



**HAL**  
open science

# The T-coercivity approach for mixed problems

Mathieu Barré, Patrick Ciarlet

► **To cite this version:**

Mathieu Barré, Patrick Ciarlet. The T-coercivity approach for mixed problems. 2022. hal-03820910v1

**HAL Id: hal-03820910**

**<https://hal.science/hal-03820910v1>**

Preprint submitted on 19 Oct 2022 (v1), last revised 24 Sep 2024 (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.



---

# The T-coercivity approach for mixed problems

## *T-coercivité et problèmes mixtes*

Mathieu Barré<sup>a, b</sup> and Patrick Ciarlet<sup>c</sup>

<sup>a</sup> Inria, 1 Rue Honoré d'Estienne d'Orves, 91120 Palaiseau, France

<sup>b</sup> LMS, École Polytechnique, CNRS, Institut Polytechnique de Paris, Route de Saclay, 91120 Palaiseau, France

<sup>c</sup> POEMS, CNRS, Inria, ENSTA Paris, Institut Polytechnique de Paris, 828 Boulevard des Maréchaux, 91120 Palaiseau, France

*E-mails:* mathieu.a.barre@inria.fr (M. Barré), patrick.ciarlet@ensta-paris.fr (P. Ciarlet)

**Abstract.** Classically, the well-posedness of variational formulations of mixed linear problems is achieved through the inf-sup condition. In this note, we propose an alternative framework to study such problems by using the T-coercivity approach. This is a constructive approach that leads to the design of suitable approximations in a simple way. In general, the derivation of the uniform discrete inf-sup condition for the approximate problems stems straightforwardly from the study of the original problem. To support our view, we solve a series of classical mixed problems with the T-coercivity approach. Among others, the celebrated Fortin Lemma appears naturally in the numerical analysis of the approximate problems.

**Résumé.** Classiquement, le caractère bien posé des formulations variationnelles de problèmes linéaires mixtes est obtenu à l'aide de la condition inf-sup. Dans cette note, nous proposons un cadre alternatif pour étudier de tels problèmes en utilisant la notion de T-coercivité. Il s'agit d'une approche constructive qui permet en outre de concevoir simplement des approximations numériques adaptées car la dérivation de la condition inf-sup discrète uniforme découle en général directement de l'étude du problème continu. Pour appuyer notre propos, nous résolvons une série de problèmes mixtes classiques grâce à la notion de T-coercivité. Entre autres, le lemme de Fortin apparaît naturellement dans l'analyse numérique des problèmes discrets.

**2020 Mathematics Subject Classification.** 65N30, 35J57, 76D07, 78M10.

*This article is a draft (not yet accepted!)*

## 1. Introduction

Traditionally, the well-posedness of variational formulations of mixed linear problems is achieved through the inf-sup condition, also called stability condition [2, 11, 34]. As a matter of fact, proving this condition allows to derive existence and uniqueness of the solution, and continuous dependence with respect to the data. On the other hand, the way this condition is established depends on the problem to be solved. The analysis of such problems can be performed either following a *monolithic* approach, namely studying the *all-in-one* bilinear form incorporating the constraint, or by studying the constrained part of the problem separately.

In this note, we focus on the monolithic approach and investigate the mixed problem's well-posedness based on the T-coercivity framework. The principle of this framework is to find an *explicit realization* of the inf-sup condition for the all-in-one bilinear form. Of equal importance, in the T-coercivity framework, is the *design of suitable approximations* of the original problem. Indeed, with the help of the explicit realization of the condition for the original problem, one can get useful insight on how to derive the so-called uniform discrete inf-sup condition for the approximate, or discrete, problems set in finite-dimensional vector spaces. Thus, convergence of the approximate solutions to the exact one follows under well-known principles in numerical analysis, such as Céa's Lemma (or a variant), and a basic approximability property of elements of the original space of solutions. To summarize, although the T-coercivity approach may not bring new result to the theory of variational formulations, it proposes an entirely constructive way to study them theoretically and also on how to approximate them efficiently.

So far, the T-coercivity approach has been mainly applied to two categories of linear problems. First, for problems involving an invertible operator and a compact perturbation, see eg. [12–14, 20, 30, 37]. Then, for problems with sign-changing coefficients, cf. [5–10, 15–19, 21, 23, 29, 36]. For the second category, we observe that well-posedness and (efficient) approximation of the variational formulations has actually been achieved with the help of the T-coercivity approach. Up to the authors' knowledge, this approach was only applied to mixed problems in [31, 33]: in the first reference, it is applied to the specific case of neutron diffusion, whereas the second one focuses on perturbed saddle-point problems.

In this note, we apply the T-coercivity approach to general mixed problems, including unperturbed and perturbed saddle-point problems. In particular, we will explain the connections with the classical theory, for which we use [4] as the reference textbook. Among those connections, we note that the celebrated Fortin Lemma will appear naturally in the (numerical) analysis of the discrete problems.

Let us introduce some notations. Given a Hilbert space  $V$ , we denote by  $(\cdot, \cdot)_V$  and  $\|\cdot\|_V$  the scalar product and the norm on  $V$ , and by  $V'$  its dual space. In a product space  $V \times W$ , we use the norm

$$\|(v, w)\|_{V \times W} = (\|v\|_V^2 + \|w\|_W^2)^{1/2},$$

and similarly for the scalar product. Vector-valued function spaces are written in boldface character. A connected, bounded, open subset of  $\mathbb{R}^d$  with a Lipschitz boundary is called a *domain*.

Let  $\Omega$  be a domain with boundary  $\partial\Omega$ . We denote by  $\mathbf{n}$  the unit outward normal vector field to  $\partial\Omega$ . Let  $L^2(\Omega)$  and  $\mathbf{L}^2(\Omega)$  be the set of square-integrable real-valued and  $\mathbb{R}^d$ -valued functions on  $\Omega$ . The natural norm in  $L^2(\Omega)$  or  $\mathbf{L}^2(\Omega)$  is denoted by  $\|\cdot\|$ , and we let

$$L_0^2(\Omega) = \left\{ v \in L^2(\Omega), \int_{\Omega} v \, dx = 0 \right\}.$$

In what follows, unless otherwise stated, the standard Sobolev space  $H_0^1(\Omega)$  is endowed with the norm  $v \mapsto \|\nabla v\|$ , that defines a norm that is equivalent to  $\|\cdot\|_{H_0^1(\Omega)}$  thanks to Poincaré's inequality. The dual space of  $H_0^1(\Omega)$  is denoted by  $H^{-1}(\Omega)$ . Similarly,  $\mathbf{H}_0^1(\Omega)$  is endowed with the norm  $\mathbf{v} \mapsto (\sum_{i=1,d} \|\nabla v_i\|^2)^{1/2}$ , that defines a norm that is equivalent to  $\|\cdot\|_{\mathbf{H}_0^1(\Omega)}$ , and its dual space

is denoted by  $\mathbf{H}^{-1}(\Omega)$ . We introduce the usual Sobolev spaces for vector-valued fields [1]

$$\begin{aligned}\mathbf{H}(\operatorname{div}; \Omega) &= \{\mathbf{v} \in \mathbf{L}^2(\Omega), \operatorname{div} \mathbf{v} \in L^2(\Omega)\}, \\ \mathbf{H}_0(\operatorname{div}; \Omega) &= \{\mathbf{v} \in \mathbf{H}(\operatorname{div}; \Omega), \mathbf{v} \cdot \mathbf{n} = 0 \text{ on } \partial\Omega\}, \\ \mathbf{H}(\operatorname{div}0; \Omega) &= \{\mathbf{v} \in \mathbf{H}(\operatorname{div}; \Omega), \operatorname{div} \mathbf{v} = 0\}, \\ \mathbf{H}(\operatorname{curl}; \Omega) &= \{\mathbf{v} \in \mathbf{L}^2(\Omega), \operatorname{curl} \mathbf{v} \in \mathbf{L}^2(\Omega)\}, \quad \text{for } d = 3, \\ \mathbf{H}_0(\operatorname{curl}; \Omega) &= \{\mathbf{v} \in \mathbf{H}(\operatorname{curl}; \Omega), \mathbf{v} \times \mathbf{n} = 0 \text{ on } \partial\Omega\}, \quad \text{for } d = 3.\end{aligned}$$

Unless otherwise specified,  $\mathbf{H}(\operatorname{div}; \Omega)$  is endowed with the norm  $\mathbf{v} \mapsto (\|\mathbf{v}\|^2 + \|\operatorname{div} \mathbf{v}\|^2)^{1/2}$  and  $\mathbf{H}(\operatorname{curl}; \Omega)$  with the norm  $\mathbf{v} \mapsto (\|\mathbf{v}\|^2 + \|\operatorname{curl} \mathbf{v}\|^2)^{1/2}$ .

The outline is as follows. In Section 2, we introduce the T-coercivity approach, and explain how it can be applied to solve the Stokes problem theoretically. Then, in Section 3, we develop the abstract framework underlying the approach for mixed problems, including saddle-point, augmented and perturbed ones. In Sections 4, 5 and 6, we propose some applications, respectively to electromagnetism, nearly-incompressible elasticity, and diffusion. Then, in Section 7, we propose the *natural* extension of the T-coercivity approach for the conforming approximation of mixed problems. As before, we begin by the Stokes problem, then we consider the numerical analysis for mixed problems in general, before describing how the approach can be applied to electromagnetism, nearly-incompressible elasticity, and diffusion. We conclude by a list of further extensions and recent applications of the T-coercivity approach.

## 2. T-coercivity for the Stokes problem

The starting point of our study is to propose a T-coercivity approach to solve Stokes problem. Let  $\Omega \subset \mathbb{R}^d$  be a domain. We consider the Stokes problem with homogeneous Dirichlet boundary conditions: given a prescribed body force  $\mathbf{f} \in \mathbf{H}^{-1}(\Omega)$ , find the velocity  $\mathbf{u} \in \mathbf{H}^1(\Omega)$  and the pressure  $p \in L_0^2(\Omega)$  such that

$$\begin{aligned}-\nu \Delta \mathbf{u} + \nabla p &= \mathbf{f}, & \text{in } \Omega, \\ \operatorname{div} \mathbf{u} &= 0, & \text{in } \Omega, \\ \mathbf{u} &= 0, & \text{on } \partial\Omega,\end{aligned}\tag{1}$$

where  $\nu > 0$  denotes the fluid's viscosity.

The standard method to solve Problem (1) – see [27] – consists in a *one-plus-one* approach. The problem is split into a coercive part

$$a(\mathbf{u}, \mathbf{v}) = \nu \int_{\Omega} \nabla \mathbf{u} : \nabla \mathbf{v} \, dx$$

and divergence constraint terms of the form

$$b(\mathbf{v}, q) = - \int_{\Omega} q \operatorname{div} \mathbf{v} \, dx,$$

so that the weak formulation of Problem (1) reads: find  $(\mathbf{u}, p) \in \mathbf{H}_0^1(\Omega) \times L_0^2(\Omega)$  such that

$$\begin{aligned}a(\mathbf{u}, \mathbf{v}) + b(\mathbf{v}, p) &= \langle \mathbf{f}, \mathbf{v} \rangle, \quad \forall \mathbf{v} \in \mathbf{H}_0^1(\Omega), \\ b(\mathbf{u}, q) &= 0, \quad \forall q \in L_0^2(\Omega),\end{aligned}\tag{2}$$

where  $\langle \cdot, \cdot \rangle$  denotes the duality product in  $\mathbf{H}_0^1(\Omega)$ . The well-posedness of Problem (2) then follows from Ladyzhenskaya–Babuška–Brezzi's theory [2, 11, 34] since the bilinear form  $a$  is coercive on  $\mathbf{H}_0^1(\Omega)$  and the bilinear form  $b$  satisfies the *inf-sup condition*

$$\inf_{q \in L_0^2(\Omega) \setminus \{0\}} \sup_{\mathbf{v} \in \mathbf{H}_0^1(\Omega) \setminus \{0\}} \frac{b(\mathbf{v}, q)}{\|\nabla \mathbf{v}\| \|q\|} \geq \underline{\beta}\tag{3}$$

for some constant  $\beta > 0$ .

Here, we are going to give an alternative proof that Problem (1) is well-posed by analysing the *all-in-one* bilinear form defined on  $\mathbf{H}_0^1(\Omega) \times L_0^2(\Omega)$

$$\mathcal{A}((\mathbf{u}, p), (\mathbf{v}, q)) = \nu \int_{\Omega} \nabla \mathbf{u} : \nabla \mathbf{v} \, dx - \int_{\Omega} p \operatorname{div} \mathbf{v} \, dx - \int_{\Omega} q \operatorname{div} \mathbf{u} \, dx$$

instead of splitting it into two bilinear forms  $a$  and  $b$  as in (2). This bilinear form is not coercive since

$$\mathcal{A}((0, p), (0, p)) = 0, \quad \forall p \in L_0^2(\Omega).$$

For this reason, we use the notion of T-coercivity [19, 20], which can be seen as a reformulation of Banach-Nečas-Babuška's theory. The definition and the main property of T-coercivity are recalled below.

**Definition 1.** *Let  $W$  be a Hilbert space and let  $\mathcal{A}(\cdot, \cdot)$  be a continuous bilinear form over  $W \times W$ . We say that  $\mathcal{A}$  is T-coercive if there exists a bijective operator  $\mathbb{T} \in \mathcal{L}(W)$  and  $\underline{\alpha} > 0$  such that*

$$|\mathcal{A}(u, \mathbb{T}u)| \geq \underline{\alpha} \|u\|_W^2, \quad \forall u \in W.$$

**Proposition 2.** *Let  $W$  be a Hilbert space. Let  $\ell(\cdot)$  be a continuous linear form over  $W$  and  $\mathcal{A}(\cdot, \cdot)$  be a continuous bilinear form over  $W \times W$ . The problem*

$$\begin{cases} \text{Find } u \in W & \text{such that} \\ \forall v \in W, & \mathcal{A}(u, v) = \ell(v) \end{cases}$$

is well-posed if and only if  $\mathcal{A}$  is T-coercive. If so, it holds that

$$\|u\|_W \leq \frac{C_\ell}{\underline{\alpha}} \|\mathbb{T}\|, \quad (4)$$

with  $C_\ell$  the continuity constant of the linear form  $\ell$ . When the bilinear form  $\mathcal{A}(\cdot, \cdot)$  is in addition symmetric, the requirement that the operator  $\mathbb{T}$  is bijective can be dropped.

### 2.1. Proving well-posedness with T-coercivity

With the T-coercivity tool in mind, we are now ready to establish the main result of this section.

**Theorem 3.** *The problem*

$$\begin{cases} \text{Find } (\mathbf{u}, p) \in \mathbf{H}_0^1(\Omega) \times L_0^2(\Omega) & \text{such that} \\ \forall (\mathbf{v}, q) \in \mathbf{H}_0^1(\Omega) \times L_0^2(\Omega), & \mathcal{A}((\mathbf{u}, p), (\mathbf{v}, q)) = \langle \mathbf{f}, \mathbf{v} \rangle \end{cases} \quad (5)$$

is well-posed and

$$\|(\mathbf{u}, p)\|_{\mathbf{H}_0^1(\Omega) \times L_0^2(\Omega)} \leq \frac{2 \max\left(\sqrt{2} \nu C_{\operatorname{div}}^2, C_{\operatorname{div}}(2 + \nu^2 C_{\operatorname{div}}^2)^{1/2}\right)}{\min(\nu^2 C_{\operatorname{div}}^2, 1)} \|\mathbf{f}\|_{\mathbf{H}^{-1}(\Omega)}. \quad (6)$$

**Proof.** The linear form defined by

$$\ell((\mathbf{v}, q)) = \langle \mathbf{f}, \mathbf{v} \rangle, \quad \forall (\mathbf{v}, q) \in \mathbf{H}_0^1(\Omega) \times L_0^2(\Omega)$$

is continuous over  $\mathbf{H}_0^1(\Omega) \times L_0^2(\Omega)$  in view of the inequality

$$\ell((\mathbf{v}, q)) \leq \|\mathbf{f}\|_{\mathbf{H}^{-1}(\Omega)} \|(\mathbf{v}, q)\|_{\mathbf{H}_0^1(\Omega) \times L_0^2(\Omega)}. \quad (7)$$

The bilinear form  $\mathcal{A}$  is continuous over  $(\mathbf{H}_0^1(\Omega) \times L_0^2(\Omega))^2$  and we observe that it is also symmetric.

Then, from Proposition 2, it is sufficient to show that the bilinear form  $\mathcal{A}$  is T-coercive. For a given  $(\mathbf{u}, p) \in \mathbf{H}_0^1(\Omega) \times L_0^2(\Omega)$ , we look for an element  $(\mathbf{v}^*, q^*)$  of  $\mathbf{H}_0^1(\Omega) \times L_0^2(\Omega)$  depending continuously on  $(\mathbf{u}, p)$  and such that

$$\mathcal{A}((\mathbf{u}, p), (\mathbf{v}^*, q^*)) \geq \underline{\alpha} \|(\mathbf{u}, p)\|_{\mathbf{H}_0^1(\Omega) \times L_0^2(\Omega)}^2$$

for some constant  $\underline{\alpha} > 0$ . In order to get an intuitive idea of the construction of  $(\mathbf{v}^*, q^*)$ , let us start with specific elements  $(\mathbf{u}, p)$ .

- If  $p = 0$ , then  $\|(\mathbf{u}, p)\|_{\mathbf{H}_0^1(\Omega) \times L_0^2(\Omega)}^2 = \|\nabla \mathbf{u}\|^2$  and

$$\mathcal{A}((\mathbf{u}, p), (\mathbf{v}^*, q^*)) = \nu \int_{\Omega} \nabla \mathbf{u} : \nabla \mathbf{v}^* \, dx - \int_{\Omega} \operatorname{div} \mathbf{u} \, q^* \, dx,$$

so that we can take  $\mathbf{v}^* = \mathbf{u}$  and  $q^* = p = 0$ .

- If  $\mathbf{u} = 0$ , then  $\|(\mathbf{u}, p)\|_{\mathbf{H}_0^1(\Omega) \times L_0^2(\Omega)}^2 = \|p\|^2$  and

$$\mathcal{A}((\mathbf{u}, p), (\mathbf{v}^*, q^*)) = - \int_{\Omega} p \operatorname{div} \mathbf{v}^* \, dx.$$

In order to recover the expected term  $\|p\|^2$  in the above expression, we have to choose  $\mathbf{v}^*$ , the divergence of which is "as close as possible" to  $p$ . To that aim, we use the result below, see for instance [27, Corollary I.2.4]. Let  $p \in L_0^2(\Omega)$ . Then, there exists  $\mathbf{v}_p \in \mathbf{H}_0^1(\Omega)$  satisfying

$$- \operatorname{div} \mathbf{v}_p = p. \quad (8)$$

In addition, there exists a constant  $C_{\operatorname{div}} > 0$  independent of  $p$  such that

$$\|\nabla \mathbf{v}_p\| \leq C_{\operatorname{div}} \|p\|. \quad (9)$$

The idea is now to choose  $\mathbf{v}^* = \mathbf{v}_p$ , where  $\mathbf{v}_p$  is as in (8)-(9). Hence, taking  $q^* = 0$ , we find

$$\mathcal{A}((\mathbf{u}, p), (\mathbf{v}^*, q^*)) = \|p\|^2,$$

and (9) ensures that the pair  $(\mathbf{v}_p, 0)$  depends continuously on  $(0, p)$  in  $\mathbf{H}_0^1(\Omega) \times L_0^2(\Omega)$ .

