}}{\partial x_2} \, d\Omega,$$

$$a^7(\mathbf{u}, \mathbf{v}) = \int_{\Omega^{2,6}} \frac{\mu_r}{2 - \mu_r} \frac{\partial \mathbf{u}}{\partial x_2} \cdot \frac{\partial \mathbf{v}}{\partial x_2} \, d\Omega,$$

$$a^8(\mathbf{u}, \mathbf{v}) = \int_{\Omega^5} \frac{1 - \mu_r}{\mu_r} \frac{\partial \mathbf{u}}{\partial x_2} \cdot \frac{\partial \mathbf{v}}{\partial x_2} \, d\Omega,$$

$$a^9(\mathbf{u}, \mathbf{v}) = \int_{\Omega^{4,8}} \frac{1 - \mu_r}{2 - \mu_r} \frac{\partial \mathbf{u}}{\partial x_2} \cdot \frac{\partial \mathbf{v}}{\partial x_2} \, d\Omega.$$

$$\theta_b^1(\mu) = 1, \quad \theta_b^2(\mu) = \mu_1, \quad \theta_b^3(\mu) = \mu_2,$$

and

$$\begin{aligned}
b^1(\mathbf{u}, q) &= \int_{\Omega^{1,9}} -q \nabla \cdot \mathbf{u} \, d\Omega + \int_{\Omega^5} \frac{-1}{1 - \mu_r} \frac{\partial u_1}{\partial x_1} \, d\Omega \\
&+ \int_{\Omega^{2,6}} \frac{-2}{2 - \mu_r} q \frac{\partial u_2}{\partial x_2} \, d\Omega + \int_{\Omega^{4,8}} \frac{-1}{1 - \mu_r} q \frac{\partial u_1}{\partial x_1} + \frac{-2}{2 - \mu_r} q \frac{\partial u_2}{\partial x_2} \, d\Omega \\
&+ \int_{\Omega^{3,7}} \frac{-2}{2 - \mu_r} q \frac{\partial u_2}{\partial x_2} - q \frac{\partial u_1}{\partial x_1} \, d\Omega + \int_{\Omega^3} \frac{-2\mu_r}{2 - \mu_r} q \frac{\partial u_1}{\partial x_2} \, d\Omega + \int_{\Omega^7} \frac{2\mu_r}{2 - \mu_r} q \frac{\partial u_1}{\partial x_2} \, d\Omega,
\end{aligned}$$

$$\begin{aligned}
b^2(\mathbf{u}, q) &= \int_{\Omega^{2,6}} \frac{1}{2 - \mu_r} q \frac{\partial u_2}{\partial x_2} \, d\Omega + \int_{\Omega^{3,7}} \frac{1}{2 - \mu_r} q \frac{\partial u_2}{\partial x_2} \, d\Omega + \int_{\Omega^{4,8}} \frac{1}{2 - \mu_r} q \frac{\partial u_2}{\partial x_2} \, d\Omega \\
&+ \int_{\Omega^5} \frac{-1}{\mu_r} q \frac{\partial u_2}{\partial x_2} \, d\Omega,
\end{aligned}$$

$$\begin{aligned}
b^3(\mathbf{u}, q) &= \int_{\Omega^{2,6}} \frac{-1}{\mu_r} \frac{\partial u_1}{\partial x_1} \, d\Omega + \int_{\Omega^3} \frac{2}{2 - \mu_r} q \frac{\partial u_1}{\partial x_2} \, d\Omega + \int_{\Omega^7} \frac{-2}{2 - \mu_r} q \frac{\partial u_1}{\partial x_2} \, d\Omega \\
&+ \int_{\Omega^{4,8}} \frac{1}{1 - \mu_r} q \frac{\partial u_1}{\partial x_1} \, d\Omega + \int_{\Omega^5} \frac{1}{1 - \mu_r} q \frac{\partial u_1}{\partial x_1} \, d\Omega,
\end{aligned}$$

### **Introduction**

#### **Le projet CHORUS**

Cette thèse a été financée dans le cadre du projet CHORUS (Commun Horizon of Open Research in Uncertainty for Simulation), de l'agence nationale de recherche (ANR). Le but de ce projet est de supporter les activités de recherche et de développement informatique dans le domaine de la quantification d'incertitude. Les trois axes de travail principaux étaient : la recherche académique, le développement informatique et le benchmarking de cas-tests industriels. De plus, toutes les applications furent implémentées sous une licence open source afin de profiter à une plus large communauté.

La quantification d'incertitude est la capacité à mesurer les effets des incertitudes sur des modèles numériques et des simulations. Dans l'industrie, ce sujet est particulièrement important durant la phase de design et de certification. Bien que ce sujet soit très bien maîtrisé dans le milieu industriel, l'application à des problèmes industriels présentent encore quelques problèmes de scalabilité. Les méthodes statistiques employées dans ce domaine sont essentiellement basées sur la sur l'exploration des espaces paramétriques et nécessitent de répéter plusieurs fois la même expérience (numérique ou non) avec des jeux de paramètres différents. Il est alors difficilement envisageable d'appliquer ce genre de techniques à des simulations numériques sur des cas-tests industriels qui nécessitent, en général, plusieurs heures de calcul. Pour remédier à ces soucis de scalabilité, nous pouvons envisager deux possibilités : la première est l'exploitation des capacités de calcul HPC (High-Performance Computing) grandissante, la deuxième solution est



le développement de modèles d'ordre réduit. Ces sujets sont fortement liés et furent largement étudié dans le cadre de ce projet ANR.



Le projet réunit différents acteurs de la recherche académique et industrielle. Le but était l'amélioration de la robustesse, de la scalabilité et de la précision de différentes méthodes de réduction d'ordre utilisée dans la cadre de la quantification d'incertitude :

- la réduction "Low Rank"
- la POD (Proper Orthogonal Decomposition)
- les bases réduites (RB)
- le krigeage multi-fidélité
- la calibration Bayésienne

En parallèle de ce travail de recherche, les industriels proposèrent différents applications concrètes. L'idée étant de confronter le travail académiques à la résolution de cas-tests complexes et réalistes.

### Le cas-test Airbus

Le groupe Airbus fut l'une des compagnies leaders du projet. Ils proposèrent différents sujets de recherche ainsi qu'un nombre conséquent de cas-tests. Ils s'intéressaient, en particulier, au développement d'un modèle réduit pour réalisé des simulations d'écoulements aérothermiques. Ce sujet ambitieux a un double objectif. Tout d'abord, il doit d'abord rendre abordable l'utilisation de la simulation numérique lors des phases de design et certification : à l'heure actuelle, la simulation numérique est peu utilisée car elle demande d'énorme ressources de calcul. Dans un second temps, Airbus serait intéressé par l'intégration d'un tel modèle réduit dans son système de de contrôle et de régulation de le l'atmosphère intérieur des avions. Ce système appelé un ECS (Environmental control system) repose sur la résolution de nombreuse équations différentielles ordinaires. L'intégration d'un modèle réduit

3D rapide et précis pourrait grandement améliorer les possibilités de leur ECS, mais il est pour cela nécessaire de réussir à réaliser ces simulations en temps réel.

Le groupe Airbus nous proposa de travailler sur un problème d'écoulements aérothermiques dans un avion. Le domaine de résolution est divisé en deux parties distinctes : en haut, la cabine avec les passagers et en bas, la baie avec différents équipements électriques. Un système de climatisation régule la température dans l'intégralité du domaine, voir figure 18

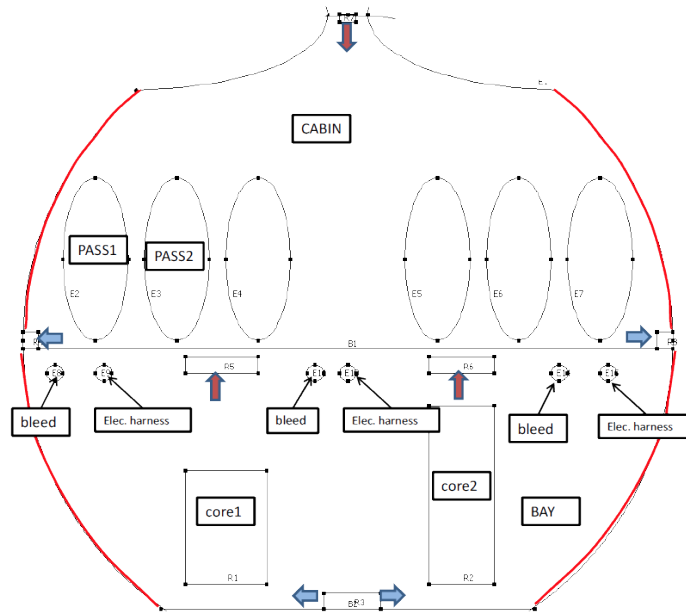


Figure 18 – Description of the Cabin Use Case Proposed by Airbus Group

Le développement d'un modèle réduit pour ce cas-test fut la motivation principale de cette thèse.

### La Librairie Feel++

Feel++ (Finite Element Embedded Library in C++) est une librairie C++ de calcul scientifique, spécialisée dans la méthode des éléments finis. Elle dispose d'un langage très pratique pour la résolution d'équations aux dérivés partielles, en s'appuyant sur un DSEL (Domain Specific Language). Le but de Feel++ est de rendre ce langage aussi proche que possible de la formulation mathématique. L'idée est de fournir un environnement performant dans lequel il est possible de résoudre des problèmes complexes dans un langage mathématique simple et qui exploite efficacement les capacités du calcul parallèle. Feel++ utilise les derniers standards du C++ ainsi que la méta programmation afin de fournir un maximum de généralité. La librairie offre également un parallélisme complètement transparent pour l'utilisateur. Ces capacités de calcul haute performance est indispensable pour

le traitement des problème industriels.

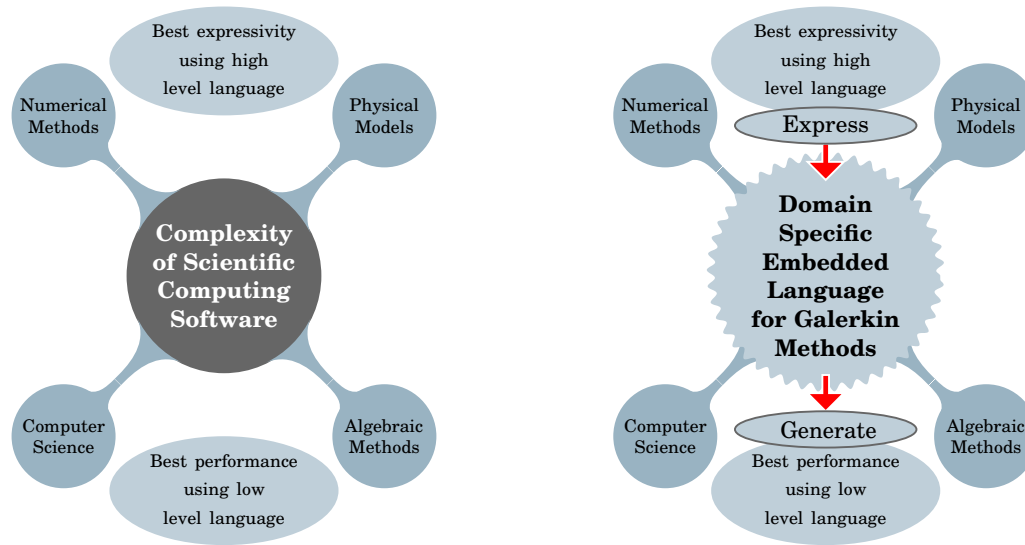


Figure 19 – Le DSEL dans Feel++ fournit un langage de haut niveau afin de cassé la complexité des logiciel de calcul scientifique, tout gardant les performances d’un langage de bas niveau.

Feel++ est une librairie open-source développée sous licence LGPL and GPL. Elle est utilisée dans un grand nombre de projets financés par différentes institutions telles que l’ANR ou la Commission Européenne.

Durant cette thèse, toutes les applications furent développées en utilisant Feel++. Nous avons également produit quelques contributions notables pour la librairie.

## Plan de la thèse

La thèse se compose de 10 chapitres, divisé en 3 parties.

