



INSTITUT DE FRANCE  
Académie des sciences

# *Comptes Rendus*

---

## *Mathématique*


Erell Jamelot and François Madiot

**Numerical analysis of the neutron multigroup  $SP_N$  equations**

Volume 359, issue 5 (2021), p. 533-545

Published online: 13 July 2021

<https://doi.org/10.5802/crmath.189>

 This article is licensed under the  
CREATIVE COMMONS ATTRIBUTION 4.0 INTERNATIONAL LICENSE.  
<http://creativecommons.org/licenses/by/4.0/>



*Les Comptes Rendus. Mathématique* sont membres du  
Centre Mersenne pour l'édition scientifique ouverte  
[www.centre-mersenne.org](http://www.centre-mersenne.org)  
e-ISSN : 1778-3569



---

Numerical analysis / *Analyse numérique*

# Numerical analysis of the neutron multigroup $SP_N$ equations

## *Analyse numérique des équations de la neutronique $SP_N$ multigroupe*

Erell Jamelot<sup>a</sup> and François Madiot<sup>b</sup>

<sup>a</sup> Université Paris-Saclay, CEA, Service de Thermo-hydraulique et de Mécanique des Fluides, 91191, Gif-sur-Yvette, France

<sup>b</sup> Université Paris-Saclay, CEA, Service d'Études des Réacteurs et de Mathématiques Appliquées, 91191, Gif-sur-Yvette, France

*E-mails:* erell.jamelot@cea.fr, francois.madiot@cea.fr

**Abstract.** The multigroup neutron  $SP_N$  equations, which are an approximation of the neutron transport equation, are used to model nuclear reactor cores. In their steady state, these equations can be written as a source problem or an eigenvalue problem. We study the resolution of those two problems with an  $H^1$ -conforming finite element method and a Discontinuous Galerkin method, namely the Symmetric Interior Penalty Galerkin method.

**Résumé.** Les équations de la neutronique  $SP_N$  multigroupe, qui sont une approximation de l'équation de transport des neutrons, sont utilisées pour la modélisation des cœurs de réacteurs nucléaires. Dans le cas stationnaire, ces équations sont soit un problème à source, soit un problème aux valeurs propres. Nous étudions l'approximation de ces deux problèmes avec une méthode d'éléments finis conformes dans  $H^1$  et une méthode d'éléments finis discontinus appelée Symmetric Interior Penalty Galerkin.

*Manuscript received 4th October 2019, revised 24th August 2020, accepted 11th February 2021.*

### Version française abrégée

La physique d'un cœur de réacteur nucléaire est décrite par l'équation de transport des neutrons, qui dépend du temps et de 6 variables liées aux neutrons : 3 pour leur position, 2 pour la direction de leur vitesse et 1 pour leur énergie. Nous nous intéressons à la formulation stationnaire de cette équation (5) où l'énergie est discrétisée par la méthode multigroupe, et la direction est discrétisée par la méthode des harmoniques sphériques simplifiées  $SP_N$ . Cette formulation de l'équation de transport des neutrons revient à un système d'équations de la diffusion couplées. Nous proposons l'analyse numérique de ces équations discrétisées par une méthode d'éléments finis conformes dans  $H^1$  (resp. de Galerkin discontinus).

Nous commençons par l'étude des équations  $SP_N$  multigroupe pour le problème à source. À l'aide du lemme d'Aubin–Nitsche, nous obtenons une estimation d'erreur a priori dans  $L^2$  pour le problème à source discrétisé (12) (resp. (15)), énoncée dans le Théorème 5 (resp. 11). Puis nous nous intéressons au problème aux valeurs propres généralisé. Nous utilisons la théorie développée par Babuška et Osborn [2] pour obtenir une estimation d'erreur a priori sur la valeur propre, énoncée dans le Théorème 12 (resp. 13). Le Théorème 13 est obtenu à partir d'une généralisation de ces travaux présentée dans [1].

## 1. Introduction

The neutron transport equation describes the neutron flux density in a reactor core. It depends on 7 variables: 3 for the space, 2 for the motion direction, 1 for the energy (or the speed), and 1 for the time.

The energy variable is discretized using the multigroup theory [10, 16]. In this method, the entire range of neutron energies is divided into  $G$  intervals, called energy groups. In each energy group, the neutron flux density is lumped and all parameters are averaged. Let us set  $\mathcal{I}_G := \{1, \dots, G\}$ , the set of energy group indices.

Concerning the motion direction, the  $P_N$  transport equations are obtained by developing the neutron flux on the spherical harmonics from order 0 to order  $N$ . This approach is very time-consuming. The simplified  $P_N$  ( $SP_N$ ) transport theory [12] was developed to address this issue. The two fundamental hypotheses to obtain the  $SP_N$  equations are that locally, the angular flux has a planar symmetry; and that the axis system evolves slowly. The neutron flux and the scattering cross sections are then developed on the Legendre polynomials. From a mathematical point of view,  $SP_N$  equations correspond to tensorized 1D  $P_N$  transport equations, so that some couplings are missing. Consequently, the  $SP_N$  equations do not converge to transport equations. Nevertheless, they are commonly used by physicists since their resolution is cheap in terms of computational cost. The order  $N$  is odd, and the number of  $SP_N$  odd (resp. even) moments is  $\tilde{N} := \frac{N+1}{2}$ . We will denote by  $\mathcal{I}_e$  (resp.  $\mathcal{I}_o$ ) the subset of even (resp. odd) integers of the integer set  $\{0, \dots, N\}$ .

Finally, the (motion direction and energy) discretization of the neutron flux is such that there are  $\tilde{N} \times G$  even and odd moments of the neutron flux.

We will denote by  $\phi = ((\phi_m^g)_{m \in \mathcal{I}_e})_{g \in \mathcal{I}_G} \in \mathbb{R}^{\tilde{N} \times G}$  the set of functions containing, for all energy group  $g$ , the even moments of the neutron flux.

Likewise, we will denote by  $\mathbf{p} = ((p_{x,m}^g)_{m \in \mathcal{I}_o})_{g \in \mathcal{I}_G}^d \in \left(\mathbb{R}^{\tilde{N} \times G}\right)^d$  the set of functions containing the odd moments of the neutron flux.