- If  $\operatorname{div} \mathbf{u} = 0$ , then

$$\mathcal{A}((\mathbf{u}, p), (\mathbf{v}^*, q^*)) = \nu \int_{\Omega} \nabla \mathbf{u} : \nabla \mathbf{v}^* \, dx - \int_{\Omega} p \operatorname{div} \mathbf{v}^* \, dx.$$

Since we need to get a term of the form  $\|\nabla \mathbf{u}\|^2$  but also of the form  $\|p\|^2$ , we combine the previous two cases by setting  $\mathbf{v}^* = \lambda \mathbf{u} + \mathbf{v}_p$ , where  $\lambda$  is a positive coefficient to be adjusted and  $\mathbf{v}_p$  is the divergence lifting from (8) – (9). Now, we compute

$$\begin{aligned} \mathcal{A}((\mathbf{u}, p), (\mathbf{v}^*, q^*)) &= \nu \lambda \int_{\Omega} \nabla \mathbf{u} : \nabla \mathbf{u} \, dx + \nu \int_{\Omega} \nabla \mathbf{u} : \nabla \mathbf{v}_p \, dx - \lambda \int_{\Omega} p \operatorname{div} \mathbf{u} \, dx - \int_{\Omega} p \operatorname{div} \mathbf{v}_p \, dx \\ &= \nu \lambda \|\nabla \mathbf{u}\|^2 + \nu \int_{\Omega} \nabla \mathbf{u} : \nabla \mathbf{v}_p \, dx + \|p\|^2 \end{aligned}$$

since  $\operatorname{div} \mathbf{u} = 0$  and  $-\operatorname{div} \mathbf{v}_p = p$ . For all  $\eta > 0$ , Young's inequality implies that

$$\begin{aligned} \int_{\Omega} \nabla \mathbf{u} : \nabla \mathbf{v}_p \, dx &\geq -\frac{\eta}{2} \|\nabla \mathbf{u}\|^2 - \frac{1}{2\eta} \|\nabla \mathbf{v}_p\|^2 \\ &\geq -\frac{\eta}{2} \|\nabla \mathbf{u}\|^2 - \frac{C_{\operatorname{div}}^2}{2\eta} \|p\|^2 \quad \text{in virtue of (9),} \end{aligned}$$

and thus

$$\mathcal{A}((\mathbf{u}, p), (\mathbf{v}^*, q^*)) \geq \nu \left( \lambda - \frac{\eta}{2} \right) \|\nabla \mathbf{u}\|^2 + \left( 1 - \frac{\nu C_{\operatorname{div}}^2}{2\eta} \right) \|p\|^2.$$

Hence, by setting  $\eta = \lambda = \nu C_{\operatorname{div}}^2$ , we obtain

$$\mathcal{A}((\mathbf{u}, p), (\mathbf{v}^*, q^*)) \geq \frac{\nu^2 C_{\operatorname{div}}^2}{2} \|\nabla \mathbf{u}\|^2 + \frac{1}{2} \|p\|^2.$$

Note that this result holds for any  $q^* \in L_0^2(\Omega)$  depending continuously on  $p$ .

In the general case, we choose  $\mathbf{v}^* = \lambda \mathbf{u} + \mathbf{v}_p$  with  $\lambda = \nu C_{\text{div}}^2$  and  $q^* = -\lambda p$  so that, even if  $\text{div } \mathbf{u} \neq 0$ , the term  $-\lambda \int_{\Omega} p \text{div } \mathbf{u} dx$  cancels with the term  $-\int_{\Omega} \text{div } \mathbf{u} q^* dx$  and we get the same results as in the case  $\text{div } \mathbf{u} = 0$ . Namely, the bilinear form  $\mathcal{A}$  is T-coercive for the mapping

$$\begin{aligned} \mathbb{T} : \mathbf{H}_0^1(\Omega) \times L_0^2(\Omega) &\longrightarrow \mathbf{H}_0^1(\Omega) \times L_0^2(\Omega) \\ (\mathbf{u}, p) &\longmapsto (\nu C_{\text{div}}^2 \mathbf{u} + \mathbf{v}_p, -\nu C_{\text{div}}^2 p), \end{aligned}$$

where  $\mathbf{v}_p$  is defined by (8) with estimate (9), and it holds that

$$\mathcal{A}((\mathbf{u}, p), \mathbb{T}(\mathbf{u}, p)) \geq \frac{\nu^2 C_{\text{div}}^2}{2} \|\nabla \mathbf{u}\|^2 + \frac{1}{2} \|p\|^2 \geq \frac{1}{2} \min(\nu^2 C_{\text{div}}^2, 1) \|(\mathbf{u}, p)\|_{\mathbf{H}_0^1(\Omega) \times L_0^2(\Omega)}^2. \quad (10)$$

Thanks to (8)-(9),  $\mathbb{T}$  belongs to  $\mathcal{L}(\mathbf{H}_0^1(\Omega) \times L_0^2(\Omega))$ . More precisely, we have

$$\begin{aligned} \|\mathbb{T}(\mathbf{u}, p)\|_{\mathbf{H}_0^1(\Omega) \times L_0^2(\Omega)}^2 &= \|\nu C_{\text{div}}^2 \mathbf{u} + \mathbf{v}_p\|_{\mathbf{H}_0^1(\Omega)}^2 + \|\nu C_{\text{div}}^2 p\|^2 \\ &\leq 2(\nu C_{\text{div}}^2)^2 \|\nabla \mathbf{u}\|^2 + 2\|\nabla \mathbf{v}_p\|^2 + (\nu C_{\text{div}}^2)^2 \|p\|^2 \\ &\leq 2(\nu C_{\text{div}}^2)^2 \|\nabla \mathbf{u}\|^2 + (2C_{\text{div}}^2 + (\nu C_{\text{div}}^2)^2) \|p\|^2 \end{aligned}$$

and thus

$$\|\mathbb{T}\| \leq \max\left(\sqrt{2}\nu C_{\text{div}}^2, C_{\text{div}}(2 + \nu^2 C_{\text{div}}^2)^{1/2}\right). \quad (11)$$

Using (7), (10) and (11) in the stability estimate (4), we finally obtain (6).  $\square$

**Remark 4.** The previous result readily extends to the case of a non-null divergence constraint

$$\begin{aligned} -\nu \Delta \mathbf{u} + \nabla p &= \mathbf{f}, & \text{in } \Omega, \\ \text{div } \mathbf{u} &= g, & \text{in } \Omega, \\ \mathbf{u} &= 0, & \text{on } \partial\Omega, \end{aligned}$$

with  $g \in L_0^2(\Omega)$ , leading to the stability estimate

$$\|(\mathbf{u}, p)\|_{\mathbf{H}_0^1(\Omega) \times L_0^2(\Omega)} \leq \frac{2 \max\left(\sqrt{2}\nu C_{\text{div}}^2, C_{\text{div}}(2 + \nu^2 C_{\text{div}}^2)^{1/2}\right)}{\min(\nu^2 C_{\text{div}}^2, 1)} \|(\mathbf{f}, g)\|_{\mathbf{H}^{-1}(\Omega) \times L_0^2(\Omega)}. \quad (12)$$

## 2.2. Comments

The stability estimates (6) and (12) are valid for all  $C_{\text{div}}$  that fulfills (9). On the other hand, one has

$$\lim_{C_{\text{div}} \rightarrow \infty} \frac{2 \max\left(\sqrt{2}\nu C_{\text{div}}^2, C_{\text{div}}(2 + \nu^2 C_{\text{div}}^2)^{1/2}\right)}{\min(\nu^2 C_{\text{div}}^2, 1)} = +\infty,$$

*i.e.* the stability estimates become meaningless for large  $C_{\text{div}}$ .

Going through the proof of Theorem 3, we observe that the constant obtained in (6) and (12) is just one of the many bounds one can achieve with T-coercivity for the Stokes problem. Indeed, one can choose any positive value of  $\lambda$ : hence, there exists a family of admissible operators  $\mathbb{T}$ , in the sense of Definition 1, which shows the flexibility of the approach.

Let us provide an illustration. For small viscosity  $\nu$  (the domain  $\Omega$  being fixed), it is well-known that the stability constant appearing in the estimate

$$\|(\mathbf{u}, p)\|_{\mathbf{H}_0^1(\Omega) \times L_0^2(\Omega)} \leq C(\nu) \|(\mathbf{f}, g)\|_{\mathbf{H}^{-1}(\Omega) \times L_0^2(\Omega)}$$

behaves like  $O(\nu^{-1})$ . For instance, for the velocity  $\mathbf{u}$ , the result is elementarily obtained by taking the test field  $(\mathbf{v}, q) = (\mathbf{u}, p)$  in (2). On the other hand, in (6) and (12), we find a behavior in  $O(\nu^{-2})$ . But, if one is interested in obtaining a less severe blowup, one can simply choose

$$\eta = \frac{\nu C_{\text{div}}^2}{2(1 - \frac{\nu}{2})} \quad \text{and} \quad \lambda = \frac{1}{2}(1 + \eta)$$

in the above proof, for all  $0 < \nu \leq 1$ . Then, one finds that

$$\underline{\alpha} = \frac{\nu}{2} \quad \text{and} \quad \|\mathbb{T}\| \leq \max\left(\frac{1}{\sqrt{2}}(1 + C_{\text{div}}^2), \left(2C_{\text{div}}^2 + \frac{1}{4}(1 + C_{\text{div}}^2)^2\right)^{1/2}\right),$$

so that (4) actually yields a stability constant in  $O(\nu^{-1})$ .

Theorem 3 provides a fully constructive proof for the well-posedness of Stokes problem, which is an emblematic example of mixed problem. In the next section, we show that the T-coercivity approach employed here is in fact very general and can be extended to a large class of saddle-point problems.

### 3. Abstract framework

We start with the classical statements regarding the definition of saddle-point problems, and the equivalent conditions to ensure an inf-sup condition. Then, we proceed with the design of abstract operators T to ensure well-posedness for saddle-problems, and for augmented saddle-point problems.

#### 3.1. Saddle-point problems in Hilbert spaces

Let  $V$  and  $Q$  be two Hilbert spaces. In the Hilbert space  $Q$ , we introduce the canonical isomorphism  $\mathbb{1}_{Q \rightarrow Q'} : Q \rightarrow Q'$  defined by

$$\langle \mathbb{1}_{Q \rightarrow Q'} p, q \rangle_{Q', Q} = (p, q)_Q, \quad \forall p \in Q, \forall q \in Q,$$

which is a bijective isometry according to Riesz Theorem. As a matter of fact, its inverse  $\mathbb{1}_{Q' \rightarrow Q}$  is also a bijective isometry, and

$$(\mathbb{1}_{Q' \rightarrow Q} g, q)_Q = \langle g, q \rangle_{Q', Q}, \quad \forall g \in Q', \forall q \in Q.$$

We then introduce two bilinear forms  $a(\cdot, \cdot)$  on  $V \times V$  and  $b(\cdot, \cdot)$  on  $V \times Q$  that are assumed to be continuous, *i.e.* there exist  $C_a > 0$  and  $C_b > 0$  such that

$$a(u, v) \leq C_a \|u\|_V \|v\|_V, \quad \forall u \in V, \forall v \in V, \quad (13)$$

$$b(v, q) \leq C_b \|v\|_V \|q\|_Q, \quad \forall v \in V, \forall q \in Q. \quad (14)$$

We denote by  $A$  and  $B$  the linear continuous operators associated with  $a$  and  $b$ , defined by

$$A \in \mathcal{L}(V, V'), \quad \langle Au, v \rangle_{V', V} = a(u, v), \quad \forall u \in V, \forall v \in V,$$

$$B \in \mathcal{L}(V, Q'), \quad \langle Bv, q \rangle_{Q', Q} = b(v, q), \quad \forall v \in V, \forall q \in Q.$$

The adjoint operator of  $B$  is given by

$$B^* \in \mathcal{L}(Q, V'), \quad \langle B^* q, v \rangle_{V', V} = \langle Bv, q \rangle_{Q', Q} = b(v, q), \quad \forall v \in V, \forall q \in Q.$$

Given  $f \in V'$  and  $g \in Q'$ , we consider the saddle-point problem: find  $(u, p) \in V \times Q$  such that

$$\begin{aligned} Au + B^* p &= f, & \text{in } V', \\ Bu &= g, & \text{in } Q'. \end{aligned} \quad (15)$$

Or, equivalently, in variational form:

$$\left\{ \begin{array}{l} \text{Find } (u, p) \in V \times Q \text{ such that} \\ \forall v \in V, \quad a(u, v) + b(v, p) = \langle f, v \rangle_{V', V}, \\ \forall q \in Q, \quad b(u, q) = \langle g, q \rangle_{Q', Q}. \end{array} \right. \quad (16)$$



As for the Stokes problem, we write Problem (15) as an *all-in-one* variational formulation

$$\begin{cases} \text{Find } (u, p) \in V \times Q & \text{such that} \\ \forall (v, q) \in V \times Q, & \mathcal{A}((u, p), (v, q)) = \langle f, v \rangle_{V', V} + \langle g, q \rangle_{Q', Q}, \end{cases} \quad (17)$$

where

$$\mathcal{A}((u, p), (v, q)) = a(u, v) + b(v, p) + b(u, q).$$

In what follows, we show that Problem (17) is well-posed using the notion of T-coercivity, with slightly different techniques depending on the assumptions made on the bilinear form  $a$ .

Regarding the form  $b(\cdot, \cdot)$  and the operator  $B$ , one has the well-known result below, see for instance [27, Lemma I.4.1]<sup>1</sup>, which can be viewed as a reformulation of Banach's Closed Range Theorem.

**Theorem 5.** *The following three statements are equivalent:*

(i) *There exists  $\underline{\beta} > 0$  such that*

$$\inf_{q \in Q \setminus \{0\}} \sup_{v \in V \setminus \{0\}} \frac{b(v, q)}{\|v\|_V \|q\|_Q} \geq \underline{\beta}. \quad (18)$$

(ii)  *$B : (\text{Ker } B)^\perp \rightarrow Q'$  is an isomorphism, and*

$$\|Bv\|_{Q'} \geq \underline{\beta} \|v\|_V, \quad \forall v \in (\text{Ker } B)^\perp.$$

(iii) *There exists an isomorphic operator  $L_B : Q' \rightarrow (\text{Ker } B)^\perp$  such that*

$$B(L_B g) = g \quad \text{and} \quad \|g\|_{Q'} \geq \underline{\beta} \|L_B g\|_V, \quad \forall g \in Q'.$$

Since our aim is to build operators T from  $V \times Q$  to itself, we first introduce the operator

$$B = \mathbb{1}_{Q' \rightarrow Q} \circ B : V \rightarrow Q.$$

For all  $v \in V$ ,  $\|Bv\|_{Q'} = \|\mathbb{1}_{Q' \rightarrow Q}(Bv)\|_Q = \|Bv\|_Q$  and, for all  $(v, q) \in V \times Q$ ,

$$b(v, q) = \langle Bv, q \rangle_{Q', Q} = \langle \mathbb{1}_{Q \rightarrow Q'}(Bv), q \rangle_{Q', Q} = (Bv, q)_Q.$$

Whenever applicable, we also introduce its *right-inverse*

$$L_B = L_B \circ \mathbb{1}_{Q \rightarrow Q'} : Q \rightarrow (\text{Ker } B)^\perp.$$

Observe that

$$b(L_B p, q) = \langle BL_B p, q \rangle_{Q', Q} = \langle \mathbb{1}_{Q \rightarrow Q'} p, q \rangle_{Q', Q} = (p, q)_Q, \quad \forall p \in Q, \forall q \in Q. \quad (19)$$

Under these notations, items (ii)-(iii) of Theorem 5 now write

(ii)  $B : (\text{Ker } B)^\perp \rightarrow Q$  is an *isomorphism*, and

$$\|Bv\|_Q \geq \underline{\beta} \|v\|_V, \quad \forall v \in (\text{Ker } B)^\perp. \quad (20)$$

(iii) There exists an isomorphic operator  $L_B : Q \rightarrow (\text{Ker } B)^\perp$  such that

$$B(L_B q) = q \quad \text{and} \quad \|q\|_Q \geq \underline{\beta} \|L_B q\|_V, \quad \forall q \in Q. \quad (21)$$

For convenience, we often use  $\beta = \underline{\beta}^{-1}$ , so that

$$\|L_B q\|_V \leq \beta \|q\|_Q, \quad \forall q \in Q.$$

<sup>1</sup>Item (iii) below is a rephrasing of the original statement, because it is better suited for our purposes. For details, see the proof of Lemma I.4.1. of [27] p. 59, item 2°. The operator  $L_B$  is a *right-inverse* of the operator  $B$ .

### 3.2. How to achieve T-coercivity for saddle-point problems?

If  $a$  is coercive on the whole space  $V$ , we can extend the proof of Theorem 3 in the following way.

**Theorem 6.** *Assume that (18) holds true and that the form  $a$  is symmetric and positive. If there exists a constant  $\alpha > 0$  such that*

$$a(u, u) \geq \alpha \|u\|_V^2, \quad \forall u \in V, \quad (22)$$

then there exists a unique solution to Problem (17) and

$$\|(u, p)\|_{V \times Q} \leq \frac{2 \max\left(\sqrt{2}C_a\beta^2, \beta(2 + C_a^2\beta^2)^{1/2}\right)}{\min(\alpha C_a\beta^2, 1)} \|(f, g)\|_{V' \times Q'}. \quad (23)$$

**Proof.** First, we note that the symmetry of the bilinear form  $a$  implies that  $\mathcal{A}$  is also symmetric. Then, we follow the same ideas as in the proof of Theorem 3, replacing  $\mathbf{v}_p$  by  $L_B p$ . We introduce the mapping

$$\begin{aligned} \mathbf{T} : V \times Q &\longrightarrow V \times Q \\ (u, p) &\longmapsto (\lambda u + L_B p, -\lambda p) \end{aligned}$$

and we compute

$$\begin{aligned} \mathcal{A}((u, p), \mathbf{T}(u, p)) &= a(u, \lambda u) + a(u, L_B p) + b(\lambda u, p) + b(L_B p, p) - b(u, \lambda p) \\ &= \lambda a(u, u) + a(u, L_B p) + \|p\|_Q^2, \end{aligned}$$

in view of (19).

Because the form  $a$  is symmetric and positive, we can apply Young's inequality: for any  $\eta > 0$ ,

$$a(u, L_B p) \geq -\frac{\eta}{2} a(u, u) - \frac{1}{2\eta} a(L_B p, L_B p).$$

Taking into account (13) and (21), the latter being equivalent to (18), we get

$$a(L_B p, L_B p) \leq C_a \|L_B p\|_V^2 \leq C_a \beta^2 \|p\|_Q^2$$

and thus

$$a(u, L_B p) \geq -\frac{\eta}{2} a(u, u) - \frac{C_a \beta^2}{2\eta} \|p\|_Q^2.$$

Hence, recalling (22), if  $\lambda - \frac{\eta}{2} > 0$  it follows that

$$\mathcal{A}((u, p), \mathbf{T}(u, p)) \geq \alpha \left(\lambda - \frac{\eta}{2}\right) \|u\|_V^2 + \left(1 - \frac{C_a \beta^2}{2\eta}\right) \|p\|_Q^2.$$

Setting in particular  $\eta = \lambda = C_a \beta^2$ , we infer that

$$\mathcal{A}((u, p), \mathbf{T}(u, p)) \geq \alpha \frac{C_a \beta^2}{2} \|u\|_V^2 + \frac{1}{2} \|p\|_Q^2 \geq \frac{1}{2} \min(\alpha C_a \beta^2, 1) \|(u, p)\|_{V \times Q}^2, \quad (24)$$

which proves that  $\mathcal{A}$  is T-coercive.

Since  $\mathbf{T}(u, p) = (C_a \beta^2 u + L_B p, -C_a \beta^2 p)$ , it holds that

$$\begin{aligned} \|\mathbf{T}(u, p)\|_{V \times Q}^2 &= \|C_a \beta^2 u + L_B p\|_V^2 + \|C_a \beta^2 p\|_Q^2 \\ &\leq 2(C_a \beta^2)^2 \|u\|_V^2 + 2\|L_B p\|_V^2 + (C_a \beta^2)^2 \|p\|_Q^2 \\ &\leq 2(C_a \beta^2)^2 \|u\|_V^2 + (2\beta^2 + (C_a \beta^2)^2) \|p\|_Q^2, \end{aligned}$$

which yields

$$\|\mathbf{T}\| \leq \max\left(\sqrt{2}C_a\beta^2, \beta(2 + C_a^2\beta^2)^{1/2}\right). \quad (25)$$

Lastly, we observe that

$$\langle f, v \rangle_{V', V} + \langle g, q \rangle_{Q', Q} \leq \|(f, g)\|_{V' \times Q'} \|(v, q)\|_{V \times Q}. \quad (26)$$

Combining (24), (25) and (26), the stability estimate (4) furnishes exactly (23).  $\square$

**Remark 7.** By applying Theorem 6 to Stokes problem, we recover stability estimates (6) and (12) from the correspondence  $\alpha = \nu$ ,  $C_a = \nu$  and  $\beta = C_{\text{div}}$ .

In Ladyzhenskaya–Babuška–Brezzi’s theory and in many applications, the bilinear form  $a$  is not coercive on the whole space  $V$  but only on the kernel of the operator  $B$ . This is for instance the case in electromagnetism, which will be detailed in Section 4. The next result shows how to address this situation in the T-coercivity framework (provided that the form  $a$  is symmetric and positive), thus establishing the equivalence between the two theories.