### Chapitre 1 :

- Notions d’analyse numérique, notamment sur les espaces de fonctions,
- Description de la méthode des éléments finis,

*La première partie est consacrée à la modélisation d’un problème aérothermie.*

### Chapitre 2 :

- Introduction des équation du problème
- Discrétisation utilisant la méthode des éléments finis

- Résolution non-linéaire et continuations

### **Chapitre 3 :**

- Méthode de Stabilisation de type Streamline Diffusion (SDM)
- Discussion sur le coefficient de stabilisation et l'évaluation de la longueur caractéristique d'une maille
- Méthode de capture de chocs

### **Chapitre 4 :**

- La méthode Reynolds Average Navier-Stokes (RANS)
- Le modèle de Spalart-Allmaras
- Le modèle  $l - \omega$  SST

*La deuxième partie est consacré la réduction d'ordre*

### **Chapitre 5 :**

- La méthode des bases réduites (RBM) pour les problèmes coercifs
- La méthode des bases réduites pour les problèmes de type point-selle

### **Chapitre 6 :**

- La méthode d'interpolation empirique (EIM)
- EIM pour les opérateurs discrets d'ordre  $r$
- Méthode de coconstruction EIM et RB (SER)

### **Chapitre 7 :**

- Réduction d'un problème d'aérothermie
- Application OPUS : refroidissement de composants électroniques
- Application CHORUS : climatisation dans une cabine d'avion

*La troisième et dernière partie est consacrée à l'implémentation des méthodes décrites précédemment.*

**Chapitre 8 :** Modèle d'aérothermie stabilisé

**Chapitre 9 :** Contribution au framework RB de Feel++

**Chapitre 10 :** Version discrète de l'algorithme DEIM

# 1 Modélisation d'un problème d'aérothermie

## Résolution du système d'équations

Les modèles que nous considérons dans cette thèse se trouvent dans des conditions physiques telles que nous pouvons faire l'hypothèse que le fluide est incompressible. Nous supposons de plus qu'il est Newtonien et que les conditions de température permettent de travailler dans l'approximation de Boussinesq pour le terme de convection naturelle. Avec ces hypothèses, nous pouvons écrire le système d'équations du problème, sous forme forte,

$$\begin{cases} \rho(\mathbf{u} \cdot \nabla \mathbf{u}) + \nabla p - 2\nabla \cdot (\mu \mathbf{S}(\mathbf{u})) = -\rho\beta(T - T_0)\mathbf{g}, \\ \nabla \cdot \mathbf{u} = 0, \\ \mathbf{u} \cdot \nabla T - \nabla \cdot (\kappa \nabla T) = 0. \end{cases}, \quad (58)$$

où  $\mathbf{u}$ ,  $p$  et  $T$  représentent respectivement la vitesse, la pression et la température du fluide considéré, avec les paramètres physiques suivant :

- $\rho$ : la densité du fluide en  $\text{kg.m}^{-3}$
- $\mu$ : la viscosité dynamique du fluide en  $\text{Pa.s}$
- $\kappa$ : la diffusivité thermique du fluide en  $\text{m}^2.\text{s}^{-1} \text{ kg.m}^{-3}$
- $\beta$ : le coefficient d'expansion thermique en  $\text{K}^{-1}$

Ce système d'équation sera clos par des conditions aux bords dépendantes du cas-test considéré.

Ce système d'équation non-linéaire sera résolu en utilisant la méthode des éléments finis (FEM). Pour cela, nous commençons par introduire les espaces de Hilbert  $V$ ,  $Q$  et  $X$  ainsi que les produits scalaires  $(\cdot, \cdot)_V$ ,  $(\cdot, \cdot)_Q$ ,  $(\cdot, \cdot)_X$  et les normes associées  $\|\cdot\|_V$ ,  $\|\cdot\|_Q$ ,  $\|\cdot\|_X$ .

Nous pouvons alors réécrire le système (58) sous forme variationnelle afin de chercher une solution faible  $(\mathbf{u}, p, T) \in V \times Q \times X$  telle que

$$\begin{aligned} \int_{\Omega} \rho(\mathbf{u} \cdot \nabla \mathbf{u}) \cdot \mathbf{v} \, d\Omega + \int_{\Omega} 2\mu \mathbf{S}(\mathbf{u}) : \mathbf{S}(\mathbf{v}) \, d\Omega - \int_{\Gamma} \mathbf{T}(\mathbf{u}, p) \mathbf{n} \cdot \mathbf{v} \, d\Gamma \\ - \int_{\Omega} p \nabla \cdot \mathbf{v} \, d\Omega + \int_{\Omega} (\nabla \cdot \mathbf{u}) q \, d\Omega = - \int_{\Omega} \rho \beta (T - T_0) \mathbf{g} \cdot \mathbf{v} \, d\Omega, \\ \int_{\Omega} (\mathbf{u} \cdot \nabla T) S \, d\Omega + \int_{\Omega} \kappa \nabla T \cdot \nabla S \, d\Omega - \int_{\Gamma} \kappa \nabla T \mathbf{n} \cdot S \, d\Gamma = 0, \end{aligned}$$

$$\forall (\mathbf{v}, q, S) \in V \times Q \times X. \quad (59)$$

Nous introduisons ensuite les espaces discrets

$$V_h \times Q_h = \text{TH}_{ch}^k \subset V \times Q, \quad X_h = \mathbb{L}_{ch}^k \subset X, \quad (60)$$

où  $\text{TH}_{ch}^k$  et  $\mathbb{L}_{ch}^k$  désignent les espaces éléments finis usuels respectivement de Taylor-Hood et de Lagrange. La projection de la forme faible (59) sur ces espaces permet de réécrire le problème sous forme d'une équation matricielle : on cherche  $(\mathbf{u}_h, p_h, T_h) \in V_h \times Q_h \times X_h$  telle que

$$\begin{bmatrix} \mathbf{D}_u + \mathbf{C}(\mathbf{u}_h) & \mathbf{B}^T & 0 \\ \alpha \mathbf{B} & 0 & 0 \\ 0 & 0 & \mathbf{A}(\mathbf{u}_h) + \mathbf{D}_T \end{bmatrix} \begin{bmatrix} \mathbf{u}_h \\ p_h \\ T_h \end{bmatrix} = \begin{bmatrix} \mathbf{F}(T_h) \\ 0 \\ G \end{bmatrix}, \quad (61)$$

Ce système d'équations non-linéaires, peut être écrit de façon générique sous la forme : trouver  $w \in X$  tel que

$$\mathcal{F}(w) = 0. \quad (62)$$

La résolution de ce système non-linéaire est réalisé à l'aide de méthodes itératives telles que l'algorithme de Newton ou de Picard. Dans les deux cas, l'idée est la même. Partant d'un état initial  $w^0$ , nous construisons une suite  $(w^n)$  telle que  $w^n \rightarrow w$ . Chaque itération pour passer de l'état  $n-1$  à  $n$  impliquant la résolution d'un nouveau système linéaire (dépendant du choix de la méthode entre Picard et Newton).

Ce type de méthode étant très sensible au choix de l'état initial, il est crucial de bien choisir ce dernier, sans quoi l'algorithme peut ne pas converger. Afin d'améliorer cette initialisation, nous proposons d'appliquer deux méthodes de continuation. La première est une continuation sur les paramètres physique du fluide  $\mu$  and  $\kappa$ . La continuation consiste en une succession de  $N$  résolutions non-linéaires, du système

$$\mathcal{F}(w; \mu^k, \kappa^k) = 0 \quad (63)$$

ou  $\mu^k = \varepsilon(k)\mu$  et  $\kappa^k = \varepsilon(k)\kappa$ .  $\varepsilon$  étant un coefficient de continuation tel que  $\varepsilon(N) = 1$ . Chaque nouvelle résolution est initialisée avec la solution précédente. Cette méthode

est particulièrement efficace pour résoudre des problèmes stationnaires dans des physiques très dynamiques (viscosité et diffusivité très faible).

La continuation physique ne suffit pourtant pas toujours. Dans ce cas, nous utilisons une continuation par pseudo pas de temps afin de “guider” les états intermédiaires de la résolution non-linéaire. La continuation par pseudo pas de temps consiste à introduire une matrice de masse (comme dans la cas d’une résolution transitoire) sauf que cette matrice est pondérée par un coefficient constant par cellule. Ce coefficient dépend de la taille de la cellule et de l’amplitude du champs de vitesse. Cela permet d’adapter la vitesse de résolution localement. Ce pseudo pas de temps local devient alors de plus en plus grand lorsque le système se rapproche de son état stationnaire.

## Stabilisation

La méthode des éléments finis est connus pour produire des oscillations non-physiques lors de la résolution de problèmes à convection dominante. En règle général, il est possible de supprimer ces instabilités en raffinant le maillage. Cependant, cela peut conduire à des coups de calculs démesurés et cette solution n’est pas envisageable pour des problèmes industriels. Une autre solution est d’utiliser des méthodes de stabilisation qui vont supprimer les oscillations non-physiques en introduisant de la diffusivité numérique. Parmi les différentes méthodes de la littérature, nous avons choisi de nous concentrer sur les méthodes de type Streamline Diffusion (SDM) qui ajoutent de la diffusion numérique dans la direction du champs de convection. Nous avons particulièrement étudié les formulations Streamline Upwind/Petrov-Galerkin (SUPG) et Galerkin Least Square (GLS).

Nous considérons la forme faible (59). Les opérateurs de stabilisation sont ajoutés aux équations de continuité et d’énergie afin d’obtenir la forme stabilisée : trouvé  $(\mathbf{u}_h, p_h, T_h) \in V_h \times Q_h \times X_h$  telle que

$$\begin{aligned}
& \int_{\Omega} \rho(\mathbf{u}_h \cdot \nabla \mathbf{u}_h) \cdot \mathbf{v}_h \, d\Omega + \int_{\Omega} 2\mu \mathbf{S}(\mathbf{u}_h) : \mathbf{S}(\mathbf{v}_h) \, d\Omega + a_K^{NS}(\mathbf{u}_h, \mathbf{v}_h; p_h) \\
& - \int_{\Omega} p_h \nabla \cdot \mathbf{v}_h \, d\Omega + \int_{\Omega} (\nabla \cdot \mathbf{u}_h) q_h \, d\Omega - \int_{\Gamma} \mathbf{T}(\mathbf{u}_h, p_h) \mathbf{n} \cdot \mathbf{v}_h \, d\Gamma \\
& = - \int_{\Omega} \rho \beta (T_h - T_0) \mathbf{g} \cdot \mathbf{v}_h \, d\Omega + l_K^{NS}(\mathbf{u}_h, \mathbf{v}_h), \\
& \int_{\Omega} (\mathbf{u}_h \cdot \nabla T_h) S_h \, d\Omega + \int_{\Omega} \kappa \nabla T_h \cdot \nabla S \, d\Omega - \int_{\Gamma} \kappa \nabla T_h \mathbf{n} \cdot S_h \, d\Gamma + a_K(T_h, S_h) = 0, \\
& \forall (\mathbf{v}_h, q_h, S_h) \in V_h \times Q_h \times X_h.
\end{aligned} \tag{64}$$

avec

$$\begin{aligned}
a_K^{NS}(\mathbf{u}_h, \mathbf{v}_h; p_h) &= \sum_{K \in \mathcal{T}_h} \int_K \tau_h^{NS} (\rho \mathbf{u}_h \nabla \mathbf{u}_h - \nabla \cdot (\mu \nabla \mathbf{u}_h) + \nabla p_h) \cdot (\rho \mathbf{u}_h \nabla \mathbf{v}_h - \gamma \nabla \cdot (\mu \nabla \mathbf{v}_h)) \, d\Omega \\
&\quad + \sum_{K \in \mathcal{T}_h} \int_K \delta_h \nabla \cdot \mathbf{u}_h \nabla \cdot \mathbf{v}_h \, d\Omega \\
l_K^{NS}(\mathbf{u}_h, \mathbf{v}_h) &= \sum_{K \in \mathcal{T}_h} \int_K \tau_h^{NS} g \cdot (\rho \mathbf{u}_h \nabla \mathbf{v}_h - \gamma \nabla \cdot (\mu \nabla \mathbf{v}_h)) \, d\Omega \\
a_K(T_h, S_h) &= \sum_{K \in \mathcal{T}_h} \int_K \tau_h (\mathbf{u}_h \cdot \nabla T_h - \nabla \cdot (\kappa \nabla T_h)) \cdot (\mathbf{u}_h \cdot \nabla S_h - \gamma \nabla \cdot (\kappa \nabla S_h)) \, d\Omega .
\end{aligned} \tag{65}$$

Le paramètre  $\gamma$  permet de choisir entre les différentes méthodes de stabilisation :

- $\gamma = 0$  : SUPG,
- $\gamma = 1$  : GLS,
- $\gamma = -1$  : DW.