Note that while modelling the core of a pressurized water reactor, the number of groups is such that  $2 \leq G \lesssim 30$ , physicists usually choose  $N = 1$  or 3, more rarely  $N = 5$ .

## 2. Setting of the model

The reactor core is modelled by a bounded, connected and open subset  $\mathcal{R}$  of  $\mathbb{R}^d$ ,  $d = 1, 2, 3$ , having a Lipschitz boundary which is piecewise regular. The coefficients are piecewise regular, so that we split  $\mathcal{R}$  into  $\tilde{N}$  open disjoint parts  $(\mathcal{R}_i)_{i=1}^{\tilde{N}}$  with Lipschitz, piecewise regular boundaries:  $\overline{\mathcal{R}} = \cup_{i=1}^{\tilde{N}} \overline{\mathcal{R}_i}$ . For this reason, we will use the following space of piecewise regular functions:

$$\mathcal{P}W^{1,\infty}(\mathcal{R}) = \left\{ D \in L^\infty(\mathcal{R}) \mid \vec{\nabla} D|_{\mathcal{R}_i} \in (L^\infty(\mathcal{R}_i))^d, i = 1, \dots, \tilde{N} \right\}.$$

For a set of functions  $\psi = (\phi_m^g)_{m,g} \in \mathbb{R}^{\tilde{N} \times G}$ , we make the following abuse of notation:  $\vec{\nabla} \psi = ((\partial_x \psi_m^g)_{m,g})_{x=1}^d \in \left(\mathbb{R}^{\tilde{N} \times G}\right)^d$ .

For a set of vector valued functions  $\mathbf{q} = \left( (q_{x,m}^g)_{m,g} \right)_{x=1}^d \in \left( \mathbb{R}^{\hat{N} \times G} \right)^d$ , we make the following abuse of notation:

$$\operatorname{div} \mathbf{q} = \left( \operatorname{div} \left( (q_{x,m}^g)_{x=1}^d \right) \right)_{m,g}, \quad \mathbf{q} \cdot \mathbf{p} = \left( \sum_{x=1}^d q_{x,m}^g p_{x,m}^g \right)_{m,g} \in \mathbb{R}^{\hat{N} \times G}.$$

Let us use these notations: for  $E \subset \mathbb{R}^d$ ,  $L(E) = L^2(E)$ ;  $L := L^2(\mathcal{R})$ ;  $V := H_0^1(\mathcal{R})$ ;  $V' := H^{-1}(\mathcal{R})$  its dual and  $Q := H(\operatorname{div}, \mathcal{R})$ . For  $W = L(E)$ ,  $L$ ,  $V$  or  $Q$  we define the product space  $\underline{W} := W^{\hat{N} \times G}$  endowed with the following scalar product and associated norm:

$$(\mathbf{u}, \mathbf{v})_{\underline{W}} = \sum_{g \in \mathcal{I}_G} \sum_{m \in \mathcal{I}_{e,o}} (\mathbf{u}_m^g, \mathbf{v}_m^g)_W, \quad \|\mathbf{u}\|_{\underline{W}}^2 = \sum_{g \in \mathcal{I}_G} \sum_{m \in \mathcal{I}_{e,o}} \|\mathbf{u}_m^g\|_W^2. \tag{1}$$

We also set  $\underline{V}' := (V')^{\hat{N} \times G}$ ,  $\underline{L}(E) = (L(E))^d$  and  $\underline{L}^p(\cdot) = (L^p(\cdot))^{\hat{N} \times G}$ .

Let  $\mathbf{q} \in \left( \mathbb{R}^{\hat{N} \times G} \right)^d$  and  $\mathbb{M} \in \left( \mathbb{R}^{\hat{N} \times \hat{N}} \right)^{G \times G}$ . We set  $\mathbf{q}_x = (q_{x,m}^g)_{m,g}$  and we use the notation  $\mathbb{M} \mathbf{q} = (\mathbb{M} \mathbf{q}_x)_{x=1}^d$ .

Given a source term  $S_f \in \underline{L}$ , the multigroup  $SP_N$  equations with zero-flux boundary conditions<sup>1</sup> read as coupled diffusion-like equations set in a mixed formulation:

$$\text{Solve in } (\phi, \mathbf{p}) \in \underline{V} \times \underline{Q} \mid \begin{cases} \mathbb{T}_o \mathbf{p} + \vec{\nabla}(\mathbb{H} \phi) = 0, \\ {}^t \mathbb{H} \operatorname{div} \mathbf{p} + \mathbb{T}_e \phi = S_f. \end{cases} \tag{2}$$

When  $S_f$  depends on  $\phi$ , the steady state multigroup  $SP_N$  equations read as the following generalized eigenproblem:

$$\text{Solve in } (\lambda, \phi, \mathbf{p}) \in \mathbb{R}^* \times \underline{V} \times \underline{Q} \mid \begin{cases} \mathbb{T}_o \mathbf{p} + \vec{\nabla}(\mathbb{H} \phi) = 0, \\ {}^t \mathbb{H} \operatorname{div} \mathbf{p} + \mathbb{T}_e \phi = \lambda^{-1} \mathbb{M}_f \phi. \end{cases} \tag{3}$$

The physical solution to Problem (3) corresponds to the eigenfunction associated with the smallest eigenvalue, which in addition is simple [8]. In neutronics, the *multiplication factor*  $k_{eff} = \max_{\lambda} \lambda$  characterizes the physical state of the core reactor: if  $k_{eff} = 1$ : the nuclear chain reaction is self-sustaining; if  $k_{eff} > 1$ : the chain reaction is diverging; if  $k_{eff} < 1$ : the chain reaction vanishes.