**Theorem 8.** *Assume that the form  $a$  is symmetric and positive.*

1. *If (18) holds true, and if there exists a constant  $\alpha_0 > 0$  such that*

$$a(u_0, u_0) \geq \alpha_0 \|u_0\|_V^2, \quad \forall u_0 \in \text{Ker } B, \quad (27)$$

*then the form  $\mathcal{A}$  is T-coercive. In other words, Problem (17) is well-posed and*

$$\|(u, p)\|_{V \times Q} \leq C \|(f, g)\|_{V' \times Q'}, \quad (28)$$

*with  $C$  a constant depending only on  $\alpha_0$ ,  $\beta$ ,  $C_a$  and  $C_b$ .*

2. *Conversely, if Problem (17) is well-posed, that is, if the form  $\mathcal{A}$  is T-coercive, then (18) and (27) both hold.*

**Proof.** 1. We consider the mapping

$$\begin{aligned} \mathbb{T}: V \times Q &\longrightarrow V \times Q \\ (u, p) &\longmapsto (\lambda u + L_B p, -\lambda p + \lambda \mu B u). \end{aligned}$$

This is almost the same mapping as the one used in the proof of Theorem 6. The only difference is the term  $\lambda \mu B u$ , which is going to help us handling the extra terms that do not belong to the kernel of  $B$  by adjusting the value of the constant  $\mu$ . We get

$$\begin{aligned} \mathcal{A}((u, p), \mathbb{T}(u, p)) &= a(u, \lambda u) + a(u, L_B p) + b(\lambda u, p) + b(L_B p, p) - b(u, \lambda p) + b(u, \lambda \mu B u) \\ &= \lambda a(u, u) + a(u, L_B p) + \|p\|_Q^2 + \lambda \mu \|B u\|_Q^2 \end{aligned}$$

because  $b(L_B p, p) = \|p\|_Q^2$  as previously, and

$$b(u, B u) = \langle B u, B u \rangle_{Q', Q} = (\mathbb{1}_{Q' \rightarrow Q}(B u), B u)_Q = \|B u\|_Q^2.$$

Since the form  $a$  is symmetric and positive, one may use Young’s inequality. By proceeding as in the proof of Theorem 6 and after setting  $\lambda = C_a \beta^2$ , we know that

$$\lambda a(u, u) + a(u, L_B p) + \|p\|_Q^2 \geq \frac{C_a \beta^2}{2} a(u, u) + \frac{1}{2} \|p\|_Q^2,$$

from which we deduce

$$\mathcal{A}((u, p), \mathbb{T}(u, p)) \geq \frac{C_a \beta^2}{2} (a(u, u) + 2\mu \|B u\|_Q^2) + \frac{1}{2} \|p\|_Q^2.$$

To compensate the lack of coercivity of  $a$  outside  $\text{Ker } B$ , we use the decomposition  $u = u_0 + \bar{u}$  with  $u_0 \in \text{Ker } B$  and  $\bar{u} \in (\text{Ker } B)^\perp$ . Following [4, p. 254], Young’s inequality yields

$$\begin{aligned} a(u, u) &= a(u_0, u_0) + 2a(u_0, \bar{u}) + a(\bar{u}, \bar{u}) \\ &\geq (1 - \theta) a(u_0, u_0) + \left(1 - \frac{1}{\theta}\right) a(\bar{u}, \bar{u}) \\ &\geq (1 - \theta) a(u_0, u_0) + \left(C_a - \frac{C_a}{\theta}\right) \|\bar{u}\|_V^2 \end{aligned}$$

for all  $0 < \theta < 1$ . Since  $u_0 \in \text{Ker} B$ , we have  $\|Bu\|_Q^2 = \|B\bar{u}\|_Q^2$ . Moreover, using (20) yields  $\|B\bar{u}\|_Q^2 \geq \beta^{-2} \|\bar{u}\|_V^2$ . Thus

$$a(u, u) + 2\mu \|Bu\|_Q^2 \geq (1 - \theta)a(u_0, u_0) + \left(C_a - \frac{C_a}{\theta} + \frac{2\mu}{\beta^2}\right) \|\bar{u}\|_V^2. \quad (29)$$

Choosing  $\theta = \frac{1}{2}$  and  $\mu = \frac{3}{4}C_a\beta^2$ , it holds that

$$a(u, u) + 2\mu \|Bu\|_Q^2 \geq \frac{1}{2}a(u_0, u_0) + \frac{C_a}{2} \|\bar{u}\|_V^2.$$

Hence, recalling (27) and using the inequality  $C_a \geq \alpha_0$ , we obtain

$$a(u, u) + 2\mu \|Bu\|_Q^2 \geq \frac{\alpha_0}{2} \|u_0\|_V^2 + \frac{\alpha_0}{2} \|\bar{u}\|_V^2 = \frac{\alpha_0}{2} \|u\|_V^2$$

and we conclude that

$$\mathcal{A}((u, p), T(u, p)) \geq \alpha_0 \frac{C_a\beta^2}{4} \|u\|_V^2 + \frac{1}{2} \|p\|_Q^2. \quad (30)$$

From the above, we have

$$T(u, p) = (C_a\beta^2 u + L_B p, -C_a\beta^2 p + \frac{3}{4}(C_a\beta^2)^2 Bu).$$

Finally,  $T$  belongs to  $\mathcal{L}(V \times Q)$  since  $\|L_B p\|_V \leq \beta \|p\|_Q$  (see (21)) and<sup>2</sup>

$$\|Bu\|_Q \leq C_b \|u\|_V. \quad (31)$$

The stability estimate (28) is then given by (4).

2. Conversely, suppose that there exist  $\alpha_V > 0$ ,  $\alpha_Q > 0$  and  $T \in \mathcal{L}(V \times Q)$  such that

$$\mathcal{A}((u, p), T(u, p)) \geq \alpha_V \|u\|_V^2 + \alpha_Q \|p\|_Q^2, \quad \forall (u, p) \in V \times Q. \quad (32)$$

Noting  $T : (u, p) \mapsto (T_V(u, p), T_Q(u, p))$ , we have

$$\mathcal{A}((u, p), T(u, p)) = a(u, T_V(u, p)) + b(T_V(u, p), p) + b(u, T_Q(u, p))$$

and, since  $T$  is bounded,

$$\|T_V(u, p)\|_V^2 + \|T_Q(u, p)\|_Q^2 \leq \|T\|^2 (\|u\|_V^2 + \|p\|_Q^2). \quad (33)$$

Now, choosing  $u = 0$  in (32) and (33) yields

$$b(T_V(0, p), p) \geq \alpha_Q \|p\|_Q^2 \quad \text{and} \quad \|T_V(0, p)\|_V \leq \|T\| \cdot \|p\|_Q, \quad \forall p \in Q.$$

Thus, for  $p \in Q \setminus \{0\}$ ,  $T_V(0, p) \neq 0$ , otherwise  $b(T_V(0, p), p) = 0$ , which contradicts  $b(T_V(0, p), p) > 0$ . Then it follows that

$$\sup_{v \in V \setminus \{0\}} \frac{b(v, p)}{\|v\|_V} \geq \frac{b(T_V(0, p), p)}{\|T_V(0, p)\|_V} \geq \frac{\alpha_Q}{\|T\|} \|p\|_Q, \quad \forall p \in Q \setminus \{0\},$$

which shows that the inf-sup condition (18) is fulfilled. Likewise, taking  $p = 0$  and  $u \in \text{Ker} B$  in (32) and (33), we get

$$a(u, T_V(u, 0)) \geq \alpha_V \|u\|_V^2 \quad \text{and} \quad \|T_V(u, 0)\|_V \leq \|T\| \|u\|_V, \quad \forall u \in \text{Ker} B.$$

By symmetry and positivity of  $a$ , it holds that

$$a(u, T_V(u, 0)) \leq (a(u, u))^{1/2} a(T_V(u, 0), T_V(u, 0))^{1/2}.$$

<sup>2</sup>Classically,

$$\begin{aligned} \|Bu\|_Q^2 &= (Bu, Bu)_Q = \langle \mathbb{1}_{Q \rightarrow Q'}(Bu), Bu \rangle_{Q', Q} \quad \text{by definition of } \mathbb{1}_{Q \rightarrow Q'}, \\ &= \langle \mathbb{1}_{Q \rightarrow Q'} \circ \mathbb{1}_{Q' \rightarrow Q}(Bu), Bu \rangle_{Q', Q} = \langle Bu, Bu \rangle_{Q', Q} \quad \text{since } \mathbb{1}_{Q \rightarrow Q'} \circ \mathbb{1}_{Q' \rightarrow Q} = \text{Id}_{Q'}, \\ &= b(u, Bu) \leq C_b \|u\|_V \|Bu\|_Q \quad \text{by definition and continuity of } b \text{ (14)}. \end{aligned}$$

Thus

$$\alpha_V \|u\|_V^2 \leq a(u, T_V(u, 0)) \leq (a(u, u))^{1/2} (C_a \|T\|^2 \|u\|_V^2)^{1/2}$$

and hence  $a(u, u) \geq \frac{\alpha_V^2}{C_a \|T\|^2} \|u\|_V^2$  for all  $u \in \text{Ker} B$ , which proves (27).  $\square$

**Remark 9.** The T-coercivity estimate (30) is very close to the case where  $a$  is coercive on the whole space  $V$ . As a matter of fact, the only difference compared to (24) is that the constant before the term  $\|u\|_V^2$  is twice as small, with  $\alpha_0 = \alpha$ .

### 3.3. Augmented saddle-point problems

Let  $c(\cdot, \cdot)$  be a positive and continuous bilinear form defined on  $Q \times Q$ , namely

$$c(p, p) \geq 0, \quad \forall p \in Q \quad \text{and} \quad \exists C_c > 0, \quad c(p, q) \leq C_c \|p\|_Q \|q\|_Q, \quad \forall p \in Q, \forall q \in Q. \quad (34)$$

We denote by  $C$  the linear operator associated with the bilinear form  $c$ , defined by

$$C \in \mathcal{L}(Q, Q'), \quad \langle Cp, q \rangle_{Q', Q} = c(p, q), \quad \forall p \in Q, \forall q \in Q.$$

The *all-in-one* approach developed previously also enables us to deal with augmented saddle-point problems: given  $f \in V'$  and  $g \in Q'$ , find  $(u, p) \in V \times Q$  such that

$$\begin{aligned} Au + B^* p &= f, & \text{in } V', \\ Bu - Cp &= g, & \text{in } Q', \end{aligned} \quad (35)$$

where the operator  $C$  possibly acts as a *small perturbation* of the original saddle-point problem (15). The weak formulation of (35) reads:

$$\left\{ \begin{array}{l} \text{Find } (u, p) \in V \times Q \quad \text{such that} \\ \forall (v, q) \in V \times Q, \quad \mathcal{A}_c((u, p), (v, q)) = \langle f, v \rangle_{V', V} + \langle g, q \rangle_{Q', Q}, \end{array} \right. \quad (36)$$

with

$$\mathcal{A}_c((u, p), (v, q)) = a(u, v) + b(v, p) + b(u, q) - c(p, q).$$

As before, the bilinear form  $a$  is supposed to be symmetric and positive.

### 3.4. How to achieve T-coercivity for augmented saddle-point problems?

Once again, we distinguish the case where the form  $a$  is coercive on  $V$  or only on  $\text{Ker} B$ . If the form  $a$  is coercive on  $V$ , the results from the un-augmented case allow straightforwardly to handle the augmented one.

**Theorem 10.** *Assume that (18) holds true and that the form  $a$  is symmetric and positive. If there exists a constant  $\alpha > 0$  such that*

$$a(u, u) \geq \alpha \|u\|_V^2, \quad \forall u \in V,$$

*then there exists a unique solution to Problem (36).*

**Proof.** With the same operator  $T$  as for the un-augmented problem, namely

$$\begin{aligned} T: V \times Q &\longrightarrow V \times Q \\ (u, p) &\longmapsto (C_a \beta^2 u + L_B p, -C_a \beta^2 p), \end{aligned}$$

it holds that

$$\mathcal{A}_c((u, p), T(u, p)) = C_a \beta^2 a(u, u) + a(u, L_B p) + \|p\|_Q^2 + C_a \beta^2 c(p, p).$$

Therefore, a similar argument as in Theorem 6 furnishes

$$\mathcal{A}_c((u, p), T(u, p)) \geq \alpha \frac{C_a \beta^2}{2} \|u\|_V^2 + \frac{1}{2} \|p\|_Q^2 + C_a \beta^2 c(p, p),$$

which shows that  $\mathcal{A}_c$  is T-coercive since  $c$  is positive.  $\square$

**Remark 11.** A particular case that appears in many applications – see Section 5 for the example of nearly-incompressible elasticity – is when  $c$  has the form

$$c(p, q) = \varepsilon(p, q)_Q, \quad \varepsilon \geq 0.$$

In this case, we obtain the estimate

$$\mathcal{A}_c((u, p), T(u, p)) \geq \alpha \frac{C_a \beta^2}{2} \|u\|_V^2 + \left( \frac{1}{2} + \varepsilon C_a \beta^2 \right) \|p\|_Q^2, \quad (37)$$

so that the augmentation  $c$  improves the constant before the term  $\|p\|_Q^2$  and thus stabilizes the bilinear form  $\mathcal{A}_c$ . Moreover, the above estimate is robust for small values of  $\varepsilon$ . Besides, it even allows to take negative values of  $\varepsilon$ . Indeed, if  $\varepsilon < 0$ , we have

$$\mathcal{A}_c((u, p), T(u, p)) \geq \alpha \frac{C_a \beta^2}{2} \|u\|_V^2 + \left( \frac{1}{2} - |\varepsilon| C_a \beta^2 \right) \|p\|_Q^2.$$

Hence, the bilinear form  $\mathcal{A}_c$  remains T-coercive whenever  $|\varepsilon| < \frac{1}{2C_a \beta^2}$ .

Let us now suppose that  $a$  is not coercive on the whole space  $V$  but only on the kernel of  $B$ . Then, two different situations occur. Either the form  $c$  can be viewed as a *small perturbation*, and we shall look for a solution of (35) that is *close* to the solution of the original problem (15). Or this is not the case, and the form  $c$  is viewed as a “fixed” augmentation, and there is no obvious connection *a priori* between the solutions of the augmented and un-augmented problems.

### 3.5. Additional results for small perturbations

We say that  $c$  is a small perturbation if it can be written as

$$c(p, q) = \varepsilon c_0(p, q), \quad \varepsilon > 0, \quad (38)$$

with  $\varepsilon$  a small parameter and  $c_0$  a symmetric, positive and continuous form on  $Q$ . We start with the simple case

$$c(p, q) = \varepsilon(p, q)_Q, \quad \varepsilon > 0, \quad (39)$$

for which the T-coercivity approach yields a shorter proof than the corresponding result stated in Ladyzhenskaya–Babuška–Brezzi’s framework, see [4, pages 247-252].

**Theorem 12.** *Assume that (18) holds true, that the form  $a$  is symmetric and positive, and that  $c$  takes the simple form of (39). If there exists a constant  $\alpha_0 > 0$  such that*

$$a(u_0, u_0) \geq \alpha_0 \|u_0\|_V^2, \quad \forall u_0 \in \text{Ker } B, \quad (40)$$

and if  $\varepsilon$  is small enough, namely

$$\varepsilon \leq \frac{1}{2C_a \beta^4 C_b^2} \left( 2 - \frac{\alpha_0}{C_a} \right), \quad (41)$$

then Problem (36) is well-posed and

$$\|(u, p)\|_{V \times Q} \leq C \|(f, g)\|_{V' \times Q'}, \quad (42)$$

with  $C$  a constant depending only on  $\alpha_0$ ,  $\beta$ ,  $C_a$  and  $C_b$ .

**Proof.** Here again, we consider the mapping

$$\begin{aligned} T: V \times Q &\longrightarrow V \times Q \\ (u, p) &\longmapsto (\lambda u + L_B p, -\lambda p + \lambda \mu B u). \end{aligned}$$

The beginning of the proof is the same as in Theorem 8. Taking into account the extra terms coming from the perturbation, we get

$$\mathcal{A}_c((u, p), T(u, p)) = \lambda a(u, u) + a(u, L_B p) + \|p\|_Q^2 + \lambda \mu \|B u\|_{Q'}^2 + \lambda c(p, p) - \lambda \mu c(p, B u).$$

Using Young's inequality and setting  $\lambda = C_a \beta^2$ , it follows that

$$\mathcal{A}_c((u, p), T(u, p)) \geq \frac{C_a \beta^2}{2} (a(u, u) + 2\mu \|B u\|_{Q'}^2) + \frac{1}{2} \|p\|_Q^2 + C_a \beta^2 c(p, p) - C_a \beta^2 \mu c(p, B u). \quad (43)$$

Now, as in (29), it holds that

$$a(u, u) + 2\mu \|B u\|_{Q'}^2 \geq (1 - \theta) a(u_0, u_0) + \left( C_a - \frac{C_a}{\theta} + \frac{2\mu}{\beta^2} \right) \|\bar{u}\|_V^2 \quad (44)$$

for all  $0 < \theta < 1$ , where  $u = u_0 + \bar{u}$  with  $u_0 \in \text{Ker} B$  and  $\bar{u} \in (\text{Ker} B)^\perp$ .

Knowing that  $c(p, q) = \varepsilon(p, q)_Q$  for all  $p$  and  $q$  in  $Q$ , Young's inequality implies that, for all  $\delta > 0$ ,

$$\begin{aligned} -c(p, B u) &= -c(p, B \bar{u}) = -\varepsilon(p, B \bar{u})_Q \geq -\varepsilon \frac{\delta}{2} \|p\|_Q^2 - \frac{\varepsilon}{2\delta} \|B \bar{u}\|_Q^2 \\ &\geq -\varepsilon \frac{\delta}{2} \|p\|_Q^2 - \varepsilon \frac{C_b^2}{2\delta} \|\bar{u}\|_V^2 \quad \text{in view of (31)}. \end{aligned}$$

Putting (43), (44) and the above inequality together, we find that

$$\begin{aligned} \mathcal{A}_c((u, p), T(u, p)) &\geq \frac{C_a \beta^2}{2} \left( (1 - \theta) a(u_0, u_0) + \left( C_a - \frac{C_a}{\theta} + \frac{2\mu}{\beta^2} - \mu \varepsilon \frac{C_b^2}{\delta} \right) \|\bar{u}\|_V^2 \right) \\ &\quad + \frac{1}{2} \|p\|_Q^2 + \varepsilon C_a \beta^2 \left( 1 - \mu \frac{\delta}{2} \right) \|p\|_Q^2. \end{aligned}$$

Hence, choosing  $\theta = \frac{1}{2}$ ,  $\mu = C_a \beta^2$  and recalling (40), it holds that

$$\mathcal{A}_c((u, p), T(u, p)) \geq \frac{C_a \beta^2}{2} \left( \frac{\alpha_0}{2} \|u_0\|_V^2 + C_a \left( 1 - \varepsilon \frac{\beta^2 C_b^2}{\delta} \right) \|\bar{u}\|_V^2 \right) + \frac{1}{2} \|p\|_Q^2 + \varepsilon C_a \beta^2 \left( 1 - C_a \beta^2 \frac{\delta}{2} \right) \|p\|_Q^2. \quad (45)$$

Finally, we set  $\delta = \frac{1}{C_a \beta^2}$  so that

$$1 - C_a \beta^2 \frac{\delta}{2} = \frac{1}{2} \quad \text{and} \quad 1 - \varepsilon \frac{\beta^2 C_b^2}{\delta} = 1 - \varepsilon C_a \beta^4 C_b^2 \geq \frac{1}{2} \cdot \frac{\alpha_0}{C_a}$$

in virtue of (41). Thus

$$\mathcal{A}_c((u, p), T(u, p)) \geq \alpha_0 \frac{C_a \beta^2}{4} \|u\|_V^2 + \left( \frac{1}{2} + \varepsilon \frac{C_a \beta^2}{2} \right) \|p\|_Q^2, \quad (46)$$

where we used that  $\|u\|_V^2 = \|u_0\|_V^2 + \|\bar{u}\|_V^2$ . All in all, we have chosen

$$T(u, p) = (C_a \beta^2 u + L_B p, -C_a \beta^2 p + (C_a \beta^2)^2 B u).$$

Then, estimate (42) follows from (4) with a stability constant independent of  $\varepsilon$  since (46) is robust for vanishing  $\varepsilon$  and since  $\|T\|$  does not depend on  $\varepsilon$  either.  $\square$

**Remark 13.** The final estimate (46) is very close to (37). The only difference between these two estimates is a factor of 2 between the constants multiplying the norms of  $u$  and  $p$ , with  $\alpha_0 = \alpha$ .

**Remark 14.** In Ladyzhenskaya–Babuška–Brezzi's framework, it is commonly assumed that  $\varepsilon \leq 1$ . On the other hand, in (41), we find a smallness condition that depends *explicitly* on the various constants of the problem.