Les coefficients de stabilisation  $\tau_h$  et  $\tau_h^{NS}$  sont constant par cellule et dépendent de plusieurs paramètres :

- l'amplitude du champs de convection,
- la taille caractéristique de la cellule,
- l'ordre polynomial de l'espace éléments finis.

Parmi les différentes formulations étudiées, la plus pertinente nous semble être la suivante :

$$\begin{aligned}
\tau_h &= \frac{h}{2|\mathbf{u}_h|_2} \xi(\alpha) \\
\xi(\alpha) &= \begin{cases} \alpha, & 0 \leq \alpha \leq 1, \\ 1, & \alpha \geq 1, \end{cases} \\
\alpha &= \frac{m|\mathbf{u}_h|_p h}{2\kappa}, \\
|\mathbf{u}_h|_p &= \begin{cases} \left( \sum_{i=1}^d |(\mathbf{u}_h)_i|^p \right)^{1/p}, & 1 \leq p < \infty, \\ \max_{i=1, \dots, d} |(\mathbf{u}_h)_i|, & p = \infty, \end{cases} \\
m &= \min\left\{ \frac{1}{3}, 2C \right\}, \\
C &= \frac{1}{\lambda_K h^2},
\end{aligned} \tag{66}$$



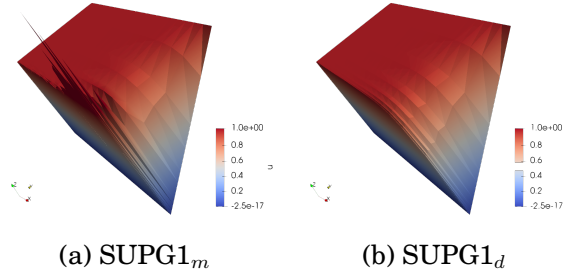


Figure 20 – Intérêt de la formulation adaptée aux maillages anisotropique pour les coefficients de stabilisation. SUPG1<sub>m</sub> évalué avec  $h = h_{\min}$ , SUPG1<sub>d</sub> évalué avec  $h = h_d$

avec  $C$  une constante par élément, vérifiant

$$C \sum_{K \in \mathcal{T}_h} h^2 \|\Delta v\|_K^2 \leq \|\nabla v\|^2, \quad v \in X_h \quad (67)$$

et  $h$  la longueur caractéristique de la cellule. Le calcul pour  $\tau_h^{NS}$  est identique en remplaçant  $\kappa$  par  $\mu$ .

Le choix de la longueur caractéristique a relativement peu d'importance sur le coefficient de stabilisation, tant que le maillage reste isotropique. Dans ce cas, on peut exemple choisir d'évaluer  $h$  comme la longueur de la plus petite arête de la cellule. Cependant, l'évaluation de cette longueur  $h$  devient critique lorsque les mailles sont étirées et donc que le maillage est fortement anisotropique. Dès lors, il est nécessaire de prendre en compte la direction du champs de convection par rapport à l'élongation de la maille. Nous utilisons pour cela l'image du champs de convection sur une maille de référence :

$$h = h_d = \frac{|\mathbf{u}_h|_p}{|\hat{\mathbf{u}}_h|_p} \hat{h} \quad (68)$$

où  $\hat{h}$  est la longueur caractéristique de l'élément de référence. Le champs  $\hat{\mathbf{u}}_h$  est évalué en utilisant la Jacobienne  $\mathbf{J}$  de la transformation géométrique entre l'élément de référence et la cellule courante

$$\hat{\mathbf{u}}_h = \mathbf{J}^{-1} \mathbf{u}_h. \quad (69)$$

Cette définition permet d'évaluer  $h$  dans la direction du champs de convection et donc d'être plus précis dans le calcul du coefficient de stabilisation. L'intérêt de cette formulation est particulièrement visible sur des problèmes avec des maillages fortement anisotropique. La figure 20 met en évidence ce phénomène avec une résolution d'un problème d'advection dans une couche limite maillée avec des simplex (rapport de forme des mailles de  $10^5$ ).

Ces méthodes de stabilisation permettent d'éviter les oscillations non-physique et de résoudre des problèmes à convection dominante, même sur des maillages relativement grossiers. Cependant, elles génèrent tout de même des undershoots et overshoots au niveau des discontinuités importantes. Ces petites oscillations

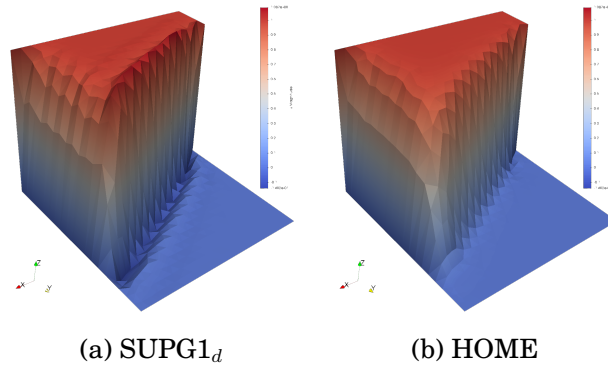


Figure 21 – Intérêt des méthodes de capture de choc sur un cas-test avec une importante discontinuité. SUPG1<sub>d</sub> sans capture de choc, HOME avec capture de choc

peuvent être critique lorsqu'il est par exemple nécessaire de conserver une solution positive. Pour remédier à ce problème, nous introduisons une nouvelle méthode de stabilisation appelé capture de choc. Ces méthodes permettent de réduire ou de faire disparaître ces undershoots. L'idée étant de rajouter de la viscosité numérique soit de façon isotropique, soit orthogonalement au champs de convection. On peut observer de façon qualitative l'effet de ces méthodes sur la figure 21. Dans cet exemple, nous introduisons une diffusion numérique orthogonale au champs de convection, avec l'ajout du terme

$$(\tau_h^{sc} \bar{D} \nabla \mathbf{u}_h, \nabla \mathbf{v}_h) \quad (70)$$

ou le coefficient  $\tau_h^{sc}$  est donné par

$$\tau_h^{sc} = \tau_h(\mathbf{a}) \max(0, \|\mathbf{a}\|_2^2 Q_K(\mathbf{u}_h) - Q_K^2(\mathbf{u}_h)). \quad (71)$$

où  $Q_K(\mathbf{u}_h) = \frac{\|R_h(\mathbf{u}_h)\|_2}{\|\nabla \mathbf{u}_h\|_2}$  et  $R_h(\mathbf{u}_h)$  est le résidu discret.

## Modelisation de la turbulence

Les méthodes de stabilisation introduites précédemment permettent de retrouver une solution stable à des problèmes à convection dominante, cependant elles ne modélisent pas correctement l'effet des plus petites échelles physiques. Une partie de l'énergie est diffusée dans ces petites échelles. Afin de correctement modéliser ce phénomène, nous allons utiliser des modèles de turbulence. Parmi les différents modèles existant dans la littérature, nous avons choisi d'étudier le modèle Reynolds Average Navier-Stokes (RANS) qui est l'un des plus courants et le mieux adapté aux simulations stationnaires.

L'idée du RANS est de faire une moyenne des quantités fluctuantes et sans avoir à les résoudre. La contribution des plus petites échelle est alors introduite par l'ajout d'une viscosité turbulente dans le système d'équation de Navier-Stokes. La diffusion turbulente dans l'équation d'énergie est alors déduite en utilisant le

nombre de Prandtl,

$$\begin{cases} \rho(\mathbf{u} \cdot \nabla \mathbf{u}) + \nabla p - 2\nabla \cdot ((\mu + \mu_t)\mathbf{S}(\mathbf{u})) = -\rho\beta(T - T_0)\mathbf{g}, \\ \nabla \cdot \mathbf{u} = 0, \\ \mathbf{u} \cdot \nabla T - \nabla \cdot ((\kappa + \frac{\mu_t}{Pr})\nabla T) = 0. \end{cases}, \quad (72)$$

L'évaluation de la viscosité turbulente  $\mu_t$  dépend du modèle considéré. Nous avons étudié et implémenté deux modèles en particulier durant cette thèse : le modèle de Spalart Allmaras et le modèle  $k - \omega$  SST.

Le modèle de Spalart Allmaras est probablement le plus utilisé des modèle RANS. La viscosité turbulente est donnée par la relation

$$\mu_t = \rho \bar{\nu} f_{v1}$$

où  $\bar{\nu}$  est solution de l'équation

$$\underbrace{\mathbf{u} \cdot \nabla \bar{\nu}}_{\text{Convection}} - \frac{1}{\sigma} \underbrace{[\nabla \cdot ((\nu + \bar{\nu})\nabla \bar{\nu}) + c_{b2}(\nabla \bar{\nu} \cdot \nabla \bar{\nu})]}_{\text{Diffusion}} - \underbrace{c_{b1}(1 - f_{t2})S\bar{\nu}}_{\text{Production}} + \underbrace{[c_{w1}f_w - \frac{c_{b1}}{\kappa^2}f_{t2}]}_{\text{Destruction}} \left(\frac{\bar{\nu}}{d}\right)^2 = 0, \quad (73)$$

avec les notations suivantes,

$$\begin{aligned} \chi &= \frac{\bar{\nu}}{\nu}, & S &= \sqrt{(\nabla \times \mathbf{u}) \cdot (\nabla \times \mathbf{u})} + \frac{\bar{\nu}}{\kappa^2 d^2} f_{v2}, \\ f_{v1} &= \frac{\chi^3}{\chi^3 + c_{v1}^3}, & f_{v2} &= 1 - \frac{\chi}{1 + \chi f_{v1}}, & f_{t2} &= c_{t3} \exp(-c_{t4}\chi^2), \\ f_w &= g \left( \frac{1 + c_{w3}^6}{g^6 + c_{w3}^6} \right)^{1/6}, & g &= r + c_{w2}(r^6 - r), & r &= \min\left(\frac{\bar{\nu}}{S\kappa^2 d^2}, 10\right), \end{aligned} \quad (74)$$

où  $d$  est la distance à la paroi la plus proche. Les constantes sont données par

$$\begin{aligned} c_{b1} &= 0.1355, & c_{b2} &= 0.622, & c_{w1} &= \frac{c_{b1}}{\kappa^2} + \frac{1 + c_{b2}}{\sigma}, & c_{w2} &= 0.2, & c_{w3} &= 2 \\ c_{v1} &= 7.1, & c_{t3} &= 1.2, & c_{t4} &= 0.5, & \kappa &= 0.41, & \sigma &= 2/3. \end{aligned} \quad (75)$$

Le modèle  $k - \omega$  SST est modèle qui semble plus adapté aux problème d'aérothermie que nous considérons. Il est aussi plus compliqué car implique la résolution d'un système de deux équations non-linéaires. La viscosité turbulente est donnée par

$$\mu_t = \frac{\rho a_1 k}{\max(a_1 \omega, SF_2)}. \quad (76)$$

Ici, l'énergie cinétique turbulente  $k$  et la dissipation spécifique  $\omega$  sont solution du système,

$$\begin{cases} \frac{\partial k}{\partial t} + \mathbf{u} \cdot \nabla k + \beta^* k \omega - \nabla \cdot \left[ \frac{\mu + \sigma_k \mu_t}{\rho} \nabla k \right] = P_k, \\ \frac{\partial \omega}{\partial t} + \mathbf{u} \cdot \nabla \omega + \beta \omega^2 - \nabla \cdot \left[ \frac{\mu + \sigma_w \mu_t}{\rho} \nabla \omega \right] = \frac{\gamma}{\mu_t} P_k + 2(1 - F1) \frac{\sigma_{w2}}{\omega} \nabla k \cdot \nabla \omega, \end{cases} \quad (77)$$

avec les relations suivantes

$$F_1 = \tanh(\arg_1^4), \quad \arg_1 = \min \left[ \max \left( \frac{\sqrt{k}}{\beta^* \omega y}, \frac{500\mu}{\rho y^2 \omega} \right), \frac{4\sigma_{w2}k}{CD_{k\omega} y^2} \right] \quad (78)$$

$$F_2 = \tanh(\arg_2^2), \quad \arg_2 = \max \left( \frac{2\sqrt{k}}{\beta^* \omega y}, \frac{500\mu}{\rho y^2 \omega} \right) \quad (79)$$

$$P_k = \min(\bar{\tau} : \nabla \mathbf{u}, 10\beta^* k\omega), \quad CD_{k\omega} = \max \left( 2\rho \frac{\sigma_{w2}}{\omega} \nabla k \nabla \omega, 10^{-10} \right), \quad (80)$$

$$\mathbf{S} = \sqrt{2(\bar{S} : \bar{S})}, \quad (81)$$

avec

$$\phi = F_1 \phi_1 + (1 - F_1) \phi_2. \quad (82)$$

Les constantes pour ce modèle sont données par

$$\begin{aligned} \sigma_{k1} &= 0.85, & \sigma_{\omega1} &= 0.5, & \beta_1 &= 0.075, & \gamma_1 &= 5/9 \\ \sigma_{k2} &= 1.0, & \sigma_{\omega2} &= 0.856, & \beta_2 &= 0.0828, & \gamma_2 &= 0.44 \\ \beta^* &= 0.09, & \kappa &= 0.41, & a_1 &= 0.31. \end{aligned} \quad (83)$$

## 2 Réduction d'un problème d'aérodynamique

### Méthode des bases réduites

Les solutions d'une équation aux dérivées partielles (EDP) paramétrisée vivent dans une variété. La topologie de cette variété varie en fonction du type d'équation. La méthode des bases réduites cherche à exploiter cette configuration particulière en construisant une approximation de la variété des solutions, quand cela est possible. Cette approximation est un espace, engendré par une base de solution, appelée base réduite.