The matrices  $\mathbb{H}$ ,  $\mathbb{T}_e$ ,  $\mathbb{T}_o$ ,  $\mathbb{M}_f \in \left( \mathbb{R}^{\hat{N} \times \hat{N}} \right)^{G \times G}$  are such that  $\forall (g, g') \in \mathcal{I}_G \times \mathcal{I}_G$  ( $\delta_{\cdot, \cdot}$  is the Kronecker symbol):

- $(\mathbb{H})_{g,g'} = \delta_{g,g'} \widehat{\mathbb{H}} \in \mathbb{R}^{\hat{N} \times \hat{N}}$ , with  $\forall (i, j) \in \{1, \dots, \hat{N}\}^2$ ,  $\widehat{\mathbb{H}}_{i,j} = \delta_{i,j} + \delta_{i,j-1}$ .
- $(\mathbb{T}_e)_{g,g} := \mathbb{T}_e^g \in \mathbb{R}^{\hat{N} \times \hat{N}}$  denotes the even removal matrix, such that:

$$\mathbb{T}_e^g = \operatorname{diag} \left( t_0 \sigma_{r,0}^g, t_2 \sigma_{r,2}^g, \dots \right),$$

$(\mathbb{T}_o)_{g,g} := \mathbb{T}_o^g \in \mathbb{R}^{\hat{N} \times \hat{N}}$  denotes the odd removal matrix, such that:

$$\mathbb{T}_o^g = \operatorname{diag} \left( t_1 \sigma_{r,1}^g, t_3 \sigma_{r,3}^g, \dots \right),$$

where  $\forall m \in \mathcal{I}_{e,o}$ ,  $\sigma_{r,m}^g := \sigma_t^g - \sigma_{s,m}^{g \leftarrow g}$ , and  $\forall m > 0$ ,  $t_m > 0$ .

The coefficient  $\sigma_t^g$  is the macroscopic total cross section of energy group  $g$ , and the coefficients  $\sigma_{s,m}^{g \leftarrow g}$  denote the  $P_N$  moments of the macroscopic self scattering cross sections from energy group  $g$  to itself.

- For  $g' \neq g$ :  
 $(\mathbb{T}_e)_{g,g'} := -\mathbb{S}_e^{g' \leftarrow g} \in \mathbb{R}^{\hat{N} \times \hat{N}}$  denotes the even scattering matrix, such that:

$$\mathbb{S}_e^{g' \leftarrow g} = \operatorname{diag} \left( t_0 \sigma_{s,0}^{g' \leftarrow g}, t_2 \sigma_{s,2}^{g' \leftarrow g}, \dots \right),$$

<sup>1</sup>ie: for  $1 \leq g \leq G$ ,  $m \in \mathcal{I}_e$ ,  $(\phi_m^g)|_{\partial \mathcal{R}} = 0$ .

$(\mathbb{T}_o)_{g,g'} := -\mathbb{S}_o^{g' \rightarrow g} \in \mathbb{R}^{\hat{N} \times \hat{N}}$  denotes the odd scattering matrix, such that:

$$\mathbb{S}_o^{g' \rightarrow g} = \text{diag} \left( t_1 \sigma_{s,1}^{g' \rightarrow g}, t_3 \sigma_{s,3}^{g' \rightarrow g}, \dots \right),$$

where  $\sigma_{s,m}^{g' \rightarrow g}$  are the  $P_N$  moments of the macroscopic scattering cross sections from energy group  $g'$  to energy group  $g$ .

- $(\mathbb{M}_f)_{g,g'} := \chi^g \mathbb{M}_f^{g'} \in \mathbb{R}^{\hat{N} \times \hat{N}}$  is such that  $\mathbb{M}_f^{g'} \phi^{g'} = {}^t(\underline{v}\sigma_f^{g'} \phi_0^{g'}, 0, \dots)$  where the coefficient  $\underline{v}\sigma_f^{g'}$  is the product of the number of neutrons emitted per fission times the macroscopic fission cross section; and the coefficient  $\chi_g$  is the fission spectrum of energy group  $g$ .

The coefficients of the matrices  $\mathbb{T}_{e,o}, \mathbb{M}_f$  are supposed to be such that:

$$\left\{ \begin{array}{l} \text{(0)} \quad \forall g, g' \in \mathcal{I}_G, \forall m \in \mathcal{I}_{e,o} : \\ \quad (\sigma_{r,m}^g, \sigma_{s,m}^{g' \rightarrow g}, \underline{v}\sigma_f^g) \in \mathcal{D}W^{1,\infty}(\mathcal{R}) \times L^\infty(\mathcal{R}) \times L^\infty(\mathcal{R}). \\ \text{(i)} \quad \exists (\sigma_{r,(e,o)})^*, (\sigma^{r,(e,o)})^* > 0 \mid \forall g \in \mathcal{I}_G, \forall m \in \mathcal{I}_{e,o} : \\ \quad (\sigma_{r,(e,o)})^* \leq t_m \sigma_{r,m}^g \leq (\sigma^{r,(e,o)})^* \text{ a.e. in } \mathcal{R}. \\ \text{(ii)} \quad \exists (\underline{v}\sigma_f)^* > 0 \mid \forall g \in \mathcal{I}_G, 0 \leq \underline{v}\sigma_f^g \leq (\underline{v}\sigma_f)^* \text{ a.e. in } \mathcal{R} \text{ and } \exists g' \mid \underline{v}\sigma_f^{g'} \neq 0. \\ \text{(iii)} \quad \exists 0 < \varepsilon < \frac{1}{G-1} \mid \forall m \in \mathcal{I}_{e,o}, \forall g, g' \in \mathcal{I}_G, g' \neq g, \\ \quad |\sigma_{s,m}^{g \rightarrow g'}| \leq \varepsilon \sigma_{r,m}^g \text{ a.e. in } \mathcal{R}. \end{array} \right. \tag{4}$$

It happens that the coefficient  $\underline{v}\sigma_f^g$  vanishes in some regions.

Hypothesis 4 (iii) is valid while modelling the core of a pressurized water reactor: the scattering cross-sections are weaker than the removal cross-sections of an order  $0 < \varepsilon \ll 1$ . Thus, the matrices  ${}^t\mathbb{T}_{e,o}$  are strictly diagonally dominant matrices: they are invertible.

Let us set  $\mathbb{D} = {}^t\mathbb{H}\mathbb{T}_o^{-1}\mathbb{H}$ .