**Remark 15.** The inf-sup condition (18) and the continuity of  $b$  imply that  $\underline{\beta} \leq C_b$ , i.e.  $C_b \beta \geq 1$ . Therefore, (41) yields in particular

$$\varepsilon \leq \frac{1}{C_a \beta^2},$$

which corresponds to the condition found in Remark 11 for negative values of  $\varepsilon$ . As a matter of fact, the non-coercivity of  $a$  on the whole space  $V$  calls for the introduction of a term  $Bu$  in the mapping  $T$ . This term induces an additional term of the form  $c(p, Bu)$  in the expression of  $\mathcal{A}_c((u, p), T(u, p))$ , that can be interpreted as a “negative perturbation” of the bilinear form  $\mathcal{A}$ .

Now, we move to the case where  $c$  is given by (38). Let us denote by  $C_{c_0}$  the continuity constant of the bilinear form  $c_0$ . The next theorem establishes the well-posedness of the perturbed problem for a very general form  $c_0$ .

**Theorem 16.** *Assume that (18) holds true, and that the bilinear forms  $a$  and  $c_0$  are both symmetric and positive. Suppose in addition that there exist  $\alpha_0 > 0$  and  $\gamma_0 > 0$  such that*

$$a(u_0, u_0) \geq \alpha_0 \|u_0\|_V^2, \quad \forall u_0 \in \text{Ker } B,$$

and

$$c_0(p_0, p_0) \geq \gamma_0 \|p_0\|_Q^2, \quad \forall p_0 \in \text{Ker } B^*. \quad (47)$$

If  $\varepsilon$  is small enough, namely

$$\varepsilon \leq \frac{1}{2C_{c_0} C_a \beta^4 C_b^2}, \quad (48)$$

then Problem (36) is well-posed and

$$\|(u, p)\|_{V \times Q} \leq C \|(f, g)\|_{V' \times Q'},$$

with  $C$  a constant depending only on  $\alpha_0$ ,  $\beta$ ,  $\gamma_0$ ,  $C_a$  and  $C_b$ .

**Proof.** First, we adapt the beginning of the proof of Theorem 12 to take into consideration the bilinear form  $c_0$ . Since  $c_0$  is symmetric and positive, we can use Young’s inequality to obtain

$$\begin{aligned} -c(p, Bu) &= -\varepsilon c_0(p, B\bar{u})_Q \geq -\varepsilon \frac{\delta}{2} c_0(p, p) - \frac{\varepsilon}{2\delta} c_0(B\bar{u}, B\bar{u}) \\ &\geq -\varepsilon \frac{\delta}{2} c_0(p, p) - \varepsilon \frac{C_{c_0} C_b^2}{2\delta} \|\bar{u}\|_V^2 \quad \text{since } \|B\bar{u}\|_Q^2 \leq C_b^2 \|\bar{u}\|_V^2, \end{aligned}$$

and thus (45) becomes

$$\begin{aligned} \mathcal{A}_c((u, p), T(u, p)) &\geq \frac{C_a \beta^2}{2} \left( \frac{\alpha_0}{2} \|u_0\|_V^2 + C_a \left( 1 - \varepsilon \frac{C_{c_0} \beta^2 C_b^2}{\delta} \right) \|\bar{u}\|_V^2 \right) \\ &\quad + \frac{1}{2} \|p\|_Q^2 + \varepsilon C_a \beta^2 \left( 1 - C_a \beta^2 \frac{\delta}{2} \right) c_0(p, p), \end{aligned}$$

where  $T$  is the mapping

$$\begin{aligned} T: V \times Q &\longrightarrow V \times Q \\ (u, p) &\longmapsto (C_a \beta^2 u + L_B p, -C_a \beta^2 p + (C_a \beta^2)^2 B u). \end{aligned}$$

Setting  $\delta = \frac{1}{C_a \beta^2}$  as before, we get the estimate

$$\mathcal{A}_c((u, p), T(u, p)) \geq \alpha_0 \frac{C_a \beta^2}{4} \|u\|_V^2 + \frac{1}{2} \|p\|_Q^2 + \varepsilon \frac{C_a \beta^2}{2} c_0(p, p),$$

as long as  $\varepsilon \leq \frac{1}{2C_{c_0} C_a \beta^4 C_b^2} (2 - \frac{\alpha_0}{C_a})$ , which is the case under the assumption (48) since  $\alpha_0 \leq C_a$ .

Then, as the bilinear form  $c_0$  is not necessarily coercive on the whole space  $Q$ , we use the decomposition  $p = p_0 + \bar{p}$  with  $p_0 \in \text{Ker } B^*$  and  $\bar{p} \in (\text{Ker } B^*)^\perp$ . From Young’s inequality, we have

$$c_0(p, p) \geq (1 - \theta) c_0(p_0, p_0) + \left( C_{c_0} - \frac{C_{c_0}}{\theta} \right) \|\bar{p}\|_Q^2$$



for all  $0 < \theta < 1$ . Setting  $\theta = \frac{1}{2}$  and using (47), it follows that

$$\mathcal{A}_c((u, p), T(u, p)) \geq \alpha_0 \frac{C_a \beta^2}{4} \|u\|_V^2 + \frac{1}{2} \|p\|_Q^2 + \varepsilon \frac{C_a \beta^2}{2} \left( \frac{\gamma_0}{2} \|p_0\|_Q^2 - C_{c_0} \|\bar{p}\|_Q^2 \right).$$

Now, we notice that

$$\frac{1}{2} \|p\|_Q^2 \geq \frac{1}{8} \|p\|_Q^2 + \frac{3}{8} \|\bar{p}\|_Q^2$$

and that, thanks to (48),

$$\frac{3}{8} \|\bar{p}\|_Q^2 = \frac{1}{2} \cdot \frac{3}{4} \|\bar{p}\|_Q^2 \geq \varepsilon C_{c_0} C_a \beta^4 C_b^2 \cdot \frac{3}{4} \|\bar{p}\|_Q^2 \geq \varepsilon \frac{C_a \beta^2}{2} \cdot \frac{3}{2} C_{c_0} \|\bar{p}\|_Q^2$$

because  $C_b \beta \geq 1$ . Hence

$$\begin{aligned} \mathcal{A}_c((u, p), T(u, p)) &\geq \alpha_0 \frac{C_a \beta^2}{4} \|u\|_V^2 + \frac{1}{8} \|p\|_Q^2 + \varepsilon \frac{C_a \beta^2}{2} \left( \frac{\gamma_0}{2} \|p_0\|_Q^2 + \left( \frac{3}{2} C_{c_0} - C_{c_0} \right) \|\bar{p}\|_Q^2 \right) \\ &\geq \alpha_0 \frac{C_a \beta^2}{4} \|u\|_V^2 + \left( \frac{1}{8} + \varepsilon \gamma_0 \frac{C_a \beta^2}{4} \right) \|p\|_Q^2 \quad \text{since } C_{c_0} \geq \gamma_0, \end{aligned}$$

which shows that  $\mathcal{A}_c$  is T-coercive.  $\square$

Lastly, we mention that an important consequence of the previous result is to estimate the distance between the solution  $(u_\varepsilon, p_\varepsilon)$  of the perturbed problem

$$\begin{aligned} Au_\varepsilon + B^* p_\varepsilon &= f, \quad \text{in } V', \\ Bu_\varepsilon - \varepsilon C_0 p_\varepsilon &= g, \quad \text{in } Q', \end{aligned} \tag{49}$$

and the solution  $(u, p)$  of the original saddle-point problem (15) as a function of the penalty parameter  $\varepsilon$ .

**Corollary 17.** *Assume that (18) holds true, that the form  $a$  is symmetric and positive, and that  $c$  takes the form of (38). If there exist  $\alpha_0 > 0$  and  $\gamma_0 > 0$  such that*

$$a(u_0, u_0) \geq \alpha_0 \|u_0\|_V^2, \quad \forall u_0 \in \text{Ker} B, \quad c_0(p_0, p_0) \geq \gamma_0 \|p_0\|_Q^2, \quad \forall p_0 \in \text{Ker} B^*,$$

and if

$$\varepsilon \leq \frac{1}{2C_{c_0} C_a \beta^4 C_b^2},$$

then we have

$$\|u - u_\varepsilon\|_V + \|p - p_\varepsilon\|_Q \leq C\varepsilon, \tag{50}$$

with  $C$  a constant depending only on  $\alpha_0$ ,  $\beta$ ,  $\gamma_0$ ,  $C_a$ ,  $C_b$  and  $C_{c_0}$ .

**Proof.** Subtracting (49) from (15), we find that  $(u - u_\varepsilon, p - p_\varepsilon)$  solves the system

$$\begin{aligned} A(u - u_\varepsilon) + B^*(p - p_\varepsilon) &= 0, \quad \text{in } V', \\ B(u - u_\varepsilon) - \varepsilon C_0(p - p_\varepsilon) &= -\varepsilon C_0 p, \quad \text{in } Q'. \end{aligned}$$

From Theorem 16, we infer that

$$\|(u - u_\varepsilon, p - p_\varepsilon)\|_{V \times Q} \leq C \|(0, -\varepsilon C_0 p)\|_{V' \times Q'}$$

with  $C$  depending only on  $\alpha_0$ ,  $\beta$ ,  $\gamma_0$ ,  $C_a$  and  $C_b$ . Thus

$$\|(u - u_\varepsilon, p - p_\varepsilon)\|_{V \times Q} \leq C C_{c_0} \varepsilon \|p\|_Q,$$

which proves (50).  $\square$

### 3.6. Case of a “fixed” augmentation

If the bilinear form  $c$  is given, the extra terms of the form  $c(p, Bu)$  arising from the previously considered T-coercivity operator can not be controlled as before, because there is no factor  $\varepsilon$  to adjust. Below, we assume that  $c$  is coercive on  $Q$ , namely that there exists  $\gamma > 0$  such that

$$c(p, p) \geq \gamma \|p\|_Q^2, \quad \forall p \in Q. \quad (51)$$

So, to control these extra terms, we introduce an operator  $C^{-1}$  in the expression of  $T$ , where  $C^{-1} \in \mathcal{L}(Q', Q)$  is defined by

$$c(C^{-1}g, q) = \langle g, q \rangle_{Q', Q}, \quad \forall g \in Q', \forall q \in Q.$$

One can easily check that the operator  $C^{-1}$  satisfies

$$(C_c)^{-1} \|g\|_{Q'} \leq \|C^{-1}g\|_Q \leq \gamma^{-1} \|g\|_{Q'}, \quad \forall g \in Q',$$

and

$$\langle g, C^{-1}g \rangle_{Q', Q} \geq \frac{\gamma}{C_c^2} \|g\|_{Q'}^2, \quad \forall g \in Q'. \quad (52)$$

**Theorem 18.** *Assume that (51) holds true and that the bilinear forms  $a$  and  $c$  are both symmetric and positive. Suppose in addition that there exists a constant  $\alpha_B > 0$  such that*

$$a(u, u) + \frac{\gamma}{2C_c^2} \|Bu\|_{Q'}^2 \geq \alpha_B \|u\|_V^2, \quad \forall u \in V, \quad (53)$$

then Problem (36) is well-posed and

$$\|(u, p)\|_{V \times Q} \leq C \|(f, g)\|_{V' \times Q'},$$

with  $C$  a constant depending only on  $\alpha_B$ ,  $\gamma$  and  $C_b$ .

**Proof.** For  $\eta, \mu > 0$ , we consider the mapping

$$\begin{aligned} T: V \times Q &\longrightarrow V \times Q \\ (u, p) &\longmapsto (u, -\eta p + \mu C^{-1}(Bu)). \end{aligned}$$

Then, using the definitions of  $C^{-1}$  and  $B$ , we compute

$$\begin{aligned} \mathcal{A}_c((u, p), T(u, p)) &= a(u, u) + b(u, p) - \eta b(u, p) + \mu b(u, C^{-1}Bu) + \eta c(p, p) - \mu c(p, C^{-1}Bu) \\ &= a(u, u) + (1 - \eta)b(u, p) + \mu \langle Bu, C^{-1}Bu \rangle_{Q', Q} + \eta c(p, p) - \mu \langle Bu, p \rangle_{Q', Q} \\ &= a(u, u) + (1 - \eta - \mu)b(u, p) + \mu \langle Bu, C^{-1}Bu \rangle_{Q', Q} + \eta c(p, p). \end{aligned}$$

Let us choose  $\eta, \mu > 0$  such that  $\eta + \mu = 1$  to cancel the second term above. To fix ideas, let  $\eta = \mu = 1/2$ , so that

$$T(u, p) = \left( u, -\frac{1}{2}p + \frac{1}{2}C^{-1}(Bu) \right) \quad (54)$$

and

$$\mathcal{A}_c((u, p), T(u, p)) = a(u, u) + \frac{1}{2} \langle Bu, C^{-1}Bu \rangle_{Q', Q} + \frac{1}{2} c(p, p).$$

Owing to (52) and (51), we deduce that

$$\mathcal{A}_c((u, p), T(u, p)) \geq a(u, u) + \frac{\gamma}{2C_c^2} \|Bu\|_{Q'}^2 + \frac{\gamma}{2} \|p\|_Q^2,$$

and the result follows.  $\square$

**Remark 19.** The T-coercivity estimate reads

$$\mathcal{A}_c((u, p), T(u, p)) \geq \alpha_B \|u\|_V^2 + \frac{\gamma}{2} \|p\|_Q^2, \quad (55)$$

so that it depends on  $\gamma$ , whereas it was independent of  $\varepsilon$  in the small perturbation case. Moreover, because of the term  $C^{-1}(Bu)$  in (54),  $\|T\|$  behaves as  $\gamma^{-1}$ . Nevertheless, the final stability estimate is robust because the value of the constant  $\gamma$  is fixed.

**Remark 20.** Note that Theorem 18 does not require the inf-sup condition (18) to be true. However, if (18) holds, then (53) is automatically satisfied. As a matter of fact, for any  $u \in V$ , using the decomposition  $u = u_0 + \bar{u}$  with  $u_0 \in \text{Ker } B$  and  $\bar{u} \in (\text{Ker } B)^\perp$ , we have seen in the proof of Theorem 8 that, for all  $0 < \theta < 1$ , it holds

$$a(u, u) \geq (1 - \theta)a(u_0, u_0) + \left(C_a - \frac{C_a}{\theta}\right) \|\bar{u}\|_V^2 \quad \text{and} \quad \|Bu\|_{Q'}^2 = \|B\bar{u}\|_Q^2 \geq \beta^{-2} \|\bar{u}\|_V^2.$$

Hence,

$$a(u, u) + \frac{\gamma}{2C_c^2} \|Bu\|_{Q'}^2 \geq (1 - \theta)a(u_0, u_0) + \left(C_a - \frac{C_a}{\theta} + \frac{\gamma}{2C_c^2} \beta^{-2}\right) \|\bar{u}\|_V^2.$$

We then observe that

$$\left(C_a - \frac{C_a}{\theta} + \frac{\gamma}{2C_c^2} \beta^{-2}\right) > 0, \quad \forall \theta \in \left(\left(1 + \frac{\gamma}{2C_c^2 C_a} \beta^{-2}\right)^{-1}, 1\right),$$

so (53) is obtained by choosing some  $\theta = \theta(C_a, \beta, C_c, \gamma)$  in the above interval.

In addition to the Stokes problem, let us see next how other typical examples of mixed formulations fall within the T-coercivity framework.

#### 4. Application to electromagnetism

Our goal is to solve the so-called quasi-static magnetic problem set in a anisotropic medium, surrounded by a perfect conductor (see [1, Section 6.4]). The medium is characterized by its dielectric permittivity  $\underline{\varepsilon}$  and its magnetic permeability  $\underline{\mu}$ .

Let  $\Omega$  be the domain of  $\mathbb{R}^3$  in which the problem is set. For simplicity, we assume that  $\Omega$  is simply connected, with a connected boundary. Moreover, we assume that  $\xi \in \{\underline{\varepsilon}, \underline{\mu}\}$  satisfy the following assumption:

$$\begin{cases} \xi \text{ is a real-valued, symmetric, measurable tensor field on } \Omega, \\ \exists \xi_-, \xi_+ > 0, \forall \mathbf{z} \in \mathbb{R}^3, \xi_- |\mathbf{z}|^2 \leq \xi \mathbf{z} \cdot \mathbf{z} \leq \xi_+ |\mathbf{z}|^2 \text{ a.e. in } \Omega. \end{cases} \quad (56)$$

Because one is dealing with symmetric tensors, if  $\xi$  fulfills (56), so does  $\xi^{-1}$ , with  $(\xi^{-1})_+ = (\xi_-)^{-1}$  and  $(\xi^{-1})_- = (\xi_+)^{-1}$ .

Given  $\mathbf{H}^* \in \mathbf{L}^2(\Omega)$ , such that  $\underline{\mu} \mathbf{H}^* \in \mathbf{H}_0(\text{div}; \Omega) \cap \mathbf{H}(\text{div } 0; \Omega)$  and  $\rho \in H^{-1}(\Omega)$ , the quasi-static magnetic problem amounts to finding  $\mathbf{E} \in \mathbf{L}^2(\Omega)$  such that

$$\begin{aligned} \underline{\mu}^{-1} \mathbf{curl} \mathbf{E} &= \mathbf{H}^*, & \text{in } \Omega, \\ \text{div}(\underline{\varepsilon} \mathbf{E}) &= \rho, & \text{in } \Omega, \\ \mathbf{E} \times \mathbf{n} &= 0, & \text{on } \partial\Omega. \end{aligned} \quad (57)$$

Under the assumptions on  $\underline{\varepsilon}$  and  $\underline{\mu}$ , on the one hand we note that  $\mathbf{E} \in \mathbf{H}_0(\mathbf{curl}; \Omega)$ . On the other hand, it is known that the problem (57) is well-posed, see for instance [1, Theorem 6.1.4]. Below, we propose to recover well-posedness using the T-coercivity approach.

##### 4.1. Proving well-posedness with T-coercivity

Classically, the electromagnetic energy is equal to  $(\underline{\varepsilon} \mathbf{E}, \mathbf{E})_{\mathbf{L}^2(\Omega)} + (\underline{\mu} \mathbf{H}, \mathbf{H})_{\mathbf{L}^2(\Omega)}$ , where  $\mathbf{H}$  is the magnetic field. For our problem with the electric field  $\mathbf{E}$  as the unknown, it can be expressed as  $(\underline{\varepsilon} \mathbf{E}, \mathbf{E})_{\mathbf{L}^2(\Omega)} + (\underline{\mu}^{-1} \mathbf{curl} \mathbf{E}, \mathbf{curl} \mathbf{E})_{\mathbf{L}^2(\Omega)}$  because  $\mathbf{H}^* \sim -\partial_t \mathbf{H}$ . Indeed, under the assumption (56) made on  $\underline{\varepsilon}$  and  $\underline{\mu}$ , we note that we can endow  $\mathbf{H}_0(\mathbf{curl}; \Omega)$  with the scalar product  $(\cdot, \cdot)_{\underline{\varepsilon}, \underline{\mu}^{-1} \mathbf{curl}} : (\mathbf{u}, \mathbf{v}) \mapsto (\underline{\varepsilon} \mathbf{u}, \mathbf{v})_{\mathbf{L}^2(\Omega)} + (\underline{\mu}^{-1} \mathbf{curl} \mathbf{u}, \mathbf{curl} \mathbf{v})_{\mathbf{L}^2(\Omega)}$ , and the associated scaled norm

$$\|\mathbf{u}\|_{\underline{\varepsilon}, \underline{\mu}^{-1} \mathbf{curl}} = \left( (\underline{\varepsilon} \mathbf{u}, \mathbf{u})_{\mathbf{L}^2(\Omega)} + (\underline{\mu}^{-1} \mathbf{curl} \mathbf{u}, \mathbf{curl} \mathbf{u})_{\mathbf{L}^2(\Omega)} \right)^{1/2}$$

is equivalent to the "natural" norm. We endow  $H_0^1(\Omega)$  with the scalar product  $(\cdot, \cdot)_{1,\underline{\varepsilon}} : (p, q) \mapsto (\underline{\varepsilon}\nabla p, \nabla q)_{L^2(\Omega)}$ , and the associated scaled norm

$$\|q\|_{1,\underline{\varepsilon}} = \left( (\underline{\varepsilon}\nabla q, \nabla q)_{L^2(\Omega)} \right)^{1/2}$$

is equivalent to  $\|\cdot\|_{H^1(\Omega)}$  according to Poincaré inequality. In this setting,  $\mathbb{1}_{H_0^1(\Omega) \rightarrow H^{-1}(\Omega)}$  is the isomorphism defined by

$$\langle \mathbb{1}_{H_0^1(\Omega) \rightarrow H^{-1}(\Omega)} p, q \rangle_{H^{-1}(\Omega), H_0^1(\Omega)} = (p, q)_{1,\underline{\varepsilon}} = (\underline{\varepsilon}\nabla p, \nabla q)_{L^2(\Omega)}, \quad \forall p, q \in H_0^1(\Omega),$$

while the norm in  $H^{-1}(\Omega)$  is

$$\|g\|_{-1,\underline{\varepsilon}^{-1}} = \sup_{q \in H_0^1(\Omega) \setminus \{0\}} \frac{\langle g, q \rangle_{H^{-1}(\Omega), H_0^1(\Omega)}}{\|q\|_{1,\underline{\varepsilon}}}, \quad \forall g \in H^{-1}(\Omega).$$

Bearing in mind that  $\mathbf{curl}(\nabla p) = 0$ , it follows that

$$\|q\|_{1,\underline{\varepsilon}} = \left( (\underline{\varepsilon}\nabla q, \nabla q)_{L^2(\Omega)} \right)^{1/2} = \|\nabla q\|_{\underline{\varepsilon}, \underline{\mu}^{-1} \mathbf{curl}}, \quad \forall q \in H_0^1(\Omega).$$

Finally, for  $\xi \in \{\underline{\varepsilon}, \underline{\varepsilon}^{-1}, \underline{\mu}, \underline{\mu}^{-1}\}$ , we use the scalar product  $(\cdot, \cdot)_\xi : (\mathbf{u}, \mathbf{v}) \mapsto (\xi \mathbf{u}, \mathbf{v})_{L^2(\Omega)}$ , and the associated scaled norm  $\|\cdot\|_\xi$  in  $L^2(\Omega)$ .