La construction de cette base se fait lors d'une première phase, offline, qui est particulièrement coûteuse puisqu'elle nécessite de résoudre plusieurs fois le système d'équations. Ces solutions sont alors stockées dans une base de données et elles serviront de nouvelle base de projection par la suite. En pratique, la base réduite est construite de façon itérative, afin de minimiser un critère d'erreur bien défini. La précision et fiabilité de ce critère d'erreur seront déterminant pour la qualité de la base engendrée. Nous disposons aujourd'hui de bornes d'erreur très efficaces pour les problèmes coercifs linéaires, mais il reste un important travail de recherche et/ou d'implémentation à réaliser pour des problèmes plus complexes tel que le notre.

Soit  $\Omega \in \mathbb{R}^d$ ,  $d = 2, 3$  un domaine borné et  $X \subset H^1(\Omega)$  un espace de Hilbert sur ce domaine. Nous notons également  $(\cdot, \cdot)_X$  and  $\|\cdot\|_X = \sqrt{(\cdot, \cdot)_X}$  le produit scalaire et la norme associée à  $X$ . Nous introduisons maintenant un espace de paramètre  $\mathcal{D} \in \mathbb{R}^P$

de dimension  $P \geq 1$  et un problème paramétrisé, écrit de façon générique sous la forme

$$\mathcal{F}(u(\boldsymbol{\mu}); \boldsymbol{\mu}) = 0. \quad (84)$$

L'idée des bases réduites est de construire un espace  $X_N$  de petite dimension, de la forme

$$X_N = \text{span}\{u_h(\boldsymbol{\mu}), \boldsymbol{\mu} \in S_N\}. \quad (85)$$

La sélection des paramètres dans  $S_N$  peut se faire de façon aléatoire ou être optimisée par l'application d'un algorithme glouton qui permet de sélectionner le meilleur paramètre possible à l'aide d'estimateurs d'erreur adaptés. Cette procédure est détaillée dans l'algorithme 18.

---

**Algorithm 18:** Algorithme Glouton pour un problème générique

---

```

1 Choose  $\Xi \subset \mathcal{D}$ ; // Super Sampling
2 Choose  $N_{\max}$ ; // Maximum size of  $X_N$ 
3 Choose  $\delta_{tol} \in (0, 1)$ ; // Algorithm tolerance
4 Choose  $\boldsymbol{\mu}_1 \in \mathcal{D}$ ; // Arbitrary chosen first parameter
5 Set  $N \leftarrow 0, S_N \leftarrow \{\}, X_N \leftarrow \{\}$ ; // Initialization
6 repeat
7    $N \leftarrow N + 1$ ;
8    $S_N \leftarrow S_{N-1} \cup \{\boldsymbol{\mu}_N\}$ ;
9   Solve  $\mathcal{F}_h(u_h(\boldsymbol{\mu}); \boldsymbol{\mu})$ ;
10   $X_N \leftarrow X_{N-1} \oplus \text{span}\{u_h(\boldsymbol{\mu}_N)\}$ ;
11  for  $\boldsymbol{\mu} \in \Xi$  do
12    | Compute error indicator,  $\Delta_N(\boldsymbol{\mu})$ ;
13  end
14   $\boldsymbol{\mu}_{N+1} = \arg \max_{\boldsymbol{\mu} \in \Xi} \Delta_N(\boldsymbol{\mu})$ ;
15 until  $\Delta_N(\boldsymbol{\mu}_{N+1}) < \delta_{tol}$  or  $N \geq N_{\max}$ ;

```

---

Dans le cadre des problèmes d'aérothermie, il est nécessaire de construire trois espaces réduits que nous noterons  $V_N$ ,  $Q_N$  et  $X_N$ . Ces espaces sont construits de la façon suivante

$$\begin{aligned} V_N &= \text{span}\{\mathbf{u}_h(\boldsymbol{\mu}_n), \mathbf{T}^\mu \eta_n\} = \text{span}\{\boldsymbol{\zeta}_n, \boldsymbol{\zeta}_2\}, \\ Q_N &= \text{span}\{p_h(\boldsymbol{\mu}_n)\} = \text{span}\{\eta_n\}, \\ X_N &= \text{span}\{T_h(\boldsymbol{\mu}_n)\} = \text{span}\{\xi_n\}, \end{aligned} \quad 1 \leq n \leq N. \quad (86)$$

L'ajout d'un vecteur de base supplémentaire dans l'espace réduite de vitesse  $V_N$  est nécessaire afin de satisfaire la condition inf-sup lors de la résolution online. La construction de ces espaces est détaillée dans l'algorithme 19.

Lors de la partie online, la forme faible des équations projeté sur ces espaces réduits (au lieu des espaces éléments finis). La résolution est alors réduite à un problème de taille  $\mathcal{O}(N)$  au lieu de  $\mathcal{O}(N)$ . Afin d'être complètement indépendant de la grande dimension (élément finis) lors de cette partie online, les opérateurs (formes

---

**Algorithm 19:** Procédure offline pour la construction des espaces réduits pour un problèmes d'aérotthermie.

---

```

1  Choose  $\Xi \subset \mathcal{D}$  ; // Super Sampling
2  Choose  $N_{\max}$  ; // Maximum size of  $X_N$ 
3  Choose  $\delta_{tol}$  ; // Algorithm tolerance
4  Set  $S_N \leftarrow \{\}$ ,  $V_N \leftarrow \{\}$ ,  $Q_N \leftarrow \{\}$ ,  $X_N \leftarrow \{\}$  ; // Initialization
5  Set  $N_0^u \leftarrow 0$ ,  $N_0^p \leftarrow 0$ ,  $N_0^T \leftarrow 0$ ;
6  repeat
7  |  $S_N \leftarrow S_{N-1} \cup \{\boldsymbol{\mu}_N\}$ ,  $N_N^u \leftarrow N_{N-1}^u$ ,  $N_N^p \leftarrow N_{N-1}^p$ ,  $N_N^T \leftarrow N_{N-1}^T$ ;
8  | Find  $(\mathbf{u}_h(\boldsymbol{\mu}_N), p_h(\boldsymbol{\mu}_N), T_h(\boldsymbol{\mu}_N))$  solution of system (7.1);
9  |  $\{\eta_1, \dots, \eta_{N^p}\} \leftarrow$  Selective Gram-Schmidt $(\{\eta_1, \dots, \eta_{N^p}\}, p_h(\boldsymbol{\mu}_N))$ ;
10 | if  $p_h(\boldsymbol{\mu}_N)$  added then
11 | |  $\{\zeta_1, \dots, \zeta_{N^u}\} \leftarrow$  Selective Gram-Schmidt $(\{\zeta_1, \dots, \zeta_{N^u}\}, \mathbf{T}^\mu \eta_{N^p})$ ;
12 | | end
13 | |  $\{\zeta_1, \dots, \zeta_{N^u}\} \leftarrow$  Selective Gram-Schmidt $(\{\zeta_1, \dots, \zeta_{N^u}\}, \mathbf{u}_h(\boldsymbol{\mu}_N))$ ;
14 | |  $\{\xi_1, \dots, \xi_{N^T}\} \leftarrow$  Selective Gram-Schmidt $(\{\xi_1, \dots, \xi_{N^T}\}, T_h(\boldsymbol{\mu}_N))$ ;
15 | |  $V_N \leftarrow \text{span}\{\zeta_1, \dots, \zeta_{N^u}\}$ ;
16 | |  $Q_N \leftarrow \text{span}\{\eta_1, \dots, \eta_{N^p}\}$ ;
17 | |  $X_N \leftarrow \text{span}\{\xi_1, \dots, \xi_{N^T}\}$ ;
18 | | Precompute online structures;
19 | | for  $\mu \in \Xi$  do
20 | | | Compute error bound,  $\Delta_N(\boldsymbol{\mu})$ ;
21 | | | end
22 | |  $\mu_{N+1} = \arg \max_{\mu \in \Xi} \Delta_N(\boldsymbol{\mu})$ ; // Or can be chosen randomly
23 until  $\Delta_N(\boldsymbol{\mu}_{N+1}) < \delta_{tol}$  or  $N \geq N_{\max}$ ;

```

---

linéaires, bilinéaires et trinéaires) doivent admettre une décomposition affine. La décomposition affine permet de séparer les variables d'espace des variables paramétrisées. Ainsi, nous pourrions pré-calculer toutes les quantités se rapportant à la variable d'espace (de grande dimension) dans la partie offline.

La décomposition affine des opérateurs est à la fois indispensable et difficile à obtenir. Il est par exemple impossible d'écrire les problèmes non-linéaire sous cette forme. Pour palier à ce problème, nous utilisons l'Empirical Interpolation Method (EIM) qui nous permet de construire une approximation affine d'une expression analytique. Nous avons également implémenté une version discrète de cette algorithm, permettant de construire une approximation affine pour des opérateurs discrets, sous forme de matrices ou de vecteurs. EIM est particulièrement utile lorsque nous travaillons avec les opérateurs de stabilisation introduits précédemment.

### La méthode d'interpolation empirique (EIM)

Afin de produire une approximation affine des opérateurs de stabilisation, il fut nécessaire d'implémenter une nouvelle version d'EIM pour des opérateurs discrets. Nous considérons ici un tenseur paramétrisé  $\mathbf{T}(u(\boldsymbol{\mu}), \boldsymbol{\mu})$ , d'ordre  $r = 1, 2, \dots$ . Le but étant de construire une approximation affine de ce tenseur,

$$\mathbf{T}(u(\boldsymbol{\mu}), \boldsymbol{\mu}) \approx \mathbf{T}_M(u(\boldsymbol{\mu}), \boldsymbol{\mu}) = \sum_{m=0}^M \beta_m(\boldsymbol{\mu}) \boldsymbol{\Phi}^m \quad (87)$$

où les  $\beta_m(\boldsymbol{\mu})$  sont des fonctions scalaires,  $\beta_m : \mathcal{D} \rightarrow \mathbb{R}$ , et  $\boldsymbol{\Phi}^m$  sont des tenseurs d'ordre  $r$  indépendants de  $\boldsymbol{\mu}$ . Nous notons  $\bar{W}_M = \text{span}\{\boldsymbol{\Phi}^m\}$ .