Problem 2 can be written as a set of coupled primal diffusion-like equations with single unknown  $\phi \in \underline{V}$ :

$$\text{Solve in } \phi \in \underline{V} \mid -\text{div}(\mathbb{D} \vec{\nabla} \phi) + \mathbb{T}_e \phi = S_f. \tag{5}$$

The variational formulation of (5) writes:

$$\text{Solve in } \phi \in \underline{V} \mid \forall \psi \in \underline{V} : c(\phi, \psi) = \ell(\psi), \tag{6}$$

where:  $\left\{ \begin{array}{l} c : \underline{V} \times \underline{V} \rightarrow \mathbb{R} \\ c(\phi, \psi) = (\mathbb{D} \vec{\nabla} \phi, \vec{\nabla} \psi)_{\underline{L}} + (\mathbb{T}_e \phi, \psi)_{\underline{L}} \end{array} \right.$ , and  $\left\{ \begin{array}{l} \ell : \underline{V} \rightarrow \mathbb{R} \\ \ell(\psi) = (S_f, \psi)_{\underline{L}} \end{array} \right.$ .

**Theorem 1.** Suppose that  $\mathbb{D}$  is positive definite. For a given source term  $S_f \in \underline{L}$ , it exists a unique  $\phi \in \underline{V}$  that solves Problem 6. In addition, it holds:  $\|\phi\|_{\underline{V}} \lesssim \|S_f\|_{\underline{L}}$ .

**Proof.** The bilinear form  $c$  and the linear form  $\ell$  are continuous and under the hypothesis on  $\mathbb{D}$ , the bilinear form  $c$  is coercive: we can apply Lax–Milgram theorem to conclude.  $\square$

In the same way, Problem 3 can be written as:

$$\text{Solve in } (\lambda, \phi) \in \mathbb{R}^* \times \underline{V} \setminus \{0\} \mid -\text{div}(\mathbb{D} \vec{\nabla} \phi) + \mathbb{T}_e \phi = \lambda^{-1} \mathbb{M}_f \phi. \tag{7}$$

The variational formulation of (7) writes:

$$\text{Solve in } (\lambda, \phi) \in \mathbb{R}^* \times \underline{V} \setminus \{0\} \mid \forall \psi \in \underline{V} : c(\phi, \psi) = \lambda^{-1} \ell_f(\phi, \psi), \tag{8}$$

where:  $\left\{ \begin{array}{l} \ell_f : \underline{L} \times \underline{L} \rightarrow \mathbb{R} \\ \ell_f(\phi, \psi) = (\mathbb{M}_f \phi, \psi)_{\underline{L}} \end{array} \right.$ .

**Theorem 2.** Suppose that  $\mathbb{D}$  is positive definite. There exists a unique compact operator  $T_f : \underline{L} \rightarrow \underline{L}$  such that  $\forall (\phi, \psi) \in \underline{L} \times \underline{V} : c(T_f \phi, \psi) = \ell_f(\phi, \psi)$ .

**Proof.** The bilinear form  $c$  is a continuous and under the hypothesis on  $\mathbb{D}$ , it is coercive onto  $\underline{V} \times \underline{V}$ . The bilinear form  $\ell_f$  is a continuous onto  $\underline{L} \times \underline{V}$ . Finally,  $\underline{V}$  is a subset of  $\underline{L}$  with a compact embedding. We can then apply the work of Babuška and Osborn in [2].  $\square$

Thus, the couple  $(\phi, \lambda^{-1})$  is a solution to Problem 8 iff the couple  $(\phi, \lambda)$  is an eigenpair of operator  $T_f$ . Moreover, Problem 8 admits a countable number of eigenvalues.

We propose first to derive conditions on the macroscopic cross sections so that Problems 5 and 7 are well-posed. Then we obtain a priori error estimates for a discretization performed with some  $H^1$ -conforming FEM and a Discontinuous Galerkin method, namely the Symmetric Interior Penalty Galerkin method (SIPG) [9, Chapter 4]. The outline is as follows: in Section 3, we exhibit some conditions so that the matrix  $\mathbb{T}_o^{-1}$  and  $\mathbb{T}_e$  are positive definite. Then we study the discretization of the source problem (5) in Section 5, and the discretization of the eigenproblem in Section 6. Finally, we perform in Section 7 a numerical study of convergence on a benchmark representative of a nuclear core.

### 3. Properties of $\mathbb{T}_e$ and $\mathbb{T}_o^{-1}$

Consider the diagonal matrix containing the even (resp. odd) removal macroscopic cross sections:  $\mathbb{T}_{r,(e,o)} = \text{diag}(\mathbb{T}_{e,o}^1, \dots, \mathbb{T}_{e,o}^G)$ . We split  $\mathbb{T}_{e,o}$  so that:  $\mathbb{T}_{e,o} = \mathbb{T}_{r,(e,o)}(\mathbb{I} - \varepsilon \mathbb{U}_{e,o})$ , where  $\mathbb{I} \in (\mathbb{R}^{\hat{N} \times \hat{N}})^{G \times G}$  is the identity matrix, and:

$$\begin{aligned} \forall g, g' \in \mathcal{I}_G, g' \neq g, \quad (\mathbb{U}_{e,o})_{g,g'} &= \text{diag} \left( \left( \begin{array}{c} \sigma_{s,m}^{g'-g} \\ \varepsilon \sigma_{r,m}^g \end{array} \right)_{m \in \mathcal{I}_{e,o}} \right) \in \mathbb{R}^{\hat{N} \times \hat{N}}, \\ \forall g \in \mathcal{I}_G, \quad (\mathbb{U}_{e,o})_{g,g} &= 0 \in \mathbb{R}^{\hat{N} \times \hat{N}}. \end{aligned}$$

We have then:  $\|\mathbb{U}_{e,o}\|_2 \lesssim \frac{\alpha_{s,(e,o)}}{\varepsilon}$  where:  $\alpha_{s,(e,o)} := (G - 1) \max_{m \in \mathcal{I}_{e,o}} \max_{g \neq g' \in \mathcal{I}_G} \sup_{\vec{x} \in \mathcal{D}} \frac{|\sigma_{s,m}^{g'-g}(\vec{x})|}{\sigma_{r,m}^g(\vec{x})}$ .