First, for  $\mathbf{H}^* \in \mathbf{H}_0(\operatorname{div}; \Omega) \cap \mathbf{H}(\operatorname{div}0; \Omega)$  and  $\rho \in H^{-1}(\Omega)$ , we observe that the equivalent weak formulation of Problem (57) reads: find  $\mathbf{E} \in \mathbf{H}_0(\mathbf{curl}; \Omega)$  such that

$$\begin{aligned} (\underline{\mu}^{-1} \mathbf{curl} \mathbf{E}, \mathbf{curl} \mathbf{v})_{L^2(\Omega)} &= (\mathbf{H}^*, \mathbf{curl} \mathbf{v})_{L^2(\Omega)}, \quad \forall \mathbf{v} \in \mathbf{H}_0(\mathbf{curl}; \Omega), \\ (\underline{\varepsilon} \mathbf{E}, \nabla q)_{L^2(\Omega)} &= -\langle \rho, q \rangle_{H^{-1}(\Omega), H_0^1(\Omega)}, \quad \forall q \in H_0^1(\Omega). \end{aligned}$$

Second, in order to fit (57) into the abstract framework (15), we introduce an artificial pressure unknown  $\tilde{p}$  by adding a term  $(\underline{\varepsilon} \mathbf{v}, \nabla \tilde{p})_{L^2(\Omega)}$  in the first equation. The previous formulation becomes: find  $(\mathbf{E}, \tilde{p}) \in \mathbf{H}_0(\mathbf{curl}; \Omega) \times H_0^1(\Omega)$  such that

$$\begin{aligned} (\underline{\mu}^{-1} \mathbf{curl} \mathbf{E}, \mathbf{curl} \mathbf{v})_{L^2(\Omega)} + (\underline{\varepsilon} \mathbf{v}, \nabla \tilde{p})_{L^2(\Omega)} &= (\mathbf{H}^*, \mathbf{curl} \mathbf{v})_{L^2(\Omega)}, \quad \forall \mathbf{v} \in \mathbf{H}_0(\mathbf{curl}; \Omega), \\ (\underline{\varepsilon} \mathbf{E}, \nabla q)_{L^2(\Omega)} &= -\langle \rho, q \rangle_{H^{-1}(\Omega), H_0^1(\Omega)}, \quad \forall q \in H_0^1(\Omega). \end{aligned} \quad (58)$$

Indeed, one can easily check that  $(\mathbf{E}, \tilde{p})$  is solution of (58) if and only if  $\tilde{p} = 0$  and  $\mathbf{E}$  is solution of (57). So, defining the bilinear forms

$$\begin{aligned} a_{\underline{\mu}^{-1}}(\mathbf{u}, \mathbf{v}) &= (\underline{\mu}^{-1} \mathbf{curl} \mathbf{u}, \mathbf{curl} \mathbf{v})_{L^2(\Omega)}, \quad \forall \mathbf{u} \in \mathbf{H}_0(\mathbf{curl}; \Omega), \forall \mathbf{v} \in \mathbf{H}_0(\mathbf{curl}; \Omega), \\ b_{\underline{\varepsilon}}(\mathbf{v}, q) &= (\underline{\varepsilon} \mathbf{v}, \nabla q)_{L^2(\Omega)}, \quad \forall \mathbf{v} \in \mathbf{H}_0(\mathbf{curl}; \Omega), \forall q \in H_0^1(\Omega), \end{aligned}$$

the *all-in-one* bilinear form of Maxwell problem is given by

$$\mathcal{A}((\mathbf{E}, \tilde{p}), (\mathbf{v}, q)) = a_{\underline{\mu}^{-1}}(\mathbf{E}, \mathbf{v}) + b_{\underline{\varepsilon}}(\mathbf{v}, \tilde{p}) + b_{\underline{\varepsilon}}(\mathbf{E}, q). \quad (59)$$

Thanks to the introduction of scaled norms, we find that the bilinear form  $a_{\underline{\mu}^{-1}}$  is continuous on  $\mathbf{H}_0(\mathbf{curl}; \Omega) \times \mathbf{H}_0(\mathbf{curl}; \Omega)$  with a continuity constant  $C_a = 1$ , while the bilinear form  $b_{\underline{\varepsilon}}$  is continuous on  $\mathbf{H}_0(\mathbf{curl}; \Omega) \times H_0^1(\Omega)$  with a continuity constant  $C_b = 1$ . Besides, we have

$$\left| (\mathbf{H}^*, \mathbf{curl} \mathbf{v})_{L^2(\Omega)} - \langle \rho, q \rangle_{H^{-1}(\Omega), H_0^1(\Omega)} \right| \leq \|\mathbf{H}^*\|_{\underline{\mu}} \|\mathbf{curl} \mathbf{v}\|_{\underline{\mu}^{-1}} + \|\rho\|_{-1,\underline{\varepsilon}^{-1}} \|q\|_{1,\underline{\varepsilon}},$$

so that  $C_\ell \leq \left( \|\mathbf{H}^*\|_{\underline{\mu}}^2 + \|\rho\|_{-1,\underline{\varepsilon}^{-1}}^2 \right)^{1/2}$ .

Let us give an explicit expression of the abstract operators

$$\mathbb{B}_{\underline{\varepsilon}} \in \mathcal{L}(\mathbf{H}_0(\mathbf{curl}; \Omega), H_0^1(\Omega)), \quad \mathbb{L}_{\mathbb{B}_{\underline{\varepsilon}}} \in \mathcal{L}(H_0^1(\Omega), (\operatorname{Ker} \mathbb{B}_{\underline{\varepsilon}})^\perp)$$

corresponding to this problem. Note that for  $\mathbf{u} \in \mathbf{H}_0(\mathbf{curl}; \Omega)$ ,

$$\mathbb{B}_{\underline{\varepsilon}} \mathbf{u} = 0 \iff (\underline{\varepsilon} \mathbf{u}, \nabla q)_{L^2(\Omega)} = 0, \quad \forall q \in H_0^1(\Omega) \iff \operatorname{div}(\underline{\varepsilon} \mathbf{u}) = 0.$$

Hence,

$$\operatorname{Ker} \mathbb{B}_{\underline{\varepsilon}} = \mathbf{K}_N(\Omega; \underline{\varepsilon}), \quad \text{where } \mathbf{K}_N(\Omega; \underline{\varepsilon}) = \{ \mathbf{v} \in \mathbf{H}_0(\mathbf{curl}; \Omega), \operatorname{div}(\underline{\varepsilon} \mathbf{v}) = 0 \}. \quad (60)$$

In addition, one knows that (see (6.16) in [1])

$$(\text{Ker } \mathbf{B}_{\underline{\varepsilon}})^\perp = \{ \mathbf{v} \in \mathbf{H}_0(\mathbf{curl}; \Omega), \exists q \in H_0^1(\Omega), \mathbf{v} = \nabla q \}. \quad (61)$$

With those results, we can explicit  $\mathbf{L}_{B_{\underline{\varepsilon}}}$ . On the one hand, by definition of  $b_{\underline{\varepsilon}}$ , we observe that

$$b_{\underline{\varepsilon}}(\mathbf{L}_{B_{\underline{\varepsilon}}} p, q) = (\underline{\varepsilon}(\mathbf{L}_{B_{\underline{\varepsilon}}} p), \nabla q)_{\mathbf{L}^2(\Omega)}, \quad \forall p, q \in H_0^1(\Omega). \quad (62)$$

On the other hand, according to (19), one has

$$b_{\underline{\varepsilon}}(\mathbf{L}_{B_{\underline{\varepsilon}}} p, q) = (p, q)_{1, \underline{\varepsilon}} = (\underline{\varepsilon} \nabla p, \nabla q)_{\mathbf{L}^2(\Omega)}, \quad \forall p, q \in H_0^1(\Omega). \quad (63)$$

Putting (61), (62) and (63) together, we deduce that

$$\mathbf{L}_{B_{\underline{\varepsilon}}} p = \nabla p.$$

Then, for all  $p \in H_0^1(\Omega)$ , it follows that

$$\|\mathbf{L}_{B_{\underline{\varepsilon}}} p\|_{\underline{\varepsilon}, \underline{\mu}^{-1} \mathbf{curl}}^2 = \|\nabla p\|_{\underline{\varepsilon}, \underline{\mu}^{-1} \mathbf{curl}}^2 = \|p\|_{1, \underline{\varepsilon}}^2,$$

hence  $\mathbf{L}_{B_{\underline{\varepsilon}}}$  is an isometry, so  $\mathbf{L}_{B_{\underline{\varepsilon}}}$  satisfies (21) with  $\beta = 1$ , *i.e.* the inf-sup condition (18) holds.

Going back to  $\text{Ker } \mathbf{B}_{\underline{\varepsilon}}$  (cf. the characterization (60)), we recall Weber inequality [38] (see also [1, Theorem 6.1.4]): there exists  $C_K > 1$  such that

$$\|\mathbf{k}\|_{\underline{\varepsilon}, \underline{\mu}^{-1} \mathbf{curl}} \leq C_K \|\mathbf{curl} \mathbf{k}\|_{\underline{\mu}^{-1}}, \quad \forall \mathbf{k} \in \mathbf{K}_N(\Omega; \underline{\varepsilon}). \quad (64)$$

The fact that  $C_K > 1$  stems from the definition of the scaled norms. Hence, Weber inequality (64) says exactly that the form  $a_{\underline{\mu}^{-1}}$  is coercive on  $\text{Ker } \mathbf{B}_{\underline{\varepsilon}}$ , so that all the conditions of Theorem 8 are fulfilled, with  $\alpha_0 = (C_K)^{-2} < 1$ . Precisely, Theorem 8 states that the bilinear form  $\mathcal{A}$  is T-coercive for the mapping

$$\begin{aligned} \mathbf{T} : \mathbf{H}_0(\mathbf{curl}; \Omega) \times H_0^1(\Omega) &\longrightarrow \mathbf{H}_0(\mathbf{curl}; \Omega) \times H_0^1(\Omega) \\ (\mathbf{E}, \tilde{p}) &\longmapsto \left( \mathbf{E} + \nabla \tilde{p}, -\tilde{p} + \frac{3}{4} \phi_{\mathbf{E}} \right), \end{aligned}$$

where  $\phi_{\mathbf{E}} = \mathbf{B}_{\underline{\varepsilon}} \mathbf{E} \in H_0^1(\Omega)$ . Note that, for any  $\mathbf{v} \in \mathbf{H}_0(\mathbf{curl}; \Omega)$ , we have

$$(\mathbf{B}_{\underline{\varepsilon}} \mathbf{v}, q)_{1, \underline{\varepsilon}} = b_{\underline{\varepsilon}}(\mathbf{v}, q) = (\underline{\varepsilon} \mathbf{v}, \nabla q)_{\mathbf{L}^2(\Omega)}, \quad \forall q \in H_0^1(\Omega).$$

Therefore,  $\mathbf{B}_{\underline{\varepsilon}} \mathbf{v}$  is the unique  $\phi_{\mathbf{v}} \in H_0^1(\Omega)$  satisfying

$$(\underline{\varepsilon} \nabla \phi_{\mathbf{v}}, \nabla q)_{\mathbf{L}^2(\Omega)} = (\underline{\varepsilon} \mathbf{v}, \nabla q)_{\mathbf{L}^2(\Omega)}, \quad \forall q \in H_0^1(\Omega). \quad (65)$$

Furthermore, following (30), it holds that

$$\mathcal{A}((\mathbf{E}, \tilde{p}), \mathbf{T}(\mathbf{E}, \tilde{p})) \geq \frac{(C_K)^{-2}}{4} \|\mathbf{E}\|_{\underline{\varepsilon}, \underline{\mu}^{-1} \mathbf{curl}}^2 + \frac{1}{2} \|\tilde{p}\|_{1, \underline{\varepsilon}}^2 \geq \underline{\alpha} (\|\mathbf{E}\|_{\underline{\varepsilon}, \underline{\mu}^{-1} \mathbf{curl}}^2 + \|\tilde{p}\|_{1, \underline{\varepsilon}}^2),$$

with  $\underline{\alpha} = \frac{(C_K)^{-2}}{4}$ .

To get the stability constant, we need to compute  $\|\mathbf{T}\|$ , that is, bound  $\|\mathbf{T}(\mathbf{E}, \tilde{p})\|_{\mathbf{H}_0(\mathbf{curl}; \Omega) \times H_0^1(\Omega)}$  for  $(\mathbf{E}, \tilde{p}) \in \mathbf{H}_0(\mathbf{curl}; \Omega) \times H_0^1(\Omega)$ . We find that

$$\begin{aligned} \|\mathbf{T}(\mathbf{E}, \tilde{p})\|_{\mathbf{H}_0(\mathbf{curl}; \Omega) \times H_0^1(\Omega)}^2 &= \|\mathbf{E} + \nabla \tilde{p}\|_{\underline{\varepsilon}, \underline{\mu}^{-1} \mathbf{curl}}^2 + \left\| -\tilde{p} + \frac{3}{4} \phi_{\mathbf{E}} \right\|_{1, \underline{\varepsilon}}^2 \\ &\leq 2\|\mathbf{E}\|_{\underline{\varepsilon}, \underline{\mu}^{-1} \mathbf{curl}}^2 + 2\|\nabla \tilde{p}\|_{\underline{\varepsilon}, \underline{\mu}^{-1} \mathbf{curl}}^2 + 2\|\tilde{p}\|_{1, \underline{\varepsilon}}^2 + 2 \cdot \left(\frac{3}{4}\right)^2 \|\phi_{\mathbf{E}}\|_{1, \underline{\varepsilon}}^2 \\ &\leq 2\|\mathbf{E}\|_{\underline{\varepsilon}, \underline{\mu}^{-1} \mathbf{curl}}^2 + 4\|\tilde{p}\|_{1, \underline{\varepsilon}}^2 + 2 \cdot \left(\frac{3}{4}\right)^2 \|\mathbf{E}\|_{\underline{\varepsilon}, \underline{\mu}^{-1} \mathbf{curl}}^2 \\ &\leq 4(\|\mathbf{E}\|_{\underline{\varepsilon}, \underline{\mu}^{-1} \mathbf{curl}}^2 + \|\tilde{p}\|_{1, \underline{\varepsilon}}^2), \end{aligned}$$

where we used that  $\|\phi_{\mathbf{E}}\|_{1, \underline{\varepsilon}} = \|\mathbf{B}_{\underline{\varepsilon}} \mathbf{E}\|_{1, \underline{\varepsilon}} \leq \|\mathbf{E}\|_{\underline{\varepsilon}, \underline{\mu}^{-1} \mathbf{curl}}$  thanks to (31). Therefore,  $\|\mathbf{T}\| \leq 2$ .

Applying (4), we conclude that

$$\|\mathbf{E}\|_{\underline{\varepsilon}, \underline{\mu}^{-1} \mathbf{curl}} \leq 8C_K^2 (\|\mathbf{H}^* \|_{\underline{\mu}}^2 + \|\rho\|_{-1, \underline{\varepsilon}^{-1}}^2)^{1/2}. \quad (66)$$

#### 4.2. Optimized bounds

To achieve T-coercivity, the abstract theory does not take into account the so-called double orthogonality property (or Helmholtz decomposition), which states that for all  $\mathbf{k} \in \mathbf{K}_N(\Omega; \underline{\varepsilon})$  and all  $q \in H_0^1(\Omega)$ , one has  $(\underline{\varepsilon} \mathbf{k}, \nabla q)_{L^2(\Omega)} = (\underline{\mu}^{-1} \mathbf{curl} \mathbf{k}, \mathbf{curl}(\nabla q))_{L^2(\Omega)} = 0$ , so that

$$\|\mathbf{k} + \nabla q\|_{\underline{\varepsilon}, \underline{\mu}^{-1} \mathbf{curl}}^2 = \|\mathbf{k}\|_{\underline{\varepsilon}, \underline{\mu}^{-1} \mathbf{curl}}^2 + \|q\|_{1, \underline{\varepsilon}}^2.$$

Indeed, given  $\mathbf{E} \in \mathbf{H}_0(\mathbf{curl}; \Omega)$ , we note that, with the help of  $\phi_E \in H_0^1(\Omega)$  solving (65), it holds that  $\mathbf{k}_E = \mathbf{E} - \nabla \phi_E \in \mathbf{K}_N(\Omega; \underline{\varepsilon})$ .

We sketch below how one can improve the estimates, see [22] for further details. Let us choose

$$\begin{aligned} \mathsf{T}_{opt} : \mathbf{H}_0(\mathbf{curl}; \Omega) \times H_0^1(\Omega) &\longrightarrow \mathbf{H}_0(\mathbf{curl}; \Omega) \times H_0^1(\Omega) \\ (\mathbf{E}, \tilde{p}) &\longmapsto (\mathbf{k}_E + \nabla \tilde{p}, \phi_E). \end{aligned}$$

Thanks to the double orthogonality property, one finds easily that  $\mathsf{T}_{opt}$  is an isometry and that

$$\begin{aligned} \mathcal{A}((\mathbf{E}, \tilde{p}), \mathsf{T}_{opt}(\mathbf{E}, \tilde{p})) &= \|\mathbf{curl} \mathbf{k}_E\|_{\underline{\mu}^{-1}}^2 + \|\tilde{p}\|_{1, \underline{\varepsilon}}^2 + \|\phi_E\|_{1, \underline{\varepsilon}}^2 \\ &\geq (C_K)^{-2} (\|\mathbf{E}\|_{\underline{\varepsilon}, \underline{\mu}^{-1} \mathbf{curl}}^2 + \|\tilde{p}\|_{1, \underline{\varepsilon}}^2). \end{aligned}$$

Applying (4), we have the optimized stability estimate

$$\|\mathbf{E}\|_{\underline{\varepsilon}, \underline{\mu}^{-1} \mathbf{curl}} \leq C_K^2 (\|\mathbf{H}^* \|_{\underline{\mu}}^2 + \|\rho\|_{-1, \underline{\varepsilon}^{-1}}^2)^{1/2}. \quad (67)$$

We conclude that, for all possible choices of coefficients  $\underline{\varepsilon}$  and  $\underline{\mu}$ , there is only a factor 8 difference between the stability constant obtained via the abstract T-coercivity approach, see (66), and the optimized stability constant which relies explicitly on the double orthogonality property, see (67). This shows the robustness of the abstract theory.

**Remark 21.** One can obtain similar results in more general geometries, such as a non-simply-connected domain, or a non-connected boundary.