Cette espace  $\bar{W}_M$  est construit en initialisant arbitrairement avec un premier paramètre  $\bar{\boldsymbol{\mu}}_1$  qui nous permet de calculer un premier indice d'interpolation

$$i_1 = \arg \max_{i \in \mathbb{I}} |\mathbf{T}(u(\bar{\boldsymbol{\mu}}_1), \bar{\boldsymbol{\mu}}_1)_i|. \quad (88)$$

Ce qui permet d'initialiser l'algorithme glouton avec

$$\boldsymbol{\Phi}_1 = \frac{1}{|\mathbf{T}(u(\bar{\boldsymbol{\mu}}_1), \bar{\boldsymbol{\mu}}_1)_{i_1}|} \mathbf{T}(u(\bar{\boldsymbol{\mu}}_1), \bar{\boldsymbol{\mu}}_1), \quad \bar{W}_1 = \text{span}\{\boldsymbol{\Phi}_1\}, \quad \mathcal{I}_1 = \{i_1\}, \quad \bar{S}_1 = \{\bar{\boldsymbol{\mu}}_1\}. \quad (89)$$

Une fois cette initialisation réalisée, le nouveau paramètre  $\bar{\boldsymbol{\mu}}_M$  est choisi pour maximiser la norme du résidu  $\mathbf{R}_{M-1}(\boldsymbol{\mu})$ ,

$$\mathbf{R}_{M-1}(\boldsymbol{\mu}) = \mathbf{T}(u(\boldsymbol{\mu}), \boldsymbol{\mu}) - \mathbf{T}_{M-1}(u(\boldsymbol{\mu}), \boldsymbol{\mu}), \quad \bar{\boldsymbol{\mu}}_M = \arg \max_{\boldsymbol{\mu} \in \bar{\Xi}} \|\mathbf{R}_M(\boldsymbol{\mu})\|_\infty. \quad (90)$$

Puis nous pouvons choisir le prochain indice d'interpolation  $i_M$ ,

$$i_M = \arg \max_{i \in \mathbb{I}} |(\mathbf{R}_{M-1}(\bar{\boldsymbol{\mu}}_M))_i|, \quad (91)$$

et nous intéressons avec

$$\begin{aligned}\Phi_M &= \frac{1}{|(\mathbf{R}_{M-1}(\bar{\boldsymbol{\mu}}_M))_{i_M}|} \mathbf{R}_{M-1}(\bar{\boldsymbol{\mu}}_M), & \bar{W}_M &= \bar{W}_{M-1} \oplus \text{span}\{\Phi_M\}, \\ \mathcal{I}_M &= \mathcal{I}_{M-1} \cup \{i_M\}, & \bar{S}_M &= \bar{S}_{M-1} \cup \{\bar{\boldsymbol{\mu}}_M\}.\end{aligned}\quad (92)$$

Une fois l'espace  $W_N$  construit, à chaque nouvelle évaluation online, nous calculons les  $\beta_m(u, \boldsymbol{\mu})$  en résolvant le système  $M \times M$  suivant

$$\mathbf{A}_M \boldsymbol{\beta}_M(u, \boldsymbol{\mu}) = \mathbf{F}_M(u, \boldsymbol{\mu}), \quad (93)$$

avec  $(\mathbf{A}_M)_{lj} = (\Phi_j)_{i_l}$ ,  $(\boldsymbol{\beta}_M(u, \boldsymbol{\mu}))_l = \beta_l(u, \boldsymbol{\mu})$  et  $(\mathbf{F}_M(u, \boldsymbol{\mu}))_l = (\mathbf{T}(u, \boldsymbol{\mu}))_{i_l}$ .

### Coconstruction EIM et RB

Durant la construction de l'espace EIM, l'algorithme glouton est obligé de réaliser un grand nombre de résolutions éléments finis pour choisir le nouveau paramètre donné par

$$\bar{\boldsymbol{\mu}}_M = \arg \max_{\boldsymbol{\mu} \in \bar{\Xi}} \|\mathbf{R}_{M-1}(\boldsymbol{\mu})\|_\infty = \arg \max_{\boldsymbol{\mu} \in \bar{\Xi}} \|\mathbf{T}(u(\boldsymbol{\mu}), \boldsymbol{\mu}) - \mathbf{T}_{M-1}(u(\boldsymbol{\mu}), \boldsymbol{\mu})\|_\infty. \quad (94)$$

Afin de palier à ce problème, une solution est d'utiliser une approximation RB à la place de  $u_h(\boldsymbol{\mu})$ . Cela implique de construire l'approximation base réduite simultanément et d'alterner entre RB et EIM. La procédure complète est résumée dans l'algorithme ???. Cette procédure sera utiliser pour la réduction des opérateurs de stabilisation.



**Algorithm 20: Algorithm SER**


---

```

1 Choose  $\Xi \subset \mathcal{D}, \bar{\Xi} \subset \mathcal{D}$ ; // Super Samplings
2 Choose  $N_{\max} \in \mathbb{N}, M_{\max} \in \mathbb{N}$ ; // Max size for  $X_N$  and  $\bar{W}_M$ 
3 Choose  $\delta_{EIM}, \delta_{RB}$ ; // Algorithm tolerances
4 Choose  $\mu_1 \in \Xi$ ; // Initialization of RB construction
5 Set  $N \leftarrow 0, S_0 \leftarrow \{\}, X_0 \leftarrow \emptyset$ ;
6 Choose  $\bar{\mu}_1 \in \bar{\Xi}$ ; // Initialization of EIM construction
7 Set  $i_M \leftarrow \arg \max_{i \in \mathbb{I}} |\mathbf{T}(u_h(\bar{\mu}_1), \bar{\mu}_1)_i|, \Phi_M \leftarrow \frac{\mathbf{T}(u_h(\bar{\mu}_1), \bar{\mu}_1)}{|(\mathbf{T}(u_h(\bar{\mu}_1), \bar{\mu}_1))_{i_1}|}$ ;
8 Set  $M \leftarrow 1, \bar{S}_1 \leftarrow \{\bar{\mu}_1\}, W_1 \leftarrow \{\Phi_1\}, \mathcal{I}_1 \leftarrow \{i_1\}$ ;
9 Set  $\text{state}_{EIM} \leftarrow \text{todo}, \text{state}_{RB} \leftarrow \text{todo}$ ;
10 repeat
11   if  $\text{state}_{RB} = \text{todo}$  then // RB Iteration
12      $N \leftarrow N + 1$ ;
13      $S_N \leftarrow S_{N-1} \cup \{\mu_N\}$ ;
14      $u_h(\mu_N) \leftarrow \text{FE solution using current EIM affine approximation}$ ;
15      $X_N \leftarrow X_{N-1} \oplus \text{span}\{u_h(\mu_N)\}$ ;
16      $\mu_{N+1} = \arg \max_{\mu \in \Xi} \Delta_N(\mu)$ ;
17     if  $\Delta_N(\mu_{N+1}) \leq \delta_{RB}$  or  $N \geq N_{\max}$  then
18       |  $\text{state}_{RB} \leftarrow \text{finished}$ ;
19     end
20   end
21   if  $\text{state}_{EIM} = \text{todo}$  then
22      $\bar{\mu}_{M+1} \leftarrow \arg \max_{\mu \in \bar{\Xi}} \|\underbrace{\mathbf{T}(u_N(\mu), \mu) - \mathbf{T}_M(u_N(\mu), \mu)}_{\mathbf{R}_{M,N}(\mu)}\|_{\infty}$ ;
23     if  $\|\mathbf{R}_{M,N}(\bar{\mu}_{M+1})\|_{\infty} \leq \delta_{EIM}$  or  $M \geq M_{\max}$  then
24       |  $\text{state}_{EIM} \leftarrow \text{finished}$ ;
25     else
26       |  $M \leftarrow M + 1$ ;
27       |  $i_M \leftarrow \arg \max_{i \in \mathbb{I}} |(\mathbf{R}_{M-1}(\bar{\mu}_M))_i|$ ;
28       |  $\Phi_M \leftarrow \frac{\mathbf{R}_{M-1}(\bar{\mu}_M)}{|(\mathbf{R}_{M-1}(\bar{\mu}_M))_{i_M}|}$ ;
29       |  $\bar{S}_M \leftarrow \bar{S}_{M-1} \cup \{\bar{\mu}_M\}, \bar{W}_M \leftarrow \bar{W}_{M-1} \oplus \text{span}\{\Phi_M\}, \mathcal{I}_M = \mathcal{I}_{M-1} \cup \{i_M\}$ ;
30     end
31   end
32 until  $\text{state}_{EIM} = \text{finished}$  and  $\text{state}_{RB} = \text{finished}$ ;

```

---

### 3 Applications numériques

Nous avons appliqué la réduction d'ordre sur deux applications proposé par Airbus Group.

#### Cas test OPUS : refroidissement de composants électroniques

Nous considérons une représentation 2D d'un système de refroidissement au voisinage d'un circuit imprimé. Deux composants (IC) sont sources de chaleurs. Le système de refroidissement souffle de l'air froid depuis le bas du domaine et ressort par le haut. Cette configuration est synthétisée sur la figure 22.

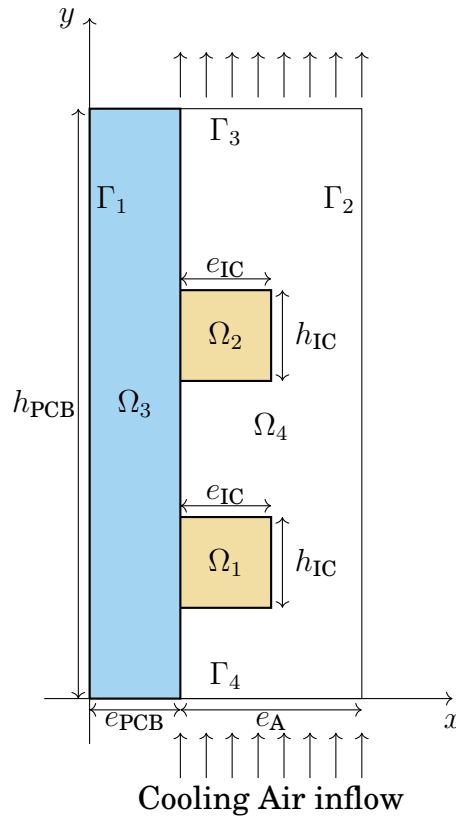


Figure 22 – Geometry of the problem

Le système d'équations décrivant ce problème est

$$\begin{aligned}
 \rho \mathbf{u} \cdot \nabla \mathbf{u} - 2\nabla \cdot (\mu \mathbf{S}(\mathbf{u})) + \nabla p &= 0, & \text{on } \Omega_4 \\
 \nabla \cdot \mathbf{u} &= 0, & \text{on } \Omega_4 \\
 \rho_i C_i \mathbf{v} \cdot \nabla T - \nabla \cdot (k_i \nabla T) &= Q_i, & \text{on } \Omega_i, \quad i = 1, 2, 3, 4.
 \end{aligned} \tag{95}$$

Nous considérons les conditions aux bords sur la température :

- sur  $\Gamma_1 \cup \Gamma_2 \cup \Gamma_3$ ,
 
$$-k\nabla T \cdot \mathbf{n} = 0, \quad (96)$$

- sur  $\Gamma_4$ ,
 
$$T = T_0; \quad (97)$$

- sur tous les bords intérieurs
 
$$T_i = T_j, \quad k_i \nabla T \cdot \mathbf{n} = -k_j \nabla T \cdot \mathbf{n}. \quad (98)$$

et les conditions sur le fluide

- sur  $\Gamma_3$ 

$$\mu \nabla \mathbf{u} \cdot \mathbf{n} - p\mathbf{n} = 0 \quad (99)$$

- sur  $\Gamma_2 \cup \Gamma_I$ 

$$\mathbf{u} = 0 \quad (100)$$

- sur  $\Gamma_4$ 

$$\mathbf{u} = \mathbf{u}_P = \left( 0, \frac{3}{2(e_A - e_{IC})} D \left( 1 - \left( \frac{x - \left( \frac{e_A + e_{IC}}{2} + e_{PCB} \right)}{\frac{e_A - e_{IC}}{2}} \right)^2 \right) \right) \quad (101)$$

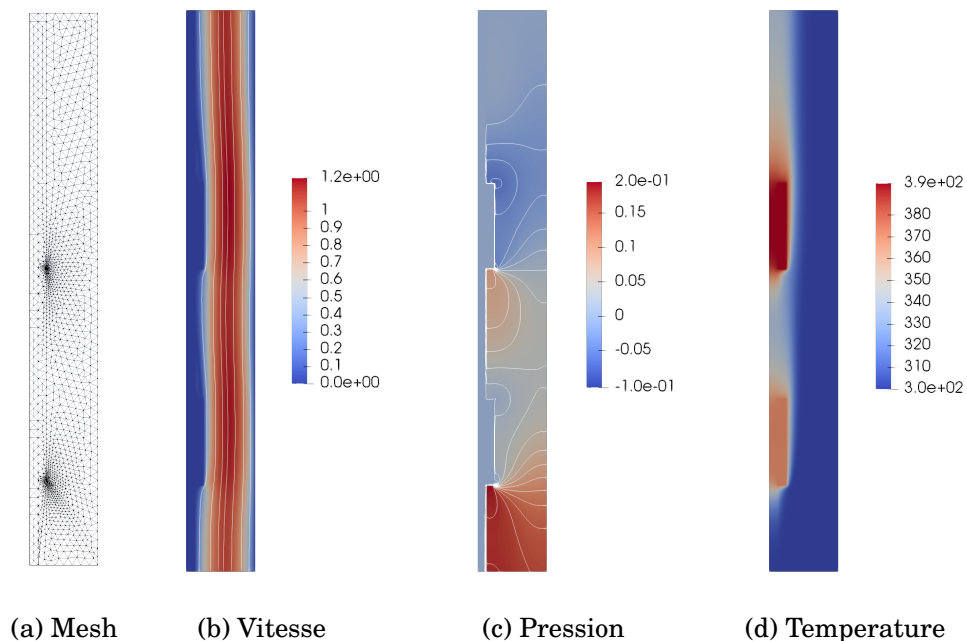
Les paramètres du modèle sont rassemblés dans le tableau 1.