Let us set  $\alpha_{r,(e,o)} = \frac{(\sigma_{r,(e,o)})^*}{(\sigma_{r,(e,o)})} > 1$ . We have the following properties.

**Property 3.** Suppose that  $\alpha_{s,e} < \frac{1}{\alpha_{r,e}}$ . The matrix  $\mathbb{T}_e$  is such that:

$$\forall X \in \mathbb{R}^{\hat{N} \times G} \quad (\mathbb{T}_e X | X) \geq \tau_e \|X\|_2^2 \quad \text{where } \tau_e = (\sigma_{r,e})^* (1 - \alpha_{r,e} \alpha_{s,e}). \tag{9}$$

**Proof.** We have:  $\forall X \in \mathbb{R}^{\hat{N} \times G}$ ,  $(\mathbb{T}_e X | X) = (\mathbb{T}_{r,e} X | X) - \varepsilon (\mathbb{U}_{e,o} X | \mathbb{T}_{r,e} X)$ , so that:

$$(\mathbb{T}_e X | X) \geq ((\sigma_{r,e})^* - \varepsilon \|\mathbb{U}_{e,o}\|_2 \|\mathbb{T}_{r,e}\|_2) \|X\|_2, \quad \text{where } \|\mathbb{T}_{r,e}\|_2 \leq (\sigma_{r,e})^*. \tag{10}$$

**Property 4.** Suppose that  $\alpha_{s,o} < \frac{1}{\alpha_{r,o} + 1}$ , the matrix  $\mathbb{T}_o^{-1}$  is such that:

$$\forall X \in \mathbb{R}^{\hat{N} \times G} \quad (\mathbb{T}_o^{-1} X | X) \geq \tau_o \|X\|_2^2 \quad \text{where } \tau_o = \frac{1}{(\sigma_{r,o})^*} \left( 1 - \frac{\alpha_{r,o} \alpha_{s,o}}{1 - \alpha_{s,o}} \right). \tag{10}$$

**Proof.** The Taylor expansion of  $\mathbb{T}_o^{-1}$  writes:  $\mathbb{T}_o^{-1} = (\mathbb{I} + \sum_{l>0} \varepsilon^l \mathbb{U}_o^l) \mathbb{T}_{r,o}^{-1}$ .

We get that  $\forall X \in \mathbb{R}^{\hat{N} \times G}$ :

$$\begin{aligned} (\mathbb{T}_o^{-1} X | X) &= (\mathbb{T}_{r,o}^{-1} X | X) + \sum_{l>0} \varepsilon^l (\mathbb{U}_o^l \mathbb{T}_{r,o}^{-1} X | X) \\ &\geq \frac{1}{(\sigma_{r,o})^*} \left( 1 - \alpha_{r,o} \sum_{l>0} \varepsilon^l \|\mathbb{U}_o\|_2^l \right) \|X\|_2^2, \\ &\geq \frac{1}{(\sigma_{r,o})^*} \left( 1 - \alpha_{r,o} \frac{\varepsilon \|\mathbb{U}_o\|_2}{1 - \varepsilon \|\mathbb{U}_o\|_2} \right) \|X\|_2^2, \\ &\geq \frac{1}{(\sigma_{r,o})^*} \left( 1 - \frac{\alpha_{r,o} \alpha_{s,o}}{1 - \alpha_{s,o}} \right) \|X\|_2^2. \end{aligned} \tag{10}$$

Under assumptions of Properties 3 and 4 the matrices  $\mathbb{T}_e$  and  $\mathbb{T}_o^{-1}$  are positive definite. Moreover, one can show that  $\|\mathbb{H} \tilde{\nabla} \phi\|_{\underline{\mathbb{L}}} \gtrsim \|\tilde{\nabla} \phi\|_{\underline{\mathbb{L}}}$  [13]. We infer that the matrix  $\mathbb{D}$  is positive definite and that there exists a constant  $C_{\mathbb{D}} > 0$  such that for all  $\xi \in \mathbb{R}^{\hat{N} \times G}$ ,

$$(\mathbb{D} \xi | \mathbb{D} \xi) \leq C_{\mathbb{D}} \|\xi\|_2^2. \tag{11}$$

From now on, we suppose that this property holds.

#### 4. Discretizations

Let  $\mathcal{T}_h$  be a shape-regular mesh of  $\mathcal{R}$ , with mesh size  $h$ . We denote by  $K$  its elements and  $F$  its facets. To simplify the presentation, we assume that the meshes are such that in every element, the cross-sections are regular. We define by  $\mathcal{F}_h^i$  the set of interior faces of  $\mathcal{T}_h$ ,  $\mathcal{F}_h^b$  the set of boundary facets and  $\mathcal{F}_h = \mathcal{F}_h^i \cup \mathcal{F}_h^b$ . We denote by  $N_{\partial}$  the maximum number of mesh faces composing the boundary of mesh elements

$$N_{\partial} := \max_{K \in \mathcal{T}_h} \text{Card}\{F \in \mathcal{F}_h, F \subset \partial K\}.$$

We will first consider an  $H^1$ -conforming finite element method (FEM). For  $k \in \mathbb{N}^*$ ,  $V_h^k \subset V$  and  $\underline{V}_h^k \subset \underline{V}$  are the finite dimension spaces defined by:

$$V_h^k = \{v_h \in V, \forall K \in \mathcal{T}_h, v_h|_K \in \mathbb{P}_k\}, \quad \underline{V}_h^k := (V_h^k)^{\hat{N} \times G}.$$

The discrete variational formulation associated with Problem (6) writes:

$$\text{Solve in } \phi_h \in \underline{V}_h^k \mid \forall \psi_h \in \underline{V}_h^k : c(\phi_h, \psi_h) = \ell(\psi_h), \tag{12}$$