## 5. Application to nearly-incompressible elasticity

In this section, we apply the T-coercivity framework to the equations of elasticity, assuming homogeneous Dirichlet boundary conditions. Let  $\Omega \subset \mathbb{R}^d$  be a domain, where  $2 \leq d \leq 3$ . For a prescribed body force  $\mathbf{f} \in \mathbf{H}^{-1}(\Omega)$ , we look for the displacement  $\mathbf{u} \in \mathbf{H}^1(\Omega)$  such that

$$\begin{aligned} -\operatorname{div}(\sigma(\mathbf{u})) &= \mathbf{f}, \quad \text{in } \Omega, \\ \mathbf{u} &= 0, \quad \text{on } \partial\Omega, \end{aligned} \quad (68)$$

where  $\sigma(\mathbf{u})$  denotes the stress tensor. We assume that it is given by Hooke's law

$$\sigma(\mathbf{u}) = 2\mu \varepsilon(\mathbf{u}) + \lambda(\operatorname{div} \mathbf{u}) \mathbf{I},$$

where  $\lambda, \mu > 0$  are the Lamé coefficients of the material and  $\varepsilon(\mathbf{u}) = \frac{1}{2}(\nabla \mathbf{u} + (\nabla \mathbf{u})^T)$  is the linearized strain tensor. Thanks to Korn inequality [25], the space  $\mathbf{H}_0^1(\Omega)$  is here endowed with the scalar product

$$(\mathbf{u}, \mathbf{v}) \longmapsto \int_{\Omega} \varepsilon(\mathbf{u}) : \varepsilon(\mathbf{v}) \, dx,$$

whose associated norm  $\mathbf{u} \mapsto \|\varepsilon(\mathbf{u})\|$  is equivalent to the  $\mathbf{H}^1(\Omega)$ -norm in  $\mathbf{H}_0^1(\Omega)$ . Introducing the new unknown  $p = \lambda \operatorname{div} \mathbf{u}$ , the elasticity system (68) can be written in mixed form as follows: find  $\mathbf{u} \in \mathbf{H}_0^1(\Omega)$  and  $p \in L_0^2(\Omega)$  such that

$$\begin{aligned} -2\mu \operatorname{div}(\varepsilon(\mathbf{u})) - \nabla p &= \mathbf{f}, & \text{in } \Omega, \\ \operatorname{div} \mathbf{u} - \frac{1}{\lambda} p &= 0, & \text{in } \Omega. \end{aligned}$$

Or equivalently, in variational form: find  $\mathbf{u} \in \mathbf{H}_0^1(\Omega)$  and  $p \in L_0^2(\Omega)$  such that

$$\begin{aligned} a(\mathbf{u}, \mathbf{v}) + b(\mathbf{v}, p) &= \langle \mathbf{f}, \mathbf{v} \rangle, & \forall \mathbf{v} \in \mathbf{H}_0^1(\Omega), \\ b(\mathbf{u}, q) - \frac{1}{\lambda} c_0(p, q) &= 0, & \forall q \in L_0^2(\Omega), \end{aligned} \quad (69)$$

with

$$a(\mathbf{u}, \mathbf{v}) = 2\mu \int_{\Omega} \varepsilon(\mathbf{u}) : \varepsilon(\mathbf{v}) \, dx, \quad b(\mathbf{v}, q) = \int_{\Omega} q \operatorname{div} \mathbf{v} \, dx \quad \text{and} \quad c_0(p, q) = \int_{\Omega} p q \, dx.$$

For nearly-incompressible materials, the first Lamé coefficient  $\lambda$  goes to infinity, so that  $\lambda^{-1}$  goes to zero. Therefore, (69) can be seen as a small perturbation of Stokes system.

Since the bilinear form  $a$  is coercive on the whole space  $\mathbf{H}_0^1(\Omega)$ , we can directly apply Theorem 10 in the special case of Remark 11. The bilinear form  $a$  is continuous and coercive, with  $C_a = \alpha = 2\mu$ . In addition, the bilinear form  $b$  is continuous and satisfies the inf-sup condition (18) with  $\beta = C_{\operatorname{div}}$  since  $b$  is the same form – except to the sign – as for Stokes problem. Then, Theorem 10 furnishes that the *all-in-one* bilinear form  $\mathcal{A}_c$  defined by

$$\mathcal{A}_c((\mathbf{u}, p), (\mathbf{v}, q)) = 2\mu \int_{\Omega} \varepsilon(\mathbf{u}) : \varepsilon(\mathbf{v}) \, dx + \int_{\Omega} p \operatorname{div} \mathbf{v} \, dx + \int_{\Omega} q \operatorname{div} \mathbf{u} \, dx - \frac{1}{\lambda} \int_{\Omega} p q \, dx \quad (70)$$

is T-coercive for the mapping

$$\begin{aligned} \mathbb{T} : \mathbf{H}_0^1(\Omega) \times L_0^2(\Omega) &\longrightarrow \mathbf{H}_0^1(\Omega) \times L_0^2(\Omega) \\ (\mathbf{u}, p) &\longmapsto (2\mu C_{\operatorname{div}}^2 \mathbf{u} + \mathbf{v}_{-p}, -2\mu C_{\operatorname{div}}^2 p), \end{aligned} \quad (71)$$

and (37) implies that

$$\mathcal{A}_c((\mathbf{u}, p), \mathbb{T}(\mathbf{u}, p)) \geq 2\mu^2 C_{\operatorname{div}}^2 \|\varepsilon(\mathbf{u})\|^2 + \left( \frac{1}{2} + \frac{2\mu}{\lambda} C_{\operatorname{div}}^2 \right) \|p\|^2.$$

Note that this estimate is robust in the incompressible limit, namely for large values of  $\lambda$ .

Finally, replacing  $\nu$  by  $2\mu$  in (11) and using (4), we get that the unique solution of (69) satisfies

$$\|(\mathbf{u}, p)\|_{\mathbf{H}_0^1(\Omega) \times L_0^2(\Omega)} \leq \frac{2 \max\left(2\sqrt{2}\mu C_{\operatorname{div}}^2, C_{\operatorname{div}}(2 + 4\mu^2 C_{\operatorname{div}}^2)^{1/2}\right)}{\min(4\mu^2 C_{\operatorname{div}}^2, 1 + 4\mu\lambda^{-1} C_{\operatorname{div}}^2)} \|\mathbf{f}\|_{(\mathbf{H}_0^1(\Omega))'},$$

where  $\|\cdot\|_{(\mathbf{H}_0^1(\Omega))'}$  denotes the dual norm of  $\|\varepsilon(\cdot)\|$ .

## 6. Application to neutron diffusion

Let  $\Omega \subset \mathbb{R}^d$  be a domain, where  $2 \leq d \leq 3$ . We consider the neutron diffusion equation with zero flux boundary condition: given a prescribed fission source  $S_f \in L^2(\Omega)$ , find  $u \in H^1(\Omega)$  such that

$$\begin{aligned} -\operatorname{div}(D\nabla u) + \sigma u &= S_f, & \text{in } \Omega, \\ u &= 0, & \text{on } \partial\Omega, \end{aligned} \quad (72)$$

where  $u$ ,  $D$ , and  $\sigma$  denote respectively the neutron flux, the diffusion coefficient and the macroscopic absorption cross section. It is assumed that the diffusion coefficient  $D$  fulfills (56), and that the macroscopic absorption cross section is such that

$$\begin{cases} \sigma \text{ is a real-valued measurable scalar field on } \Omega, \\ \exists \sigma_-, \sigma_+ > 0, \sigma_- \leq \sigma \leq \sigma_+ \text{ a.e. in } \Omega. \end{cases} \quad (73)$$

Because  $S_f \in L^2(\Omega)$ , one has  $D\nabla u \in \mathbf{H}(\text{div}; \Omega)$ . This problem can be recast equivalently in mixed form, introducing the auxiliary unknown  $\mathbf{p} = -D\nabla u$ , called the neutron current. It reads: find  $(u, \mathbf{p}) \in H_0^1(\Omega) \times \mathbf{H}(\text{div}; \Omega)$  such that

$$\begin{aligned} \text{div } \mathbf{p} + \sigma u &= S_f, & \text{in } \Omega, \\ D^{-1} \mathbf{p} + \nabla u &= 0, & \text{in } \Omega. \end{aligned} \quad (74)$$

It can be shown that equivalent weak form is: find  $(u, \mathbf{p}) \in L^2(\Omega) \times \mathbf{H}(\text{div}; \Omega)$  such that

$$\int_{\Omega} (v \text{div } \mathbf{p} + \sigma uv - D^{-1} \mathbf{p} \cdot \mathbf{q} + u \text{div } \mathbf{q}) \, dx = \int_{\Omega} S_f v \, dx \quad \forall (v, \mathbf{q}) \in L^2(\Omega) \times \mathbf{H}(\text{div}; \Omega). \quad (75)$$

**Remark 22.** Among other things, one can recover that the solution  $u \in L^2(\Omega)$  from the weak form (75) is such that  $u \in H^1(\Omega)$ , and that  $u = 0$  on  $\partial\Omega$ .

### 6.1. Proving well-posedness with T-coercivity

Defining the bilinear forms

$$\begin{aligned} a_{D^{-1}}(\mathbf{p}, \mathbf{q}) &= (D^{-1} \mathbf{p}, \mathbf{q})_{L^2(\Omega)}, & \forall \mathbf{p} \in \mathbf{H}(\text{div}; \Omega), \forall \mathbf{q} \in \mathbf{H}(\text{div}; \Omega), \\ b(\mathbf{q}, v) &= -(\text{div } \mathbf{q}, v)_{L^2(\Omega)}, & \forall \mathbf{q} \in \mathbf{H}(\text{div}; \Omega), \forall v \in L^2(\Omega), \\ c_{\sigma}(u, v) &= (\sigma u, v)_{L^2(\Omega)}, & \forall u \in L^2(\Omega), \forall v \in L^2(\Omega), \end{aligned}$$

the *all-in-one* bilinear form of the diffusion problem is given by

$$\mathcal{A}_c((\mathbf{p}, u), (\mathbf{q}, v)) = a_{D^{-1}}(\mathbf{p}, \mathbf{q}) + b(\mathbf{q}, u) + b(\mathbf{p}, v) - c_{\sigma}(u, v). \quad (76)$$

Here, we are in the case of a “fixed” augmentation, as treated in Section 3.6.

Let us check below that all the conditions of Theorem 18 are fulfilled. First,  $c_{\sigma}$  is coercive on  $L^2(\Omega)$  with  $\gamma = \sigma_-$ . Then,  $a_{D^{-1}}$  fulfills (13) with  $C_a = (D_-)^{-1}$ , whereas  $b$  fulfills (14) with  $C_b = 1$ . Finally, we look for the condition (53). It is straightforward to check that, for all  $\mathbf{p} \in \mathbf{H}(\text{div}; \Omega)$ ,  $B\mathbf{p} = B\mathbf{p} = -\text{div } \mathbf{p}$ . Hence

$$\begin{aligned} a_{D^{-1}}(\mathbf{p}, \mathbf{p}) + \frac{\gamma}{2C_c^2} \|B\mathbf{p}\|^2 &= (D^{-1} \mathbf{p}, \mathbf{p}) + \frac{\sigma_-}{2\sigma_+^2} \|\text{div } \mathbf{p}\|^2 \\ &\geq \min\left((D_+)^{-1}, \frac{\sigma_-}{2\sigma_+^2}\right) \|\mathbf{p}\|_{\mathbf{H}(\text{div}; \Omega)}^2. \end{aligned}$$

Then, Theorem 18 establishes that the bilinear form  $\mathcal{A}_c$  is T-coercive for the mapping (54)

$$\begin{aligned} \mathbb{T} : \mathbf{H}(\text{div}; \Omega) \times L^2(\Omega) &\longrightarrow \mathbf{H}(\text{div}; \Omega) \times L^2(\Omega) \\ (\mathbf{p}, u) &\longmapsto \left(\mathbf{p}, \frac{1}{2}(-u - \sigma^{-1} \text{div } \mathbf{p})\right). \end{aligned}$$

Furthermore, using the estimate (55), it holds that

$$\mathcal{A}_c((\mathbf{p}, u), \mathbb{T}(\mathbf{p}, u)) \geq \min\left((D_+)^{-1}, \frac{\sigma_-}{2\sigma_+^2}\right) \|\mathbf{p}\|_{\mathbf{H}(\text{div}; \Omega)}^2 + \frac{\sigma_-}{2} \|u\|^2 \geq \underline{\alpha} \|(\mathbf{p}, u)\|_{\mathbf{H}(\text{div}; \Omega) \times L^2(\Omega)}^2, \quad (77)$$

with  $\underline{\alpha} = \frac{1}{2} \min(2(D_+)^{-1}, \sigma_-(\sigma_+)^{-2}, \sigma_-)$ .

There remains to estimate  $\|\mathbb{T}\|$ . One has

$$\begin{aligned} \|\mathbb{T}(\mathbf{p}, u)\|_{\mathbf{H}(\text{div}; \Omega) \times L^2(\Omega)}^2 &= \|\mathbf{p}\|_{\mathbf{H}(\text{div}; \Omega)}^2 + \frac{1}{4} \|-u - \sigma^{-1} \text{div } \mathbf{p}\|_{L^2(\Omega)}^2 \\ &\leq \|\mathbf{p}\|_{\mathbf{H}(\text{div}; \Omega)}^2 + \frac{1}{4} \left( (1+3) \|u\|_{L^2(\Omega)}^2 + \left(1 + \frac{1}{3}\right) (\sigma_-)^{-2} \|\text{div } \mathbf{p}\|_{L^2(\Omega)}^2 \right) \\ &\leq \left(1 + \frac{1}{3} (\sigma_-)^{-2}\right) \|(\mathbf{p}, u)\|_{\mathbf{H}(\text{div}; \Omega) \times L^2(\Omega)}^2, \end{aligned}$$



so that  $\|\mathbb{T}\| \leq \left(1 + \frac{1}{3}(\sigma_-)^{-2}\right)^{1/2}$ . Applying (4), we conclude that

$$\|(\mathbf{p}, \mathbf{u})\|_{\mathbf{H}(\operatorname{div}; \Omega) \times L^2(\Omega)} \leq \frac{2\left(1 + \frac{1}{3}(\sigma_-)^{-2}\right)^{1/2}}{\min(2(D_+)^{-1}, \sigma_-(\sigma_+)^{-2}, \sigma_-)} \|S_f\|_{L^2(\Omega)}.$$

**Remark 23.** Some of those computations can be found in [24, 33]. Here, we see them as a consequence of the general result stated in Theorem 18. Note that in [24, 33], the T-coercivity estimate (77) is obtained with a constant  $\underline{\alpha}' = \frac{1}{2} \min(2(D_+)^{-1}, (\sigma_+)^{-1}, \sigma_-)$ , which is very close to  $\underline{\alpha}$  since  $\sigma_-(\sigma_+)^{-2} = (\sigma_+)^{-1} \cdot \frac{\sigma_-}{\sigma_+}$  and  $\frac{\sigma_-}{\sigma_+} \leq 1$ .

**Remark 24.** If one wants to obtain estimates without the bounding factors  $\sigma_{\pm}$  and  $D_{\pm}$ , a standard path is to imbed the parameters  $D$  and  $\sigma$  into the definition of the norms, like it is done in Section 4. Namely, one chooses the norms:

$$\begin{aligned} \|v\|_{\sigma} &= \left( (\sigma v, v)_{L^2(\Omega)} \right)^{1/2}, \\ \|\mathbf{q}\|_{D^{-1}, \sigma^{-1} \operatorname{div}} &= \left( (D^{-1} \mathbf{q}, \mathbf{q})_{L^2(\Omega)} + (\sigma^{-1} \operatorname{div} \mathbf{q}, \operatorname{div} \mathbf{q})_{L^2(\Omega)} \right)^{1/2}, \\ \|(\mathbf{q}, v)\|_V &= \left( \|v\|_{\sigma}^2 + \|\mathbf{q}\|_{D^{-1}, \sigma^{-1} \operatorname{div}}^2 \right)^{1/2}. \end{aligned}$$

On the one hand, all norms are “fixed” once the parameters are given. On the other hand, one can easily check that the stability constant is now independent of the bounding factors, by using the same mapping  $\mathbb{T}$  as before.

## 7. T-coercivity at the discrete level

Previously, we demonstrated the robustness and the flexibility of the T-coercivity approach to study mixed problems at the continuous level. In this section, we are going to see how T-coercivity also enables us to provide a stable discretization of such problems with mixed finite elements. Let us recall the simple results below [19, 20].

**Definition 25.** Let  $W$  be a Hilbert space,  $\mathcal{A}(\cdot, \cdot)$  be a continuous bilinear form over  $W \times W$  and  $(W_h)_h$  be conforming approximations of  $W$ . We say that  $\mathcal{A}$  is uniformly  $\mathbb{T}_h$ -coercive if

$$\exists \alpha^*, \beta^* > 0, \forall h > 0, \exists \mathbb{T}_h \in \mathcal{L}(W_h), \quad |\mathcal{A}(u_h, \mathbb{T}_h u_h)| \geq \alpha^* \|u_h\|_W^2, \quad \forall u_h \in W_h, \quad \text{and} \quad \|\mathbb{T}_h\| \leq \beta^*.$$

**Proposition 26.** Let  $W$  be a Hilbert space,  $f$  be an element of  $W'$ ,  $\mathcal{A}(\cdot, \cdot)$  be a continuous bilinear form over  $W \times W$  and  $(W_h)_h$  be conforming approximations of  $W$ . Denote by  $\mathbf{A}_h \in \mathcal{L}(W_h, W_h')$  the discrete operator associated to  $\mathcal{A}|_{W_h}$ . The problem

$$\begin{cases} \text{Find } u_h \in W_h & \text{such that} \\ \forall v_h \in W_h, & \mathcal{A}(u_h, v_h) = \langle f, v_h \rangle \end{cases}$$

is well-posed and  $(\mathbf{A}_h^{-1})_h$  is uniformly bounded if and only if  $\mathcal{A}$  is uniformly  $\mathbb{T}_h$ -coercive. In that case, denoting by  $C_{\mathcal{A}}$  the continuity constant of the bilinear form  $\mathcal{A}$ , it holds that

$$\|u - u_h\|_W \leq C \inf_{v_h \in W_h} \|u - v_h\|_W, \quad (78)$$

with  $C = 1 + \frac{C_{\mathcal{A}} \beta^*}{\alpha^*}$  independent of  $h$ .

**Remark 27.** Proposition 26 can be extended to the case where the discrete forms  $\mathcal{A}_h$  and  $f_h$  differs from the continuous forms  $\mathcal{A}$  and  $f$ . In that case, Céa’s lemma (78) becomes

$$\|u - u_h\|_W \leq C \inf_{v_h \in W_h} \left( \|u - v_h\|_W + \operatorname{Cons}_{f,h} + \operatorname{Cons}_{\mathcal{A},h}(v_h) \right),$$

with

$$\operatorname{Cons}_{f,h} = \sup_{v_h \in W_h \setminus \{0\}} \frac{|\langle f - f_h, v_h \rangle|}{\|v_h\|_W} \quad \text{and} \quad \operatorname{Cons}_{\mathcal{A},h}(v_h) = \sup_{w_h \in W_h \setminus \{0\}} \frac{|(\mathcal{A} - \mathcal{A}_h)(v_h, w_h)|}{\|w_h\|_W}, \quad \forall v_h \in W_h.$$

As before, we start with the leading example of Stokes problem.

### 7.1. Stokes problem

For a given  $h$ , the natural discretization of Problem (5) reads:

$$\begin{cases} \text{Find } (\mathbf{u}_h, p_h) \in \mathbf{V}_h \times Q_h & \text{such that} \\ \forall (\mathbf{v}_h, q_h) \in \mathbf{V}_h \times Q_h, & \mathcal{A}((\mathbf{u}_h, p_h), (\mathbf{v}_h, q_h)) = \langle \mathbf{f}, \mathbf{v}_h \rangle, \end{cases} \quad (79)$$

where  $\mathbf{V}_h \subset \mathbf{H}_0^1(\Omega)$  and  $Q_h \subset L_0^2(\Omega)$  are two finite dimensional spaces constituting a *conforming* approximation of  $\mathbf{H}_0^1(\Omega) \times L_0^2(\Omega)$ .