Un exemple de solution est présenté en figure 23 ainsi que le maillage utilisé pour la simulation.

Cette simulation a été réalisée dans deux configurations différentes. La première n'a pas de stabilisation et est réalisée sur un maillage plus fin. La deuxième a nécessité l'ajout de stabilisation, sur un maillage plus grossier et avec des vitesses d'entrée plus élevées. La convergence de la méthode pour ces deux cas test illustré respectivement sur les figures 24 et 25 Ces résultats sont très satisfaisant et les seuls réserves apportés sur notre modèle, pour l'instant, porte sur des optimisations techniques qui doivent être apporté à l'actuelle implémentation.

Name	Description	Nominal Value	Range	Units
IC				
$\rho_1 C_1 = \rho_2 C_2$	Heat Capacity	$1.4 \cdot 10^6$		$\text{J.m}^{-1}.\text{K}^{-1}$
$Q_1, Q_2$	Heat source	$10^6$	$[0, 10^6]$	$\text{W.m}^{-3}$
$k_1, k_2$	Thermal conductivity	2	$[0.2, 150]$	$\text{W.m}^{-1}.\text{K}^{-1}$
$e_{\text{IC}}$	Thickness	$2 \cdot 10^{-3}$		m
$h_{\text{IC}}$	Height	$70 \cdot 10^{-3}$		m
PCB				
$\rho_3 C_3$	Heat Capacity	$2 \cdot 10^6$		$\text{J.m}^{-1}.\text{K}^{-1}$
$Q_3$	Heat source	0		$\text{W.m}^{-3}$
$k_3$	Thermal conductivity	0.2		$\text{W.m}^{-1}.\text{K}^{-1}$
$e_{\text{PCB}}$	Thickness	$2 \cdot 10^{-3}$		m
$h_{\text{PCB}}$	Height	$130 \cdot 10^{-3}$		m
Air				
$\rho_4 C_4$	Heat Capacity	1100		$\text{J.m}^{-1}.\text{K}^{-1}$
$Q_4$	Heat source	0		$\text{W.m}^{-3}$
$k_4$	Thermal conductivity	$3 \cdot 10^{-2}$		$\text{W.m}^{-1}.\text{K}^{-1}$
$\rho$	Density	1.0		$\text{kg.m}^{-3}$
$\mu$	Viscosity	$1.8 \cdot 10^{-5}$		$\text{kg.m}^{-1}.\text{s}^{-1}$
$e_{\text{A}}$	Thickness	$14 \cdot 10^{-3}$		m
$T_0$	Inflow temperature	300		K
$D$	Inflow rate	7	$[0.5, 50]$	$10^{-3} \text{ m}^2.\text{s}^{-1}$

Table 1 – Application OPUS: paramètres du modèle

Figure 23 – Application OPUS: mesh et exemple de solution pour  $\mu = (150, 150, 10^6, 10^6, 10^{-2})$ .

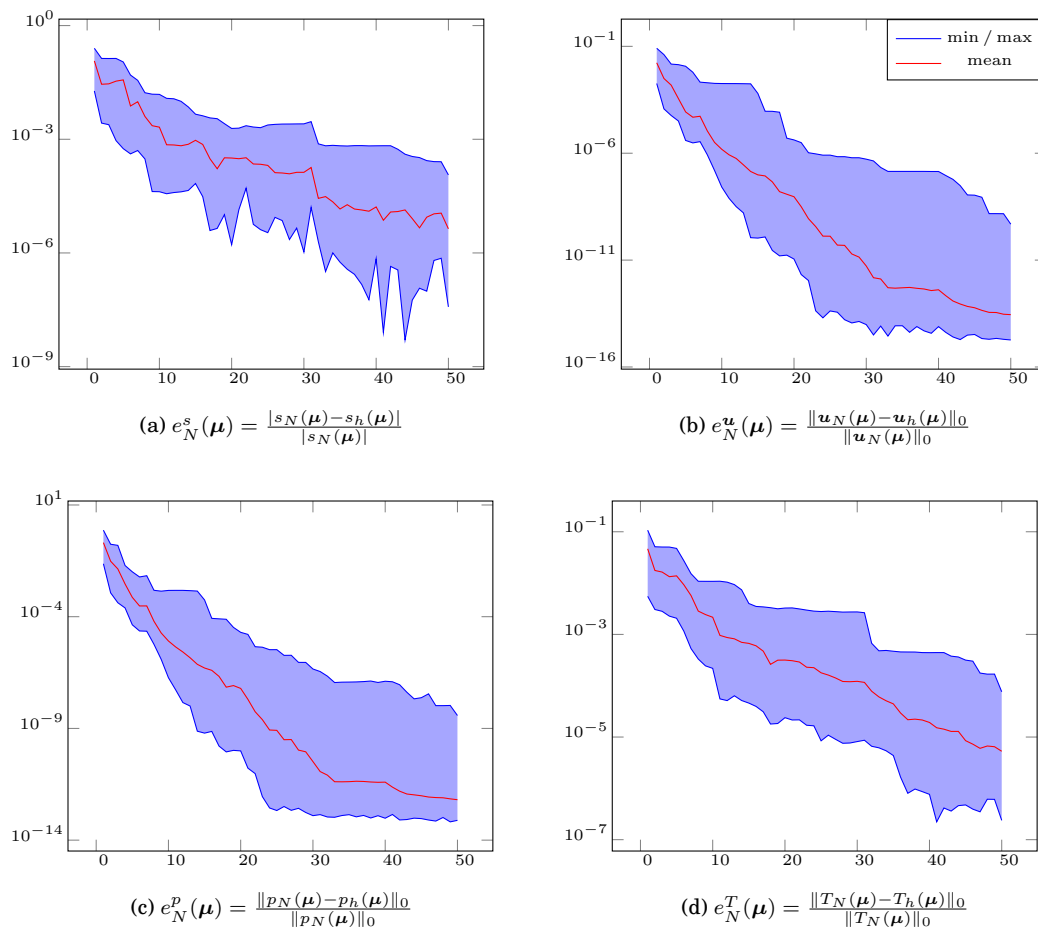


Figure 24 – Application OPUS sans stabilisation : convergence de l’erreur en fonction de la taille de l’espace réduit utilisé. Maximum, minimum et moyenne de l’erreur, évalué sur 50 approximation RB

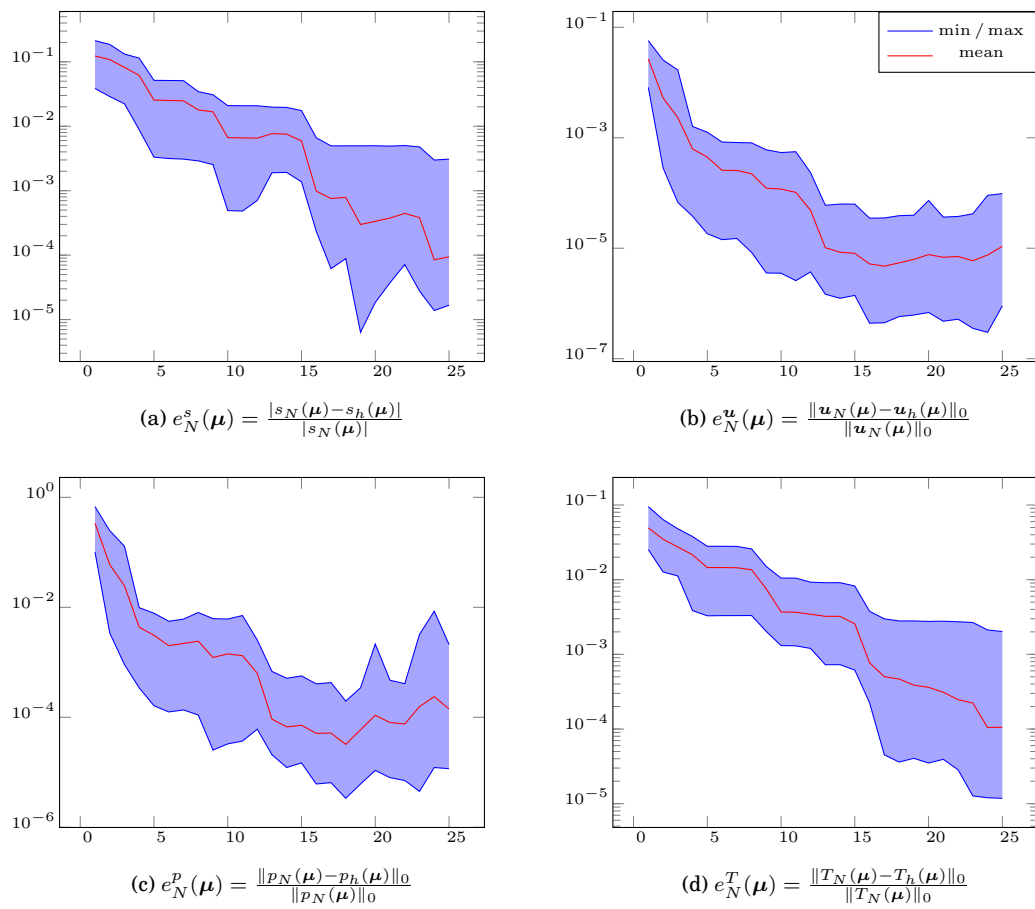


Figure 25 – Application OPUS sans stabilisation : convergence de l’erreur en fonction de la taille de l’espace réduit utilisé. Maximum, minimum et moyenne de l’erreur, évalué sur 50 approximation RB



---

## Bibliography

---

- [1] Regina C Almeida and Renato S Silva. A stable petrov-galerkin method for convection-dominated problems. *Computer methods in applied mechanics and engineering*, 140(3-4):291–304, 1997.
- [2] Etienne Balmès. Parametric families of reduced finite element models. theory and applications. *Mechanical Systems and Signal Processing*, 10(4):381–394, 1996.
- [3] Maxime Barrault, Yvon Maday, Ngoc Cuong Nguyen, and Anthony T Patera. An empirical interpolation method: application to efficient reduced-basis discretization of partial differential equations. *Comptes Rendus Mathématique*, 339(9):667–672, 2004.
- [4] George Keith Batchelor. *An introduction to fluid dynamics*. Cambridge university press, 2000.
- [5] Franco Brezzi. On the existence, uniqueness and approximation of saddle-point problems arising from lagrangian multipliers. *Revue française d'automatique, informatique, recherche opérationnelle. Analyse numérique*, 8(R2):129–151, 1974.
- [6] Franco Brezzi and Michel Fortin. *Mixed and hybrid finite element methods*, volume 15. Springer Science & Business Media, 2012.
- [7] A. Brooks and T. J R Hughes. Streamline-upwind/petrov-galerkin methods for advection dominated flows. *IN: PROC. 3RD INT. CONF. ON FINITE ELEMENTS IN FLOW PROBLEMS*, , D.H. NORRIE (ED.), 2 , Calgary, Canada, Calgary Univ., no date, 1980.