Similarly, the discrete variational formulation associated with Problem (7) writes:

$$\text{Solve in } (\lambda_h, \phi_h) \in \mathbb{R}^* \times \underline{V}_h^k \setminus \{0\} \mid \forall \psi \in \underline{V}_h^k : c(\phi_h, \psi_h) = \lambda_h^{-1} \ell_f(\phi_h, \psi_h). \tag{13}$$

Then, we will consider a non-conforming FEM. We define the broken spaces:

$$V_{\text{NC}} = \{v \in L^2(\mathcal{R}) \mid \forall K \in \mathcal{T}_h, v \in H^1(K)\}, \quad \underline{V}_{\text{NC}} = (V_{\text{NC}})^{\hat{N} \times G}.$$

For  $(\phi, \psi) \in \underline{V}_{\text{NC}} \times \underline{V}_{\text{NC}}$ , and  $\mathbb{T} \in \mathbb{R}^{\hat{N} \times G}$ , we set:

$$(\mathbb{D} \tilde{\nabla}_h \phi, \tilde{\nabla}_h \psi)_{\mathcal{F}_h} = \sum_{K \in \mathcal{T}_h} (\mathbb{D} \tilde{\nabla} \phi, \tilde{\nabla} \psi)_{\underline{\mathbb{L}}(K)}, \quad \text{and} \quad \|\tilde{\nabla}_h \psi\|_{\mathcal{F}_h} = (\tilde{\nabla}_h \psi, \tilde{\nabla}_h \psi)_{\mathcal{F}_h}^{1/2}.$$

For  $F \in \mathcal{F}_h^i$  such that  $F = \partial K_1 \cap \partial K_2$ , we define the average  $\{\mathbb{D} \tilde{\nabla}_h \psi\}$  and the jump  $\llbracket \psi \rrbracket$  as:

$$\begin{aligned} \{\mathbb{D} \tilde{\nabla}_h \psi\}|_F &= \frac{1}{2} \left( (\mathbb{D}_1 \tilde{\nabla} \psi_1)|_F + (\mathbb{D}_2 \tilde{\nabla} \psi_2)|_F \right) \in \left( \mathbb{R}^{\hat{N} \times G} \right)^d, \\ \llbracket \psi \rrbracket|_F &= \psi_1|_F \mathbf{n}_1 + \psi_2|_F \mathbf{n}_2 \in \left( \mathbb{R}^{\hat{N} \times G} \right)^d. \end{aligned}$$

where  $\mathbf{n}_i$  is the unit outward normal to  $K_i$  at face  $F$  and  $\mathbb{D}_i = \mathbb{D}|_{K_i}$ ,  $\psi_i = \psi|_{K_i}$ .

For  $F \in \mathcal{F}_h^b$  such that  $F \in K$ , we set  $\{\mathbb{D} \tilde{\nabla}_h \psi\}|_F = \mathbb{D}|_K \tilde{\nabla} \psi|_K$  and  $\llbracket \psi \rrbracket|_F = (\psi_K)|_F \mathbf{n}$ , where  $\psi_K = \psi|_K$  and  $\mathbf{n}$  is the unit outward normal to  $K$  at face  $F$ .

For  $k \in \mathbb{N}^*$ ,  $V_{h,\text{NC}}^k \subset H^1(\mathcal{T}_h)$  and  $\underline{V}_{h,\text{NC}}^k$  are the finite dimension spaces defined by:

$$V_{h,\text{NC}}^k = \{v_h \in L^1(\mathcal{R}); \forall K \in \mathcal{T}_h, v_h|_K \in \mathbb{P}_k\}, \quad \underline{V}_{h,\text{NC}}^k := \left( V_{h,\text{NC}}^k \right)^{\hat{N} \times G}.$$

For  $\phi_h, \psi_h \in \underline{V}_{h,\text{NC}}^k$ , we set:  $(\{\mathbb{D} \tilde{\nabla}_h \phi_h\}, \llbracket \psi_h \rrbracket)_{\mathcal{F}_h^i} = \sum_{F \in \mathcal{F}_h^i} (\{\mathbb{D} \tilde{\nabla}_h \phi_h\}, \llbracket \psi_h \rrbracket)_{\underline{\mathbb{L}}(F)}$ .

Let us set

$$c_h(\phi_h, \psi_h) = c_{\mathcal{F}_h}(\phi_h, \psi_h) + c_{\mathcal{F}_h^i}(\phi_h, \psi_h), \tag{14}$$

with

$$c_{\mathcal{T}_h}(\phi_h, \psi_h) = (\mathbb{D} \vec{\nabla}_h \phi_h, \vec{\nabla}_h \psi_h)_{\mathcal{T}_h} + (\mathbb{T}_e \phi_h, \psi_h)_{\underline{L}},$$

$$c_{\mathcal{F}_h}(\phi_h, \psi_h) = \sum_{F \in \mathcal{F}_h} \frac{\alpha}{h_F} (\llbracket \phi_h \rrbracket, \llbracket \psi_h \rrbracket)_{\underline{L}(F)} - (\{\mathbb{D} \vec{\nabla}_h \psi_h\}, \llbracket \phi_h \rrbracket)_{\mathcal{F}_h^i} - (\{\mathbb{D} \vec{\nabla}_h \phi_h\}, \llbracket \psi_h \rrbracket)_{\mathcal{F}_h^i},$$

where  $\alpha$  is a stabilization parameter.