From Proposition 26, we know that Problem (79) is well-posed if and only if  $\mathcal{A}$  is uniformly  $T_h$ -coercive. To build a suitable mapping  $T_h \in \mathcal{L}(\mathbf{V}_h \times Q_h)$ , a natural idea is to reproduce the continuous mapping from the proof of Theorem 3

$$\begin{aligned} T : \mathbf{H}_0^1(\Omega) \times L_0^2(\Omega) &\longrightarrow \mathbf{H}_0^1(\Omega) \times L_0^2(\Omega) \\ (\mathbf{u}, p) &\longmapsto (\lambda \mathbf{u} + \mathbf{v}_p, -\lambda p) \end{aligned}$$

at the discrete level. The operator  $T$  above depends on the divergence lifting  $\mathbf{v}_p \in \mathbf{H}_0^1(\Omega)$  of the pressure  $p \in L_0^2(\Omega)$  defined by, see (8)-(9),

$$-\operatorname{div} \mathbf{v}_p = p \quad \text{and} \quad \|\nabla \mathbf{v}_p\| \leq C_{\operatorname{div}} \|p\|.$$

To obtain a similar lifting in the discrete setting, we consider the continuous lifting of the discrete pressure  $p_h \in Q_h \subset L_0^2(\Omega)$ , namely  $\mathbf{v}_{p_h} \in \mathbf{H}_0^1(\Omega)$  such that

$$-\operatorname{div} \mathbf{v}_{p_h} = p_h \quad \text{and} \quad \|\nabla \mathbf{v}_{p_h}\| \leq C_{\operatorname{div}} \|p_h\|. \quad (80)$$

This lifting  $\mathbf{v}_{p_h}$  does not necessarily belong to the discrete space  $\mathbf{V}_h \subset \mathbf{H}_0^1(\Omega)$ , so we need an operator  $\Pi_h : \mathbf{H}_0^1(\Omega) \longrightarrow \mathbf{V}_h$  to project it on  $\mathbf{V}_h$ . Therefore, we consider a discrete mapping of the form

$$\begin{aligned} T_h : \mathbf{V}_h \times Q_h &\longrightarrow \mathbf{V}_h \times Q_h \\ (\mathbf{u}_h, p_h) &\longmapsto (\lambda \mathbf{u}_h + \Pi_h(\mathbf{v}_{p_h}), -\lambda p_h). \end{aligned} \quad (81)$$

Now, let us precise under which conditions the bilinear form  $\mathcal{A}$  is uniformly  $T_h$ -coercive by mimicking the proof of Theorem 3. We compute

$$\mathcal{A}((\mathbf{u}_h, p_h), T_h(\mathbf{u}_h, p_h)) = \nu \lambda \|\nabla \mathbf{u}_h\|^2 + \nu \int_{\Omega} \nabla \mathbf{u}_h : \nabla (\Pi_h(\mathbf{v}_{p_h})) \, dx - \int_{\Omega} p_h \operatorname{div} (\Pi_h \mathbf{v}_{p_h}) \, dx.$$

In order to get a term of the form  $\|p_h\|^2$ , we assume that

$$\int_{\Omega} p_h \operatorname{div} (\Pi_h \mathbf{v}_{p_h}) \, dx = \int_{\Omega} p_h \operatorname{div} \mathbf{v}_{p_h} \, dx, \quad (82)$$

so that

$$\mathcal{A}((\mathbf{u}_h, p_h), T_h(\mathbf{u}_h, p_h)) = \nu \lambda \|\nabla \mathbf{u}_h\|^2 + \nu \int_{\Omega} \nabla \mathbf{u}_h : \nabla (\Pi_h(\mathbf{v}_{p_h})) \, dx + \|p_h\|^2$$

in view of (80). Then, for any  $\eta > 0$ , Young inequality yields

$$\begin{aligned} \int_{\Omega} \nabla \mathbf{u}_h : \nabla (\Pi_h(\mathbf{v}_{p_h})) \, dx &\geq -\frac{\eta}{2} \|\nabla \mathbf{u}_h\|^2 - \frac{1}{2\eta} \|\nabla (\Pi_h(\mathbf{v}_{p_h}))\|^2 \\ &\geq -\frac{\eta}{2} \|\nabla \mathbf{u}_h\|^2 - \frac{C_{\operatorname{div}}^2 C_{\pi}^2}{2\eta} \|p_h\|^2 \end{aligned}$$

provided that there exists a constant  $C_{\pi} > 0$ , independent of  $h$  and of  $p_h$ , such that

$$\|\nabla (\Pi_h(\mathbf{v}_{p_h}))\| \leq C_{\pi} \|\nabla \mathbf{v}_{p_h}\|. \quad (83)$$

Hence, it holds that

$$\mathcal{A}((\mathbf{u}_h, p_h), T_h(\mathbf{u}_h, p_h)) \geq \nu \left( \lambda - \frac{\eta}{2} \right) \|\nabla \mathbf{u}_h\|^2 + \left( 1 - \frac{\nu C_{\operatorname{div}}^2 C_{\pi}^2}{2\eta} \right) \|p_h\|^2.$$

Setting  $\eta = \lambda = \nu C_{\text{div}}^2 C_{\pi}^2$ , we obtain

$$\begin{aligned} \mathcal{A}((\mathbf{u}_h, \mathbf{p}_h), \mathbf{T}_h(\mathbf{v}_h, \mathbf{p}_h)) &\geq \frac{\nu^2 C_{\text{div}}^2 C_{\pi}^2}{2} \|\nabla \mathbf{u}_h\|^2 + \frac{1}{2} \|\mathbf{p}_h\|^2 \\ &\geq \frac{1}{2} \min(\nu^2 C_{\text{div}}^2 C_{\pi}^2, 1) \|(\mathbf{u}_h, \mathbf{p}_h)\|_{\mathbf{H}_0^1(\Omega) \times L_0^2(\Omega)}. \end{aligned} \quad (84)$$

Moreover, taking into account (83) and mimicking the continuous case (see (11)), we have

$$\|\mathbf{T}_h\| \leq \max\left(\sqrt{2}\nu C_{\text{div}}^2 C_{\pi}^2, C_{\text{div}} C_{\pi} (2 + \nu^2 C_{\text{div}}^2 C_{\pi}^2)^{1/2}\right). \quad (85)$$

So, with the help of the operator  $\Pi_h : \mathbf{H}_0^1(\Omega) \rightarrow \mathbf{V}_h$ , we have proven the following result.

**Theorem 28.** *If there exist a family of operators  $(\Pi_h)_h$  and a constant  $C_{\pi} > 0$  such that, for all  $h$ ,*

$$\int_{\Omega} q_h \operatorname{div}(\Pi_h \mathbf{v}) \, dx = \int_{\Omega} q_h \operatorname{div} \mathbf{v} \, dx, \quad \forall \mathbf{v} \in \mathbf{H}_0^1(\Omega), \forall q_h \in Q_h, \quad (86)$$

$$\|\nabla(\Pi_h(\mathbf{v}))\| \leq C_{\pi} \|\nabla \mathbf{v}\|, \quad \forall \mathbf{v} \in \mathbf{H}_0^1(\Omega), \quad (87)$$

then Problem (79) is well-posed for all  $h$  and

$$\|(u - u_h, \mathbf{p} - \mathbf{p}_h)\|_{\mathbf{H}_0^1(\Omega) \times L_0^2(\Omega)} \leq C \inf_{(v_h, q_h) \in \mathbf{V}_h \times Q_h} \|(u - v_h, \mathbf{p} - q_h)\|_{\mathbf{H}_0^1(\Omega) \times L_0^2(\Omega)}, \quad (88)$$

with

$$C = 1 + \frac{2 \max(\nu, 2) \max\left(\sqrt{2}\nu C_{\text{div}}^2 C_{\pi}^2, C_{\text{div}} C_{\pi} (2 + \nu^2 C_{\text{div}}^2 C_{\pi}^2)^{1/2}\right)}{\min(\nu^2 C_{\text{div}}^2 C_{\pi}^2, 1)}.$$

**Proof.** The previous reasoning shows that the bilinear form  $\mathcal{A}$  is uniformly  $\mathbf{T}_h$ -coercive for the mapping

$$\begin{aligned} \mathbf{T}_h : \mathbf{V}_h \times Q_h &\longrightarrow \mathbf{V}_h \times Q_h \\ (\mathbf{u}_h, \mathbf{p}_h) &\longmapsto (\nu C_{\text{div}}^2 C_{\pi}^2 \mathbf{u}_h + \Pi_h(\mathbf{v}_{\mathbf{p}_h}), -\nu C_{\text{div}}^2 C_{\pi}^2 \mathbf{p}_h) \end{aligned}$$

as long as the two conditions (82) and (83) are fulfilled for all  $\mathbf{p}_h \in Q_h$ , which is the case if (86) and (87) hold true. The stability estimate (88) then follows by using (84) and (85) in (78).  $\square$

The conditions (86) and (87) correspond exactly to the assumptions of an abstract result known as Fortin's lemma [26]. Above, the T-coercivity approach allowed us to recover these two conditions in a fully constructive way. Moreover, we recall that, since the form  $b$  fulfills an inf-sup condition (3), those conditions (86)-(87) are equivalent to the so-called *uniform discrete inf-sup condition*

$$\exists \underline{\beta}' > 0, \quad \forall h, \quad \inf_{q_h \in Q_h \setminus \{0\}} \sup_{\mathbf{v}_h \in \mathbf{V}_h \setminus \{0\}} \frac{\int_{\Omega} q_h \operatorname{div} \mathbf{v}_h \, dx}{\|\nabla \mathbf{v}_h\| \|q_h\|} \geq \underline{\beta}',$$

see for instance [27, Lemma II.1.1].

Finally, we recall that, provided there is a basic approximability property (*i.e.* any element of  $\mathbf{V} \times Q$  can be approximated by a sequence of elements of  $(\mathbf{V}_h \times Q_h)_h$ ), the convergence of the discrete solutions to the exact one is a consequence of (88).

## 7.2. Approximation of saddle-point problems

We now derive a *conforming* approximation of the abstract problem (15), starting from the variational expressions (16) or (17), the latter with the form

$$\mathcal{A}((u, \mathbf{p}), (v, \mathbf{q})) = a(u, v) + b(v, \mathbf{p}) + b(u, \mathbf{q}).$$

So, let  $(V_h)_h$ , resp.  $(Q_h)_h$ , be two families of finite dimensional subspaces of  $V$ , resp.  $Q$ . Starting from (16), the discrete variational formulation writes

$$\begin{cases} \text{Find } (u_h, p_h) \in V_h \times Q_h \text{ such that} \\ \forall v_h \in V_h, \quad a(u_h, v_h) + b(v_h, p_h) = \langle f, v_h \rangle_{V',V} \\ \forall q_h \in Q_h, \quad b(u_h, q_h) = \langle g, q_h \rangle_{Q',Q}. \end{cases}$$

while, starting from (17), the *all-in-one* discrete variational formulation writes

$$\begin{cases} \text{Find } (u_h, p_h) \in V_h \times Q_h \text{ such that} \\ \forall (v_h, q_h) \in V_h \times Q_h, \quad \mathcal{A}((u_h, p_h), (v_h, q_h)) = \langle f, v_h \rangle_{V',V} + \langle g, q_h \rangle_{Q',Q}. \end{cases}$$

In abstract form, the *uniform discrete inf-sup condition* writes

$$\exists \underline{\beta}' > 0, \quad \forall h, \quad \inf_{q_h \in Q_h \setminus \{0\}} \sup_{v_h \in V_h \setminus \{0\}} \frac{b(v_h, q_h)}{\|v_h\|_V \|q_h\|_Q} \geq \underline{\beta}'. \quad (89)$$

We suppose that the discrete version of the operator  $B$  is the restriction of  $B$  to  $V_h$ , namely

$$B(V_h) \subset Q'_h. \quad (90)$$

Introducing the discrete operators  $B_h : V_h \rightarrow Q_h$  such that for all  $h$ ,

$$(B_h v_h, q_h)_Q = b(v_h, q_h), \quad \forall (v_h, q_h) \in V_h \times Q_h,$$

the straightforward discrete counterpart of Theorem 5 is

**Theorem 29.** *The following three statements are equivalent:*

- (i) *There exists  $\underline{\beta}' > 0$  such that  $b(\cdot, \cdot)$  fulfills the uniform discrete inf-sup condition (89).*
- (ii) *For all  $h$ ,  $B_h : (\text{Ker } B_h)^\perp \rightarrow Q_h$  is an isomorphism, and*

$$\|B_h v_h\|_Q \geq \underline{\beta}' \|v_h\|_V, \quad \forall v_h \in (\text{Ker } B_h)^\perp. \quad (91)$$

- (iii) *For all  $h$ , there exists an isomorphic operator  $L_{B,h} : Q_h \rightarrow (\text{Ker } B_h)^\perp$  such that*

$$B_h(L_{B,h} q_h) = q_h \quad \text{and} \quad \|q_h\|_Q \geq \underline{\beta}' \|L_{B,h} q_h\|_V, \quad \forall q_h \in Q_h. \quad (92)$$

**Remark 30.** Obviously, this result also holds if the value of the constant in the discrete inf-sup condition depends on  $h$ , i.e. for each  $h$  it holds for some  $\underline{\beta}'(h) > 0$ , with  $\lim_{h \rightarrow 0} \underline{\beta}'(h) = 0$ . In this case however, getting error estimates can be more intricate.

As mentioned above for the Stokes system, one has the Fortin lemma (cf. [27, Lemma II.1.1]).

**Theorem 31.** *Assume that the form  $b$  fulfills an inf-sup condition (18). The uniform discrete inf-sup condition (89) holds if, and only if, there exist a family of operators  $(\Pi_h)_h$ , with  $\Pi_h : V \rightarrow V_h$ , and a constant  $C_\pi > 0$  such that, for all  $h$ ,*

$$\begin{aligned} b(\Pi_h v, q_h) &= b(v, q_h), \quad \forall v \in V, \forall q_h \in Q_h, \\ \sup_h \|\Pi_h\| &\leq C_\pi. \end{aligned} \quad (93)$$

Let us now proceed with the derivation of conditions to ensure that the form  $\mathcal{A}$  is uniformly  $T_h$ -coercive. As a general rule, the proofs of the results follow very closely the proofs that were given in the exact case. The straightforwardness of the procedure when going from the continuous to the discrete level is one of the main features of the T-coercivity approach. We give next the discrete counterparts of Theorems 6 and 8.

**Theorem 32.** *Assume that the form  $a$  is symmetric and positive, that there exists a constant  $\alpha' > 0$  such that*

$$a(u_h, u_h) \geq \alpha' \|u_h\|_V^2, \quad \forall u_h \in V_h, \quad (94)$$

*and that the uniform discrete inf-sup condition (89) on the form  $b$  holds true. Then the form  $\mathcal{A}$  is uniformly  $T_h$ -coercive.*

The property (94) is sometimes called the the uniform discrete coercivity property.

**Proof.** Let  $h$  be given. We introduce the mapping

$$\begin{aligned} T_h : V_h \times Q_h &\longrightarrow V_h \times Q_h \\ (u_h, p_h) &\longmapsto (\lambda u_h + L_{B,h} p_h, -\lambda p_h). \end{aligned}$$

We then compute

$$\begin{aligned} \mathcal{A}((u_h, p_h), T_h(u_h, p_h)) &= a(u_h, \lambda u_h) + a(u_h, L_{B,h} p_h) + b(\lambda u_h, p_h) + b(L_{B,h} p_h, p_h) - b(u_h, \lambda p_h) \\ &= \lambda a(u_h, u_h) + a(u_h, L_{B,h} p_h) + \|p_h\|_Q^2, \text{ according to (92)-left.} \end{aligned}$$

Because the form  $a$  is symmetric and positive, we can apply Young's inequality: for any  $\eta > 0$ ,

$$a(u_h, L_{B,h} p_h) \geq -\frac{\eta}{2} a(u_h, u_h) - \frac{1}{2\eta} a(L_{B,h} p_h, L_{B,h} p_h).$$

According now to (92)-right, we find

$$a(L_{B,h} p_h, L_{B,h} p_h) \leq C_a \|L_{B,h} p_h\|_V^2 \leq C_a (\underline{\beta}')^{-2} \|p_h\|_Q^2.$$

Using assumption (94), if  $\lambda - \frac{\eta}{2} > 0$ , it follows that

$$\mathcal{A}((u_h, p_h), T_h(u_h, p_h)) \geq \alpha' \left( \lambda - \frac{\eta}{2} \right) \|u_h\|_V^2 + \left( 1 - \frac{C_a (\underline{\beta}')^{-2}}{2\eta} \right) \|p_h\|_Q^2.$$

Setting  $\eta = \lambda = C_a (\underline{\beta}')^{-2}$  as in the exact case, we infer that

$$\mathcal{A}((u_h, p_h), T_h(u_h, p_h)) \geq \frac{1}{2} \min(\alpha' C_a (\underline{\beta}')^{-2}, 1) \| (u_h, p_h) \|_{V \times Q}^2$$

which proves that  $\mathcal{A}$  is  $T_h$ -coercive, with a T-coercivity constant  $\frac{1}{2} \min(\alpha' C_a (\underline{\beta}')^{-2}, 1) > 0$  that is independent of  $h$ .

Since  $T_h(u_h, p_h) = (C_a (\underline{\beta}')^{-2} u_h + L_{B,h} p_h, -C_a (\underline{\beta}')^{-2} p_h)$ , one finds that

$$\begin{aligned} \|T_h(u_h, p_h)\|_{V \times Q}^2 &\leq 2(C_a (\underline{\beta}')^{-2})^2 \|u_h\|_V^2 + 2\|L_{B,h} p_h\|_V^2 + (C_a (\underline{\beta}')^{-2})^2 \|p_h\|_Q^2 \\ &\leq 2(C_a (\underline{\beta}')^{-2})^2 \|u_h\|_V^2 + (2(\underline{\beta}')^{-2} + (C_a (\underline{\beta}')^{-2})^2) \|p_h\|_Q^2, \end{aligned}$$

where the last inequality follows from (92)-right. The bound is valid for all  $h$ , which yields

$$\sup_h \|T_h\| \leq \max\left(\sqrt{2} C_a (\underline{\beta}')^{-2}, \beta(2 + C_a^2 (\underline{\beta}')^{-2})^{1/2}\right),$$

so the form  $\mathcal{A}$  is uniformly  $T_h$ -coercive.  $\square$

**Remark 33.** As for the Stokes problem, the discrete right-inverse  $L_{B,h}$  is connected to the Fortin operator  $\Pi_h$ . As a matter of fact, if there exists a family of discrete projectors  $(\Pi_h)_h$  verifying (93), the operator defined by  $L_{B,h} = \Pi_h(L_B)$  satisfies (92) with  $\underline{\beta}' = (C_\pi \beta)^{-1}$  since for all  $q_h \in Q_h$

$$\|\Pi_h(L_B q_h)\|_V \leq C_\pi \|L_B q_h\|_V \leq C_\pi \beta \|q_h\|_Q,$$

according to (21). As a consequence, to perform stability estimates at the discrete level using  $T_h$ -coercivity, one has only to replace  $\beta$  by  $C_\pi \beta$  in the computations done at the continuous level.

**Theorem 34.** Assume that the form  $a$  is symmetric and positive, that there exists a constant  $\alpha'_0 > 0$  such that

$$a(u_{0,h}, u_{0,h}) \geq \alpha'_0 \|u_{0,h}\|_V^2, \quad \forall u_{0,h} \in \text{Ker } B_h, \quad (95)$$

and that the uniform discrete inf-sup condition (89) on the form  $b$  holds true.

Then the form  $\mathcal{A}$  is uniformly  $T_h$ -coercive.

The property (95) is sometimes called the uniform discrete coercivity property on the kernels.