- [8] Alexander N. Brooks and Thomas J.R. Hughes. Streamline upwind/petrov-galerkin formulations for convection dominated flows with particular emphasis on the incompressible navier-stokes equations. *Computer Methods in Applied Mechanics and Engineering*, 32(1):199 – 259, 1982.
- [9] A. Buffa, Y. Maday, A. T. Patera, C. Prud'homme, and G. Turinici. A priori convergence of the greedy algorithm for the parametrized reduced basis. *M2AN Math. Model. Numer. Anal.*, 46:595–605, 2012.
- [10] P. Burda, J. Novotny, and J. Sistek. On a modification of gls stabilized fem for solving incompressible viscous flows. *International Journal for Numerical Methods in Fluids*, 51:1001–1016, 2005.
- [11] E. Burman, M. A. Fernandez, and P. Hansbo. Continuous interior penalty finite element method for oseen's equations. *SIAM Journal of Numerical Analysis*, 44:1248–1274, 2006.
- [12] Erik Burman and Alexandre Ern. Nonlinear diffusion and discrete maximum principle for stabilized galerkin approximations of the convection–diffusion–reaction equation. *Computer Methods in Applied Mechanics and Engineering*, 191(35):3833–3855, 2002.
- [13] Erik Burman and Alexandre Ern. Stabilized galerkin approximation of convection-diffusion-reaction equations: discrete maximum principle and convergence. *Mathematics of computation*, 74(252):1637–1652, 2005.
- [14] N Cardoso and P Bicudo. Time dependent simulation of the driven lid cavity at high reynolds number. *arXiv preprint arXiv:0809.3098*, 2008.
- [15] Vincent Chabannes. *Vers la simulation numérique des écoulements sanguins*. Theses, Université de Grenoble, July 2013.
- [16] Saifon Chaturantabut and Danny C Sorensen. Nonlinear model reduction via discrete empirical interpolation. *SIAM Journal on Scientific Computing*, 32(5):2737–2764, 2010.
- [17] Philippe G Ciarlet. *The finite element method for elliptic problems*. SIAM, 2002.
- [18] Ramon Codina. A discontinuity-capturing crosswind-dissipation for the finite element solution of the convection-diffusion equation. *Computer Methods in Applied Mechanics and Engineering*, 110(3-4):325–342, 1993.
- [19] Todd S Coffey, CT Kelley, and David E Keyes. Pseudotransient continuation and differential-algebraic equations. *SIAM Journal on Scientific Computing*, 25(2):553–569, 2003.
- [20] Thierry Coupez and Elie Hachem. Solution of high-reynolds incompressible flow with stabilized finite element and adaptive anisotropic meshing. *Computer Methods in Applied Mechanics and Engineering*, 267:65–85, 2013.

- [21] Nguyen Ngoc Cuong, Karen Veroy, and Anthony T Patera. Certified real-time solution of parametrized partial differential equations. In *Handbook of materials modeling*, pages 1529–1564. Springer, 2005.
- [22] Cécile Daversin and Christophe Prud’Homme. Simultaneous empirical interpolation and reduced basis method for non-linear problems. *Comptes Rendus Mathématique*, 353(12):1105–1109, 2015.
- [23] Cécile Daversin, Stéphane Veys, Christophe Trophime, and Christophe Prud’Homme. A reduced basis framework: Application to large scale non-linear multi-physics problems. In *ESAIM: Proceedings*, volume 43, pages 225–254. EDP Sciences, 2013.
- [24] Cecile Daversin Catty. *Reduced basis method applied to large non-linear multi-physics problems : application to high field magnets design*. Theses, Université de Strasbourg, September 2016.
- [25] D De Valhl Davis. Natural convection of air in a square cavity: a benchmark solution. *International Journal for Numerical Methods in Fluids*, 3:249–264, 1983.
- [26] S. Deparis and G. Rozza. Reduced basis method for multi-parameter-dependent steady navier–stokes equations: applications to natural convection in a cavity. *Journal of Computational Physics*, 228(12):4359–4378, 2009.
- [27] Ibrahima Dione. *Analyse théorique et numérique des conditions de glissement pour les fluides et les solides par la méthode de pénalisation*. PhD thesis, Université Laval, 2013.
- [28] Eduardo Gomes Dutra Do Carmo and Augusto Cesar Galeão. Feedback petrov-galerkin methods for convection-dominated problems. *Computer methods in applied mechanics and engineering*, 88(1):1–16, 1991.
- [29] Jim Douglas and Jun Ping Wang. An absolutely stabilized finite element method for the stokes problem. *Mathematics of computation*, 52(186):495–508, 1989.
- [30] Vincent Doyeux. *Modélisation et simulation de systèmes multi-fluides. Applications aux écoulements sanguins*. PhD thesis, 2014. Thèse de doctorat dirigée par Peyla, Philippe et Ismail, Mourad Physique Grenoble 2014.
- [31] Ercan Erturk and C Gökçöl. Fourth-order compact formulation of navier–stokes equations and driven cavity flow at high reynolds numbers. *International Journal for Numerical Methods in Fluids*, 50(4):421–436, 2006.
- [32] JP Fink and WC Rheinboldt. On the error behavior of the reduced basis technique for nonlinear finite element approximations. *ZAMM-Journal of Applied Mathematics and Mechanics / Zeitschrift für Angewandte Mathematik und Mechanik*, 63(1):21–28, 1983.

- [33] Leopoldo P Franca and F Valentin. On an improved unusual stabilized finite element method for the advective–reactive–diffusive equation. *Computer Methods in Applied Mechanics and Engineering*, 190(13):1785–1800, 2000.
- [34] A.-L. Gerner and K. Veroy. Reduced basis a posteriori error bounds for the stokes equations in parametrized domains: A penalty approach. *Math. Models Methods Appl. Sci.*, 21:2103–2134, 2011.
- [35] A.-L. Gerner and K. Veroy. Certified reduced basis methods for parametrized saddle point problems. *SIAM, Journal of Scientific Computing*, 2012.
- [36] Swetlana Giere, Traian Iliescu, Volker John, and David Wells. Supg reduced order models for convection-dominated convection–diffusion–reaction equations. *Computer Methods in Applied Mechanics and Engineering*, 289:454–474, 2015.
- [37] Martin A Grepl, Yvon Maday, Ngoc C Nguyen, and Anthony T Patera. Efficient reduced-basis treatment of nonaffine and nonlinear partial differential equations. *ESAIM: Mathematical Modelling and Numerical Analysis*, 41(3):575–605, 2007.
- [38] PM Gresho, DF Griffiths, and DJ Silvester. Adaptive time-stepping for incompressible flow part i: scalar advection-diffusion. *MIMS EPrint*, 2008.
- [39] H. Grotjans and F. Menter. Wall functions for general application CFD codes. In Papailou, editor, *ECCOMAS 98*, pages 1112–1117, 1998.
- [40] Elie Hachem, Benjamin Rivaux, Thibaud Kloczko, Hugues Dignonnet, and Thierry Coupez. Stabilized finite element method for incompressible flows with high reynolds number. *Journal of Computational Physics*, 229(23):8643–8665, 2010.
- [41] Yannick Hoarau. *Analyse physique par simulation numérique et modélisation des écoulements décollés instationnaires autour de surfaces portantes*. PhD thesis, 2002. Thèse de doctorat dirigée par Braza, Marianna Dynamique des fluides Toulouse, INPT 2002.
- [42] Thomas J.R. Hughes, Leopoldo P. Franca, and Marc Balestra. A new finite element formulation for computational fluid dynamics: V. circumventing the babuska-brezzi condition: a stable petrov-galerkin formulation of the stokes problem accommodating equal-order interpolations. *Computer Methods in Applied Mechanics and Engineering*, 59(1):85 – 99, 1986.
- [43] Thomas J.R. Hughes, Leopoldo P. Franca, and Gregory M. Hulbert. A new finite element formulation for computational fluid dynamics: Viii. the galerkin/least-squares method for advective-diffusive equations. *Computer Methods in Applied Mechanics and Engineering*, 73(2):173 – 189, 1989.

- [44] Thomas J.R. Hughes, Leopoldo P. Franca, and Michel Mallet. A new finite element formulation for computational fluid dynamics: Vi. convergence analysis of the generalized supg formulation for linear time-dependent multidimensional advective-diffusive systems. *Computer Methods in Applied Mechanics and Engineering*, 63(1):97 – 112, 1987.
- [45] Thomas J.R. Hughes, Michel Mallet, and Mizukami Akira. A new finite element formulation for computational fluid dynamics: Ii. beyond supg. *Computer Methods in Applied Mechanics and Engineering*, 54(3):341 – 355, 1986.
- [46] T.J.R. Hughes and A.N. Brooks. A multidimensional upwind scheme with no crosswind diffusion. *Finite Element Methods for Convection Dominated Flows*, 34:19–35, 1979.
- [47] Dinh Bao Phuong Huynh, Gianluigi Rozza, Sugata Sen, and Anthony T Patera. A successive constraint linear optimization method for lower bounds of parametric coercivity and inf–sup stability constants. *Comptes Rendus Mathematique*, 345(8):473–478, 2007.
- [48] T. J.R. Hughes I. Harari. What are c and h?: Inequalities for the analysis and design of finite element methods. *Computer Methods in Applied Mechanics and Engineering*, 97:157–192, 1992.
- [49] Kazufumi Ito and SS Ravindran. A reduced-order method for simulation and control of fluid flows. *Journal of computational physics*, 143(2):403–425, 1998.
- [50] Teymour Javaherchi. Review of spalart-allmaras turbulence model and its modifications. *University of Washington*, 2010.
- [51] Volker John and Petr Knobloch. On discontinuity—capturing methods for convection—diffusion equations. *Numerical mathematics and advanced applications*, pages 336–344, 2006.
- [52] C Johnson, AH Schatz, and LB Wahlbin. Crosswind smear and pointwise errors in streamline diffusion finite element methods. *mathematics of computation*, 49(179):25–38, 1987.
- [53] Claes Johnson, Uno Nävert, and Juhani Pitkäranta. Finite element methods for linear hyperbolic problems. *Computer methods in applied mechanics and engineering*, 45(1-3):285–312, 1984.
- [54] DA Kay, PM Gresho, DF Griffiths, and DJ Silvester. Adaptive time-stepping for incompressible flow part ii: Navier-stokes equations. *SIAM Journal on Scientific Computing*, 2010.
- [55] Carl Timothy Kelley and David E Keyes. Convergence analysis of pseudo-transient continuation. *SIAM Journal on Numerical Analysis*, 35(2):508–523, 1998.

- [56] D. J. Knezevic, N-C. Nguyen, and A. T. Patera. Reduced basis approximation and a posteriori error estimation for the parametrized unsteady boussinesq equations. *Mathematical Models and Methods in Applied Sciences*, 2010.
- [57] Andrey Nikolaevich Kolmogorov. The local structure of turbulence in incompressible viscous fluid for very large reynolds numbers. In *Dokl. Akad. Nauk SSSR*, volume 30, pages 299–303, 1941.
- [58] L. Kovasznay. Laminar flow behind a two-dimensional grid. *Mathematical Proceedings of the Cambridge Philosophical Society*, 44:58–62, 1948.
- [59] Dmitri Kuzmin, Otto Mierka, and Stefan Turek. On the implementation of the  $\kappa$ - $\varepsilon$  turbulence model in incompressible flow solvers based on a finite element discretisation. *International Journal of Computing Science and Mathematics*, 1(2-4):193–206, 2007.
- [60] P Lax and A Milgram. Ix. parabolic equations. *Contributions to the Theory of Partial Differential Equations.(AM-33)*, 33:167, 2016.
- [61] Sérgio L. Frey Leopoldo P. Franca. Stabilized finite element methods:i. application to the advective-diffusive model. *Computer Methods in Applied Mechanics and Engineering*, 95:253–276, 1992.
- [62] Sérgio L. Frey Leopoldo P. Franca. Stabilized finite element methods:ii. the incompressible navier-stokes equations. *Computer Methods in Applied Mechanics and Engineering*, 99:209–233, 1992.
- [63] Luc Machiels, Yvon Maday, Ivan B. Oliveira, Anthony T. Patera, and Dimitrios V. Rovas. Output bounds for reduced-basis approximations of symmetric positive definite eigenvalue problems. *Comptes Rendus de l'Académie des Sciences - Series I - Mathematics*, 331(2):153 – 158, 2000.
- [64] Yvon Maday, Andrea Manzoni, and Alfio Quarteroni. An online intrinsic stabilization strategy for the reduced basis approximation of parametrized advection-dominated problems. *Comptes Rendus Mathématique*, 354(12):1188–1194, 2016.
- [65] Yvon Maday, Ngoc Cuong Nguyen, Anthony T. Patera, and George S.H. Pau. A general, multipurpose interpolation procedure: the magic points. working paper or preprint, September 2007.
- [66] Yvon Maday, A.T. Patera, and Dimitrios Rovas. A blackbox reduced-basis output bound method for noncoercive linear problems. *Studies in Mathematics and Its Applications*, 31, 12 2002.
- [67] MT Manzari. An explicit finite element algorithm for convection heat transfer problems. *International Journal of Numerical Methods for Heat & Fluid Flow*, 9(8):860–877, 1999.