The Symmetric Interior Penalty Galerkin method (SIPG) associated with Problem (6) writes:

$$\text{Solve in } \phi_h \in \underline{V}_{h,\text{NC}}^k \mid \forall \psi_h \in \underline{V}_{h,\text{NC}}^k : c_h(\phi_h, \psi_h) = \ell(\psi_h). \tag{15}$$

Similarly, the SIPG method associated with Problem (8) writes:

$$\text{Solve in } (\lambda_h, \phi_h) \in \mathbb{R}^* \times \underline{V}_{h,\text{NC}}^k \setminus \{0\} \mid \forall \psi_h \in \underline{V}_{h,\text{NC}}^k : c_h(\phi_h, \psi_h) = \lambda_h^{-1} \ell_f(\phi_h, \psi_h). \tag{16}$$

## 5. The source problem

### 5.1. Conforming discretization

**Theorem 5.** *Suppose that there exists  $r_{\max}$  in  $[0, 1]$  such that  $\forall r \in [0, r_{\max}[$ ,  $\phi \in (H^{1+r}(\mathcal{R}))^{\hat{N} \times G}$  (cf. [6, Proposition 1]). Let us set  $\mu = \min(r_{\max}, k)$ . The solution of (12),  $\phi_h$  is such that:  $\|\phi - \phi_h\|_{\underline{V}} \lesssim h^\mu \|S_f\|_{\underline{L}}$  and  $\|\phi - \phi_h\|_{\underline{L}} \lesssim h^{2\mu} \|S_f\|_{\underline{L}}$ .*

**Proof.** From Céa’s lemma and Aubin–Nitsche lemma as detailed in [11, Section 2.3]. □

### 5.2. SIPG discretization

**Assumption 6 (Regularity of exact solution and space  $V^*$ ).** *Let us denote by  $W^{2,p}(\mathcal{T}_h)$  the broken Sobolev space spanned by those functions  $v$  such that for all  $K \in \mathcal{T}_h$ ,  $v|_K \in W^{2,p}(K)$ . We set  $\underline{W}^{2,p}(\mathcal{T}_h) = (W^{2,p}(\mathcal{T}_h))^{\hat{N} \times G}$ . We assume that  $d \geq 2$  and that there is  $2d/(d+2) < p \leq 2$  such that, for the exact solution  $\phi \in \underline{V}^* := \underline{V} \cap \underline{W}^{2,p}(\mathcal{T}_h)$ . This holds for our assumptions on the coefficients, which are piecewise constant with respect to the triangulation [17].*

This assumption requires  $p > 1$  for  $d = 2$  and  $p > 6/5$  for  $d = 3$ . In particular, we observe that, in two space dimensions,  $\phi \in \underline{W}^{2,p}(\mathcal{T}_h)$  in polygonal domains. Moreover, using Sobolev embeddings [4, Section IX.3] [7], this implies

$$\phi \in (H^{1+\alpha_p}(\mathcal{R}))^{\hat{N} \times G}, \quad \alpha_p = \frac{d+2}{2} - \frac{d}{p} > 0.$$

We state the following lemma [9, Lemma 1.46, p. 27].

**Lemma 7.** *Suppose that  $(\mathcal{T}_h)_h$  is a shape- and contact-regular mesh sequence. Then, we have for all  $h > 0$ :*

$$\forall \psi_h \in \underline{V}_{h,\text{NC}}^k, \forall K \in \mathcal{T}_h, \forall F \in \partial K, \quad h_K^{1/2} \|\psi_h\|_{\underline{L}^2(F)} \leq C_{\text{tr}} \|\psi_h\|_{\underline{L}^2(K)}, \tag{17}$$

where  $h_K$  is the diameter of element  $K$ .

We aim at asserting the discrete coercivity using the following norm:

$$\forall \psi_h \in \underline{V}_{h,\text{NC}}^k, \quad \|\|\| \psi_h \|\|\|_{\text{sip}}^2 := c_{\mathcal{T}_h}(\psi_h, \psi_h) + \|\psi_h\|_J^2,$$

with the jump semi-norm

$$\|\psi_h\|_J^2 := \sum_{F \in \mathcal{F}_h} \frac{1}{h_F} \|\llbracket \psi_h \rrbracket\|_{\underline{L}(F)}^2.$$

Under assumption (4), there exists  $\beta > 0$  we have for all  $\psi_h \in \underline{V}_{h,\text{NC}}^k$

$$c_{\mathcal{T}_h}(\psi_h, \psi_h) \geq \beta \left( \|\vec{\nabla}_h \psi_h\|_{\mathcal{T}_h}^2 + \|\psi_h\|_{\underline{L}}^2 \right), \tag{18}$$



so that

$$\|\psi_h\|_{sip}^2 \geq \beta \left( \|\tilde{\nabla}_h \psi_h\|_{\mathcal{T}_h}^2 + \|\psi_h\|_{\underline{L}}^2 + \|\psi_h\|_J^2 \right).$$

**Lemma 8 (Discrete coercivity).** *Let  $\underline{\alpha} := C_{tr}^2 N_\partial \frac{C_D}{\beta}$  where*

- $C_{tr}$  results from the discrete trace inequality (17),
- $N_\partial$  is defined in Section 4,
- $C_D$  is defined in (11).

For all  $\alpha \geq \underline{\alpha}$ , the SIP bilinear form defined by (14) is coercive on  $\underline{V}_{h,NC}^k$  with respect to the  $\|\cdot\|_{sip}$ -norm, i.e.,

$$c_h(\psi_h, \psi_h) \geq C_\alpha \|\psi_h\|_{sip}^2,$$

with  $C_\alpha := \left( \alpha - C_{tr}^2 N_\partial \frac{C_D}{\beta} \right) \min \left\{ \frac{1}{2}, \beta \left( \alpha + C_{tr}^2 N_\partial \frac{C_D}{\beta} \right)^{-1} \right\}$ .

**Proof.** We follow the proof of [9, Lemma 4.12]. For all  $\psi_h \in \underline{V}_{h,NC}^k$ ,

$$\begin{aligned} c_h(\psi_h, \psi_h) &= c_{\mathcal{T}_h}(\psi_h, \psi_h) + c_{\mathcal{F}_h}(\psi_h, \psi_h) \\ &= c_{\mathcal{T}_h}(\psi_h, \psi_h) + \sum_{F \in \mathcal{F}_h} \frac{\alpha}{h_F} \|\llbracket \psi_h \rrbracket\|_{\underline{L}(F)}^2 - 2 \left( \{\mathbb{D} \tilde{\nabla}_h \psi_h\}, \llbracket \psi_h \rrbracket \right)_{\mathcal{F}_h^i} \\ &\geq c_{\mathcal{T}_h}(\psi_h, \psi_h) + \alpha \|\psi_h\|_J^2 - 2C_{tr}(N_\partial)^{1/2} \|\mathbb{D} \tilde{\nabla}_h \psi_h\|_{\mathcal{T}_h} \|\psi_h\|_J \end{aligned}$$