**Proof.** Let  $h$  be given. We consider the mapping

$$\begin{aligned} T_h : V_h \times Q_h &\longrightarrow V_h \times Q_h \\ (u_h, p_h) &\longmapsto (\lambda u_h + L_{B,h} p_h, -\lambda p_h + \lambda \mu B_h u_h). \end{aligned}$$

As in the proof of Theorem 8, we can compute

$$\mathcal{A}((u_h, p_h), T_h(u_h, p_h)) = \lambda a(u_h, u_h) + a(u_h, L_{B,h} p_h) + \|p_h\|_Q^2 + \lambda \mu \|B_h u_h\|_Q^2$$

because  $b(L_{B,h} p_h, p_h) = \|p_h\|_Q^2$ . Since the form  $a$  is symmetric and positive, one may use Young's inequality. By proceeding as in the proof of Theorem 32 and after setting  $\lambda = C_a(\underline{\beta}')^{-2}$ , we find that

$$\lambda a(u_h, u_h) + a(u_h, L_{B,h} p_h) + \|p_h\|_Q^2 \geq \frac{1}{2} C_a(\underline{\beta}')^{-2} a(u_h, u_h) + \frac{1}{2} \|p_h\|_Q^2,$$

and

$$\mathcal{A}((u_h, p_h), T_h(u_h, p_h)) \geq \frac{1}{2} C_a(\underline{\beta}')^{-2} (a(u_h, u_h) + 2\mu \|B_h u_h\|_Q^2) + \frac{1}{2} \|p_h\|_Q^2.$$

Then, we use the decomposition  $u_h = u_{0,h} + \bar{u}_h$  with  $u_{0,h} \in \text{Ker} B_h$  and  $\bar{u}_h \in (\text{Ker} B_h)^\perp$ . As before, Young's inequality yields

$$a(u_h, u_h) \geq (1 - \theta) a(u_{0,h}, u_{0,h}) + \left( C_a - \frac{C_a}{\theta} \right) \|\bar{u}_h\|_V^2$$

for all  $0 < \theta < 1$ . Moreover,  $\|B_h u_h\|_Q^2 = \|B_h \bar{u}_h\|_Q^2 \geq (\underline{\beta}')^2 \|\bar{u}_h\|_V^2$  according to (91), so that

$$a(u_h, u_h) + 2\mu \|B_h u_h\|_Q^2 \geq (1 - \theta) a(u_{0,h}, u_{0,h}) + \left( C_a - \frac{C_a}{\theta} + 2\mu (\underline{\beta}')^2 \right) \|\bar{u}_h\|_V^2.$$

Choosing  $\theta = \frac{1}{2}$  and  $\mu = \frac{3}{4} C_a (\underline{\beta}')^{-2}$ , it holds that

$$\begin{aligned} a(u_h, u_h) + 2\mu \|B_h u_h\|_Q^2 &\geq \frac{1}{2} a(u_{0,h}, u_{0,h}) + \frac{C_a}{2} \|\bar{u}_h\|_V^2 \\ &\geq \frac{\alpha'_0}{2} \|u_{0,h}\|_V^2 + \frac{\alpha'_0}{2} \|\bar{u}_h\|_V^2 = \frac{\alpha'_0}{2} \|u_h\|_V^2, \end{aligned}$$

where we used assumption (95) and  $C_a \geq \alpha'_0$  on the second line.

Finally, we conclude that

$$\mathcal{A}((u_h, p_h), T_h(u_h, p_h)) \geq \frac{1}{4} \alpha'_0 C_a (\underline{\beta}')^{-2} \|u_h\|_V^2 + \frac{1}{2} \|p_h\|_Q^2,$$

which yields that  $\mathcal{A}$  is  $T_h$ -coercive, with a T-coercivity constant  $\min(\frac{1}{4} \alpha'_0 C_a (\underline{\beta}')^{-2}, \frac{1}{2}) > 0$  that is independent of  $h$ .

From the above, we have  $T_h(u_h, p_h) = (C_a (\underline{\beta}')^{-2} u_h + L_{B,h} p_h, -C_a (\underline{\beta}')^{-2} p_h + \frac{3}{4} (C_a (\underline{\beta}')^{-2})^2 B_h u_h)$ , and, noting that  $\|B_h u_h\|_Q \leq C_b \|u_h\|_V$ , one concludes that

$$\sup_h \|T_h\| \leq \infty,$$

so the form  $\mathcal{A}$  is uniformly  $T_h$ -coercive.  $\square$

**Remark 35.** Note that replacing  $Bu$  by  $B_h u_h$  when going from the continuous operator  $T$  to the discrete operator  $T_h$  is possible because we assumed that  $B_h$  is the restriction of  $B$  to  $V_h$ , see (90). If (90) does not hold, one introduces  $\Phi_h : Q \rightarrow Q_h$  defined by

$$b(v_h, \Phi_h q) = b(v_h, q), \quad \forall q \in Q, \forall v_h \in V_h,$$

like in [4, Proposition 5.1.2]. Then, one has to replace  $B_h u_h$  by  $\Phi_h(B_h u_h)$  in the previous proof, *i.e.*

$$\begin{aligned} T_h : V_h \times Q_h &\longrightarrow V_h \times Q_h \\ (u_h, p_h) &\longmapsto \left( C_a (\underline{\beta}')^{-2} u_h + L_{B,h} p_h, -C_a (\underline{\beta}')^{-2} p_h + \frac{3}{4} (C_a (\underline{\beta}')^{-2})^2 \Phi_h(B_h u_h) \right). \end{aligned}$$

Again, provided a basic approximability property holds, that is, any element of  $V \times Q$  can be approximated by a sequence of elements of  $(V_h \times Q_h)_h$ , convergence will follow under the assumptions of Theorem 32 or Theorem 34.

### 7.3. Approximation of augmented saddle-point problems

We now approximate the abstract problem (35), starting from the variational expression (36), with the form

$$\mathcal{A}_c((u, p), (v, q)) = a(u, v) + b(v, p) + b(u, q) - c(p, q),$$

where  $c(\cdot, \cdot)$  is a form defined on  $Q \times Q$  that fulfills (34); in particular,  $c(\cdot, \cdot)$  is *positive*. So, let again  $(V_h)_h$ , resp.  $(Q_h)_h$ , be two families of finite dimensional subspaces of  $V$ , resp.  $Q$ . The discrete variational formulation writes

$$\begin{cases} \text{Find } (u_h, p_h) \in V_h \times Q_h \text{ such that} \\ \forall (v_h, q_h) \in V_h \times Q_h, \quad \mathcal{A}_c((u_h, p_h), (v_h, q_h)) = \langle f, v_h \rangle_{V', V} + \langle g, q_h \rangle_{Q', Q}, \end{cases}$$

To ensure that the form  $\mathcal{A}_c$  is uniformly  $T_h$ -coercive, the proofs once more follow very closely those that were given in the exact case. We give next the discrete counterparts of Theorems 10, 12 and 18.

**Theorem 36.** *Assume that the form  $a$  is symmetric, positive, fulfills the uniform discrete coercivity property (94), and that the uniform discrete inf-sup condition (89) on the form  $b$  holds true. Then the form  $\mathcal{A}_c$  is uniformly  $T_h$ -coercive.*

**Theorem 37.** *Assume that the form  $a$  is symmetric and positive, fulfills the uniform discrete coercivity property on the kernels (95), and that the uniform discrete inf-sup condition (89) on the form  $b$  holds true. If moreover the form  $c$  is like in (39), where  $\varepsilon$  is small enough, namely*

$$\varepsilon \leq \frac{1}{2C_a(C_\pi\beta)^4 C_b^2} \left(2 - \frac{\alpha_0}{C_a}\right),$$

*then the form  $\mathcal{A}_c$  is uniformly  $T_h$ -coercive.*

**Theorem 38.** *Assume that (51) holds true and that the bilinear forms  $a$  and  $c$  are both symmetric and positive. If there exists a constant  $\alpha'_B > 0$  such that*

$$a(u_h, u_h) + \frac{\gamma}{2C_c^2} \|B_h u_h\|_Q^2 \geq \alpha'_B \|u_h\|_V^2, \quad \forall u_h \in V_h, \quad (96)$$

*then the form  $\mathcal{A}_c$  is uniformly  $T_h$ -coercive.*

As before, provided a basic approximability property holds, that is, any element of  $V \times Q$  can be approximated by a sequence of elements of  $(V_h \times Q_h)_h$ , convergence will follow under the assumptions of Theorem 36, Theorem 37 or Theorem 38.

### 7.4. Applications

Let us briefly see how the T-coercivity approach can be used to discretize the mixed problems, that is for Stokes, electromagnetism, nearly-incompressible elasticity and finally neutron diffusion. For each problem, we propose one or several possibilities. Note that, since there is a vast literature on this topic, there is no need to devise new approximation techniques. On the contrary, the simple framework of the T-coercivity approach provides elementary guidelines to help us choose among existing techniques.

In each case, the first step is to choose a *conforming* finite element discretization adapted to the space  $V$  under consideration. We assume for simplicity that  $\Omega$  is a polyhedron for  $d = 3$ , or a polygon for  $d = 2$ , so one can use meshes made of simplices for the discretization by finite elements. For  $k \geq 1$ ,  $\mathcal{P}_k$  stands for the Lagrange finite elements of order  $k$ .

For Stokes and elasticity, we note that the space  $\mathbf{H}_0^1(\Omega)$  may be approximated using  $(\mathcal{P}_k)^d$  finite elements with  $k \geq 2$ . For electromagnetism, we have to deal with the space  $\mathbf{H}_0(\mathbf{curl}; \Omega)$ , which can be discretized using the (first-kind) Nédélec finite elements of order  $k \geq 1$ , denoted by  $\mathcal{N}_k$ . Lastly, for neutron diffusion, we have to deal with the space  $\mathbf{H}(\text{div}; \Omega)$ , discretized with the help of the Raviart-Thomas elements of order  $k \geq 0$ , denoted by  $\mathcal{RT}_k$ . We refer to [4] for details.

The next step is to choose the *conforming* finite element discretization in the space  $Q$  in such a way that convergence of the discrete solutions to the exact one is guaranteed.

First, for Stokes and elasticity, and for  $k = 2$ , setting  $Q_h = \mathcal{P}_1$  leads to Fortin operators  $\Pi_h^\mathcal{P} : \mathbf{H}_0^1(\Omega) \rightarrow (\mathcal{P}_k)^d$  satisfying (86)-(87), or the abstract counterpart (93): the pair  $((\mathcal{P}_2)^d, \mathcal{P}_1)$  is called the Taylor-Hood finite element. Then, for electromagnetism and for  $k \geq 1$ , setting  $Q_h = \mathcal{P}_k$  leads to Fortin operators  $\Pi_h^\mathcal{N} : \mathbf{H}_0(\mathbf{curl}; \Omega) \rightarrow \mathcal{N}_k$  satisfying (93). Last, for neutron diffusion and for  $k \geq 1$ , one only needs to select  $Q_h$  in such a way that (90) is fulfilled, namely  $\text{div}(\mathcal{RT}_k) \subset Q_h'$ . To do so, we set  $Q_h = \mathcal{P}_k^{pw}$ , for  $k \geq 0$ , where the superscript  $^{pw}$  stands for piecewise Lagrange finite elements of order  $k$ . We again refer to [4] for details and possible extensions, such as the generalized Taylor-Hood elements for Stokes or elasticity, for  $k \geq 3$ .

On the other hand, as demonstrated earlier, those properties may be also recovered straightforwardly thanks to the T-coercivity approach. As a matter of fact, we can build the discrete operators  $T_h$  similarly as in the continuous case but changing the constant  $\beta$  to take into account (when applicable) the influence of the projection operators, see Remark 33.

For Stokes, we refer to (81).

For electromagnetism, we infer from Theorem 34 that the bilinear form defined in (59) is uniformly  $T_h$ -coercive for the mapping

$$T_h : \mathcal{N}_k \times \mathcal{P}_k \longrightarrow \mathcal{N}_k \times \mathcal{P}_k$$

$$(\mathbf{E}_h, \tilde{p}_h) \longmapsto \left( (C_{\pi, \mathcal{N}})^2 \mathbf{E}_h + \Pi_h^\mathcal{N}(\nabla \tilde{p}_h), -(C_{\pi, \mathcal{N}})^2 \tilde{p}_h + \frac{3}{4} (C_{\pi, \mathcal{N}})^4 \phi_{E_h} \right),$$

where  $\phi_{E_h} \in Q_h$  satisfies the discrete counterpart of (65), namely

$$(\underline{\varepsilon} \nabla \phi_{E_h}, \nabla q_h)_{L^2(\Omega)} = (\underline{\varepsilon} \mathbf{E}_h, \nabla q_h)_{L^2(\Omega)}, \quad \forall q_h \in \mathcal{P}_k.$$

For nearly-incompressible elasticity, the bilinear form defined in (70) is uniformly  $T_h$ -coercive for the mapping

$$T_h : \mathcal{P}_2 \times \mathcal{P}_1 \longrightarrow \mathcal{P}_2 \times \mathcal{P}_1$$

$$(\mathbf{u}_h, p_h) \longmapsto \left( 2\mu(C_{\pi, \mathcal{P}} C_{\text{div}})^2 \mathbf{u}_h + \Pi_h^\mathcal{P}(\mathbf{v}_{-p_h}), -2\mu(C_{\pi, \mathcal{P}} C_{\text{div}})^2 p_h \right),$$

according to Theorem 36 and (71).

Finally, for neutron diffusion, assuming for simplicity that  $\sigma$  restricted to any simplex is constant, we introduce the discrete mapping

$$T_h : \mathcal{RT}_k \times \mathcal{P}_k^{pw} \longrightarrow \mathcal{RT}_k \times \mathcal{P}_k^{pw}$$

$$(\mathbf{p}_h, u_h) \longmapsto \left( \mathbf{p}_h, \frac{1}{2}(-u_h + \sigma^{-1} \text{div} \mathbf{p}_h) \right),$$

and the property  $\text{div}(\mathcal{RT}_k) \subset \mathcal{P}_k^{pw}$  guarantees the uniform  $T_h$ -coercivity of the bilinear form (76) in virtue of Theorem 38.



All basic approximability properties are established in [4], which guarantees convergence in each case.

## 8. Conclusion and perspectives

We have demonstrated the flexibility of the T-coercivity approach, here applied to classical linear mixed problems, both for the theoretical study of the problems and for their numerical approximation by finite elements. Let us mention some possible extensions, such as nonconforming discretization methods for Stokes [32], multigroup diffusion [28] or DDM for diffusion [24].

It is our belief that numerous applications can be studied with the T-coercivity approach, both theoretically and numerically. Recent works include application in poromechanics [3], time-harmonic Maxwell's equations with impedance surfaces [35], and the applications listed in [31].

## References

- [1] F. Assous, P. Ciarlet Jr, S. Labrunie, *Mathematical foundations of computational electromagnetism*, Springer, 2018.
- [2] I. Babuška, "The finite element method with Lagrangian multipliers", *Numerische Mathematik* **20** (1973), no. 3, p. 179-192.
- [3] M. Barré, C. Grandmont, P. Moireau, "Analysis of a linearized poromechanics model for incompressible and nearly incompressible materials", Technical Report HAL, 2021, <https://hal.inria.fr/hal-03501526>.
- [4] D. Boffi, F. Brezzi, M. Fortin *et al.*, *Mixed finite element methods and applications*, vol. 44, Springer, 2013.
- [5] A.-S. Bonnet-Ben Dhia, C. Carvalho, P. Ciarlet, Jr, "Mesh requirements for the finite element approximation of problems with sign-changing coefficients", *Numer. Math.* **138** (2018), p. 801–838.
- [6] A.-S. Bonnet-Ben Dhia, L. Chesnel, P. Ciarlet, Jr, "T-coercivity for scalar interface problems between dielectrics and metamaterials", *Math. Mod. Num. Anal.* **46** (2012), p. 1363-1387.
- [7] ———, "T-coercivity for the Maxwell problem with sign-changing coefficients", *Communications in Partial Differential Equations* **39** (2014), p. 1007-1031.
- [8] ———, "Two-dimensional Maxwell's equations with sign-changing coefficients", *Appl. Numer. Math.* **79** (2014), p. 29-41.
- [9] A.-S. Bonnet-Ben Dhia, L. Chesnel, X. Claeys, "Radiation condition for a non-smooth interface between a dielectric and a metamaterial", *Mathematical Models and Methods in Applied Sciences* **23** (2013), no. 09, p. 1629-1662.
- [10] A.-S. Bonnet-Ben Dhia, P. Ciarlet Jr, C. M. Zwölf, "Time harmonic wave diffraction problems in materials with sign-shifting coefficients", *Journal of Computational and Applied Mathematics* **234** (2010), no. 6, p. 1912-1919.
- [11] F. Brezzi, "On the existence, uniqueness and approximation of saddle-point problems arising from Lagrangian multipliers", *Publications mathématiques et informatique de Rennes* (1974), no. S4, p. 1-26.
- [12] A. Buffa, "Remarks on the discretization of some noncoercive operator with applications to heterogeneous Maxwell equations", *SIAM J. Numer. Anal.* **43** (2005), p. 1-18.
- [13] A. Buffa, S. H. Christiansen, "The electric field integral equation on Lipschitz screens: definitions and numerical approximation", *Numerische Mathematik* **94** (2003), no. 2, p. 229-267.
- [14] A. Buffa, M. Costabel, C. Schwab, "Boundary element methods for Maxwell's equations on non-smooth domains", *Numerische Mathematik* **92** (2002), no. 4, p. 679-710.
- [15] R. Bunoiu, L. Chesnel, K. Ramdani, M. Rihani, "Homogenization of Maxwell's equations and related scalar problems with sign-changing coefficients", in *Annales de la Faculté des Sciences de Toulouse. Mathématiques.*, 2020.
- [16] R. Bunoiu, K. Ramdani, "Homogenization of materials with sign changing coefficients", *Communications in Mathematical Sciences* **14** (2016), no. 4, p. 1137-1154.
- [17] R. Bunoiu, K. Ramdani, C. Timofte, "T-coercivity for the asymptotic analysis of scalar problems with sign-changing coefficients in thin periodic domains", *Electronic Journal of Differential Equations* (2021), p. 1-22.
- [18] L. Chesnel, "Bilaplacian problems with a sign-changing coefficient", *Mathematical Methods in the Applied Sciences* **39** (2016), no. 17, p. 4964-4979.
- [19] L. Chesnel, P. Ciarlet Jr, "T-coercivity and continuous Galerkin methods: application to transmission problems with sign changing coefficients", *Numerische Mathematik* **124** (2013), no. 1, p. 1-29.
- [20] P. Ciarlet Jr, "T-coercivity: Application to the discretization of Helmholtz-like problems", *Computers & Mathematics with Applications* **64** (2012), no. 1, p. 22-34.
- [21] ———, "Mathematical and numerical analyses for the div-curl and div-curlcurl problems with a sign-changing coefficient", Technical Report HAL, 2020, <https://hal.inria.fr/hal-02651682>.

- [22] ———, “Lecture notes on Maxwell’s equations and their approximation (in French)”, Master’s degree Analysis, Modelling and Simulation from Paris-Saclay University and Institut Polytechnique de Paris, 2021, <https://hal.inria.fr/hal-03153780>.
- [23] ———, “On the approximation of electromagnetic fields by edge finite elements – Part 4: analysis of the model with one sign-changing coefficient”, *Numer. Math.* **152** (2022), p. 223-257.
- [24] P. Ciarlet Jr, E. Jamelot, F. D. Kpadonou, “Domain Decomposition Methods for the diffusion equation with low-regularity solution”, *Computers Math. Applic.* **74** (2017), p. 2369-2384.
- [25] G. Duvaut, J. L. Lions, *Les inéquations en mécanique et en physique*, Dunod, 1972.
- [26] M. Fortin, “An analysis of the convergence of mixed finite element methods”, *RAIRO. Analyse numérique* **11** (1977), no. 4, p. 341-354.
- [27] V. Girault, P.-A. Raviart, *Finite element methods for Navier-Stokes equations: theory and algorithms*, vol. 5, Springer Science & Business Media, 2012.
- [28] L. Giret, “Numerical analysis of a non-conforming Domain Decomposition for the multigroup SPN equations”, PhD Thesis, Paris-Saclay University, 2018, <https://pastel.archives-ouvertes.fr/tel-01936967>.
- [29] M. Halla, “Galerkin approximation of holomorphic eigenvalue problems: weak T-coercivity and T-compatibility”, *Numerische Mathematik* **148** (2021), no. 2, p. 387-407.
- [30] R. Hiptmair, “Finite elements in computational electromagnetics”, *Acta Numerica* (2002), p. 237-339.
- [31] Q. Hong, J. Kraus, M. Lybery, F. Philo, “A new framework for the stability analysis of perturbed saddle-point problems and applications in poromechanics”, arXiv preprint, 2022, <https://arxiv.org/pdf/2103.09357.pdf>.
- [32] E. Jamelot, “T-coercivity for solving Stokes problem with nonconforming finite elements”, 2022.
- [33] E. Jamelot, P. Ciarlet Jr, “Fast non-overlapping Schwarz domain decomposition methods for solving the neutron diffusion equation”, *J. Comput. Phys.* **241** (2013), p. 445-463.
- [34] O. A. Ladyzhenskaya, *The mathematical theory of viscous incompressible flow*, vol. 2, Gordon and Breach New York, 1969.
- [35] D. P. Levadoux, “Analyse numérique de la formulation intégrodifférentielle d’un problème de Maxwell harmonique impliquant un diélectrique traversé de surfaces exfoliées métalliques et impédantes”, Technical Report HAL, 2022, <https://hal.archives-ouvertes.fr/hal-03644547>.
- [36] S. Nicaise, J. Venel, “A posteriori error estimates for a finite element approximation of transmission problems with sign changing coefficients”, *J. Comput. Appl. Math.* **235** (2011), p. 4272-4282.
- [37] E.-J. Sayas, T. S. Brown, M. E. Hassell, *Variational techniques for elliptic partial differential equations*, CRC Press, 2019.
- [38] C. Weber, “A local compactness theorem for Maxwell’s equations”, *Mathematical Methods in the Applied Sciences* **2** (1980), no. 1, p. 12-25.