- [68] Andrea Manzoni. An efficient computational framework for reduced basis approximation and a posteriori error estimation of parametrized navier–stokes flows. *ESAIM: Mathematical Modelling and Numerical Analysis*, 48(4):1199–1226, 2014.
- [69] David A Mayne, Asif S Usmani, and Martin Crapper. h-adaptive finite element solution of high rayleigh number thermally driven cavity problem. *International Journal of Numerical Methods for Heat & Fluid Flow*, 10(6):598–615, 2000.
- [70] F Menter. Zonal two equation kw turbulence models for aerodynamic flows. In *23rd fluid dynamics, plasmadynamics, and lasers conference*, page 2906, 1993.
- [71] Florian R Menter. Two-equation eddy-viscosity turbulence models for engineering applications. *AIAA journal*, 32(8):1598–1605, 1994.
- [72] FR Menter, M Kuntz, and R Langtry. Ten years of industrial experience with the sst turbulence model. *Turbulence, heat and mass transfer*, 4(1):625–632, 2003.
- [73] S. Mittal. On the performance of high aspect ratio elements for incompressible flows. *Computer Methods in Applied Mechanics and Engineering*, 188(1–3):269–287, 2000.
- [74] Akira Mizukami and Thomas JR Hughes. A petrov-galerkin finite element method for convection-dominated flows: an accurate upwinding technique for satisfying the maximum principle. *Computer Methods in Applied Mechanics and Engineering*, 50(2):181–193, 1985.
- [75] Federico Negri, Andrea Manzoni, and David Amsallem. Efficient model reduction of parametrized systems by matrix discrete empirical interpolation. *Journal of Computational Physics*, 303:431–454, 2015.
- [76] Ahmed K Noor. Recent advances in reduction methods for nonlinear problems. *Computers & Structures*, 13(1-3):31–44, 1981.
- [77] Y. Papadopoulos. A driven cavity exploration. <https://www.acenumerics.com/the-cavity-sessions.html>.
- [78] Janet S Peterson. The reduced basis method for incompressible viscous flow calculations. *SIAM Journal on Scientific and Statistical Computing*, 10(4):777–786, 1989.
- [79] Olivier Pironneau and Olivier Pironneau. *Finite element methods for fluids*. Wiley Chichester, 1989.
- [80] C. Prud’homme, D. V. Rovas, K. Veroy, L. Machiels, Y. Maday, A. T. Patera, and G. Turinici. Reliable real-time solution of parametrized partial differential equations: Reduced-basis output bound methods. *Journal of Fluids Engineering*, 124(1):70–80, 2002.

- [81] Christophe Prud'homme and AT Patera. Reduced-basis output bounds for approximately parametrized elliptic coercive partial differential equations. *Computing and Visualization in Science*, 6(2):147–162, 2004.
- [82] Alfio Quarteroni and Gianluigi Rozza. Numerical solution of parametrized navier–stokes equations by reduced basis methods. *Numerical Methods for Partial Differential Equations: An International Journal*, 23(4):923–948, 2007.
- [83] Alfio Quarteroni, Gianluigi Rozza, and Andrea Manzoni. Certified reduced basis approximation for parametrized partial differential equations and applications. *Journal of Mathematics in Industry*, 1(1):3, 2011.
- [84] Werner C Rheinboldt. On the theory and error estimation of the reduced basis method for multi-parameter problems. *Nonlinear Analysis: Theory, Methods & Applications*, 21(11):849–858, 1993.
- [85] D. V. Rovas. Reduced-basis output bound methods for parametrized partial differential equations. *Ph.D. thesis, Massachusetts Institute of Technology, Cambridge*, 2003.
- [86] Dimitrios Vasileios Rovas. *Reduced-basis output bound methods for parametrized partial differential equations*. PhD thesis, Massachusetts Institute of Technology, 2003.
- [87] Gianluigi Rozza, DB Phuong Huynh, and Andrea Manzoni. Reduced basis approximation and a posteriori error estimation for stokes flows in parametrized geometries: roles of the inf-sup stability constants. *Numerische Mathematik*, 125(1):115–152, 2013.
- [88] Gianluigi Rozza, Dinh Bao Phuong Huynh, and Anthony T Patera. Reduced basis approximation and a posteriori error estimation for affinely parametrized elliptic coercive partial differential equations. *Archives of Computational Methods in Engineering*, 15(3):1, 2007.
- [89] Gianluigi Rozza, Dinh Bao Phuong Huynh, and Anthony T Patera. Reduced basis approximation and a posteriori error estimation for affinely parametrized elliptic coercive partial differential equations. *Archives of Computational Methods in Engineering*, 15(3):229, 2008.
- [90] Gianluigi Rozza and Karen Veroy. On the stability of the reduced basis method for stokes equations in parametrized domains. *Computer methods in applied mechanics and engineering*, 196(7):1244–1260, 2007.
- [91] E. Schenone, S. Veys, and C. Prud'Homme. High Performance Computing for the Reduced Basis Method. Application to Natural Convection. *ESAIM: Proceedings*, pages 255 – 273, December 2013.
- [92] PN Shankar and MD Deshpande. Fluid mechanics in the driven cavity. *Annual Review of Fluid Mechanics*, 32(1):93–136, 2000.

- [93] Ratnesh K Shukla, Mahidhar Tatineni, and Xiaolin Zhong. Very high-order compact finite difference schemes on non-uniform grids for incompressible navier–stokes equations. *Journal of Computational Physics*, 224(2):1064–1094, 2007.
- [94] Pavel E Smirnov and Florian R Menter. Sensitization of the sst turbulence model to rotation and curvature by applying the spalart–shur correction term. *Journal of Turbomachinery*, 131(4):041010, 2009.
- [95] Philippe R Spalart and Christopher L Rumsey. Effective inflow conditions for turbulence models in aerodynamic calculations. *AIAA journal*, 45(10):2544–2553, 2007.
- [96] PRaA Spalart and S1 Allmaras. A one-equation turbulence model for aerodynamic flows. In *30th aerospace sciences meeting and exhibit*, page 439, 1992.
- [97] T Tonn. Reduced-basis method (rbm) for non-affine elliptic parametrized pdes.(phd). *Ulm University*, 2012.
- [98] David J Tritton. *Physical fluid dynamics*. Springer Science & Business Media, 2012.
- [99] C. T. Shin U. Ghia, K. N. Ghia. High-re solutions for incompressible flow using the navier-stokes equations and a multigrid method. *Journal of Computational Physics*, 48:387–411, 1982.
- [100] José M Urquiza, André Garon, and Marie-Isabelle Farinas. Weak imposition of the slip boundary condition on curved boundaries for stokes flow. *Journal of Computational Physics*, 256:748–767, 2014.
- [101] Sylvain Vallaghé, Michel Fouquembergh, Annabelle Le Hyaric, and Christophe Prud’Homme. A successive constraint method with minimal off-line constraints for lower bounds of parametric coercivity constant. working paper or preprint, June 2011.
- [102] K. Veroy and A.T. Patera. Certified real-time solution of the parametrized steady incompressible navier-stokes equations: Rigorous reduced-basis a posteriori error bounds. *Int. J. Numer. Meth. Fluids*, 47:773–788, 2005.
- [103] Karen Veroy, Christophe Prud’Homme, Dimitrios V Rovas, and Anthony T Patera. A posteriori error bounds for reduced-basis approximation of parametrized noncoercive and nonlinear elliptic partial differential equations. In *Proceedings of the 16th AIAA computational fluid dynamics conference*, volume 3847, pages 23–26. Orlando, FL, 2003.
- [104] Karen Veroy, Dimitrios V Rovas, and Anthony T Patera. A posteriori error estimation for reduced-basis approximation of parametrized elliptic coercive partial differential equations: “convex inverse” bound conditioners. *ESAIM: Control, Optimisation and Calculus of Variations*, 8:1007–1028, 2002.



- [105] Stéphane Veys. *A computational reduced basis framework : applications to nonlinear multiphysics problems*. Theses, Université Joseph Fourier (Grenoble I), November 2014.
- [106] D.C. Wan, B.S.V. Patnaik, and G.W. Wei. A new benchmark quality solution for the buoyancy-driven cavity by discrete singular convolution. *Numerical Heat Transfer: Part B: Fundamentals*, 40(3):199–228, 2001.
- [107] Zhu Wang, Imran Akhtar, Jeff Borggaard, and Traian Iliescu. Proper orthogonal decomposition closure models for turbulent flows: a numerical comparison. *Computer Methods in Applied Mechanics and Engineering*, 237:10–26, 2012.
- [108] K Yapici and Y Uludag. Finite volume simulation of 2-d steady square lid driven cavity flow at high reynolds numbers. *Brazilian Journal of Chemical Engineering*, 30(4):923–937, 2013.
- [109] Olivier Zahm, Marie Billaud-Friess, and Anthony Nouy. Projection-based model order reduction methods for the estimation of vector-valued variables of interest. *SIAM Journal on Scientific Computing*, 39(4):A1647–A1674, 2017.



**Summary:** We present in this thesis our work on model order reduction for aerothermal simulations. We consider the coupling between the incompressible Navier-Stokes equations and an advection-diffusion equation for the temperature. Since the physical parameters induce high Reynolds and Peclet numbers, we have to introduce stabilization operators in the formulation to deal with the well known numerical stability issue. The chosen stabilization, applied to both fluid and heat equations, is the usual Streamline-Upwind/Petrov-Galerkin (SUPG) which add artificial diffusivity in the direction of the convection field. We also introduce our order reduction strategy for this model, based on the Reduced Basis Method (RBM). To recover an affine decomposition for this complex model, we implemented a discrete variation of the Empirical Interpolation Method (EIM) which is a discrete version of the original EIM. This variant allows building an approximated affine decomposition for complex operators such as in the case of SUPG. We also use this method for the non-linear operators induced by the shock capturing method. The construction of a EIM basis for non-linear operators involves a potentially huge number of non-linear FEM resolutions - depending on the size of the sampling. Even if this basis is built during an offline phase, we usually can not afford such expensive computational cost. We took advantage of the recent development of the Simultaneous EIM Reduced basis algorithm (SER) to tackle this issue.

**Keywords:** Model order reduction, Computational fluid dynamic, Reduced Basis Method, Empirical Interpolation Method, Finite element method, HPC computing

AGENCE NATIONALE DE LA RECHERCHE  
**ANR**  
 UNIVERSITÉ DE STRASBOURG  
 CNRS  
 IRMA  
 Institut de Recherche  
 Mathématique Avancée

INSTITUT DE RECHERCHE MATHÉMATIQUE AVANCÉE  
 UMR 7501  
 Université de Strasbourg et CNRS  
 7 Rue René Descartes  
 67 084 STRASBOURG CEDEX

Tél. 03 68 85 01 29  
 Fax 03 68 85 03 28  
[www-irma.u-strasbg.fr](http://www-irma.u-strasbg.fr)  
[irma@math.unistra.fr](mailto:irma@math.unistra.fr)

IRMA 2018/003  
<http://tel.archives-ouvertes.fr/tel-01869098>  
 ISSN 0755-3390