where we used Cauchy–Schwarz and Lemma 7 in the last line. Using the inequality  $2ab \leq \varepsilon a + \varepsilon^{-1}b$  for any  $\varepsilon > 0$ , we obtain

$$\begin{aligned} 2C_{tr}(N_\partial)^{1/2} \|\mathbb{D} \tilde{\nabla}_h \psi_h\|_{\mathcal{T}_h} \|\psi_h\|_J &\leq \varepsilon C_{tr}^2 N_\partial \|\mathbb{D} \tilde{\nabla}_h \psi_h\|_{\mathcal{T}_h}^2 + \varepsilon^{-1} \|\psi_h\|_J^2 \\ &\leq \varepsilon C_{tr}^2 N_\partial C_D \|\tilde{\nabla}_h \psi_h\|_{\mathcal{T}_h}^2 + \varepsilon^{-1} \|\psi_h\|_J^2. \end{aligned}$$

Using (18), we obtain that there exists a constant  $\beta > 0$  such that

$$c_h(\psi_h, \psi_h) \geq \beta(1 - \varepsilon \underline{\alpha}) \|\tilde{\nabla}_h \psi_h\|_{\mathcal{T}_h}^2 + \beta \|\psi_h\|_{\underline{L}}^2 + (\alpha - \varepsilon^{-1}) \|\psi_h\|_J^2.$$

Choosing  $\varepsilon = 2(\alpha + \underline{\alpha})^{-1}$  yields the assertion. □

Thus, it only remains to prove boundedness. To this purpose, we need to define  $\underline{V}^{*,h} = \underline{V}^* + \underline{V}_{h,NC}^k$  and the following norm

$$\|\psi\|_{sip, \star} := \left( \|\psi\|_{sip}^p + \sum_{K \in \mathcal{T}_h} h_K^{1+\gamma_p} \|\tilde{\nabla} \psi|_K \cdot \mathbf{n}_K\|_{\underline{L}^p(\partial K)} \right)^{1/p},$$

where  $\gamma_p = \frac{d(p-2)}{2}$  and  $\mathbf{n}_K$  is the unit outward normal to  $K$ . Following [9, Section 4.2], we obtain the following results.

**Lemma 9 (Boundedness).** *There is  $C_{bnd}$ , independent of  $h$ , such that for all  $(\phi, \psi_h) \in \underline{V}^{*,h} \times \underline{V}_h$*

$$c_h(\phi, \psi_h) \leq C_{bnd} \|\phi\|_{sip, \star} \|\psi_h\|_{sip}.$$

**Theorem 10 (Convergence).** *Suppose that there exists  $r_{\max}$  in  $(0, 1]$  such that  $\forall r \in [0, r_{\max}]$ ,  $\phi \in (H^{1+r}(\mathcal{R}))^{\tilde{N} \times G}$  (cf. [6, Proposition 1]). Then the solution of (15),  $\phi_h$  is such that:*

$$\|\phi - \phi_h\|_{sip} \lesssim C \inf_{\psi_h \in \underline{V}_{h,NC}} \|\phi - \psi_h\|_{sip, \star},$$

where  $C$  is a constant independent of  $h$ . Moreover, under Assumption 6, there holds

$$\|\phi - \phi_h\|_{sip} \leq C |\phi|_{\underline{W}^{2,p}(\mathcal{T}_h)} h^\mu,$$

where  $\mu = r_{\max}$ ,  $C$  is a constant independent of  $h$  and  $p$  is such that  $\mu = \frac{d+2}{2} - \frac{d}{p}$ .

**Theorem 11 ( $L^2$ -norm estimate).** *Suppose that there exists  $r_{\max}$  in  $(0, 1]$  such that  $\forall r \in [0, r_{\max}]$ ,  $\phi_m^s \in H^{1+r}(\mathcal{R})$  (cf. [6, Proposition 1]). Under Assumption 6, the solution of (15),  $\phi_h$  is such that:  $\|\phi - \phi_h\|_{\underline{L}} \lesssim h^{2\mu} \|S_f\|_{\underline{L}}$ , where  $\mu = r_{\max}$ .*

**Proof.** We apply the Aubin–Nitsche similarly as in [9, Theorem 4.25]. □

## 6. The eigenproblem

### 6.1. Conforming discretization

**Theorem 12.** *Let  $\mu$  be the regularity of the eigenfunction  $\phi$  associated with  $\lambda$ , and  $\omega = \min(\mu, k)$ . Let  $\lambda_h$  be the discrete eigenvalue associated with Problem (13). The following a priori error estimate holds:  $|\lambda - \lambda_h| \lesssim h^{2\omega}$ .*

**Proof.** As in the continuous case (Theorem 2), since the discretization is conforming, there exists a unique compact operator  $T_h : \underline{V}_h^k \rightarrow \underline{V}_h^k$  such that  $\forall (\phi_h, \psi_h) \in \underline{V}_h^k \times \underline{V}_h^k: c(T_h \phi_h, \psi_h) = \ell_f(\phi_h, \psi_h)$ . According to Theorem 5, the sequence of the operators  $(T_h)_h$  is pointwise converging towards  $T$ . As  $T_h$  and  $T$  are compact operators, the sequence of operators  $(T_h)_h$  is then converging in  $\mathcal{L}(\underline{V})$  towards  $T$ :  $\|T_h - T\|_{\mathcal{L}(\underline{V})} \rightarrow 0$ . The norm convergence guarantees that there is no spectral pollution (see [18]). Moreover, we can apply Theorem 8.3 in [2] to state the error estimate on the eigenvalue. We remark that  $(\mathbb{M}_f \phi, \phi)_{\underline{L}}$  is a norm over  $\underline{V}_\lambda := \{\phi \in \underline{V} \mid \forall \psi \in \underline{V}, c(\phi, \psi) = \lambda \ell_f(\phi, \psi)\}$  [13, Section 5.2.2 p. 78]. □

### 6.2. SIPG discretization

We recall that, in this section, we work under the assumption 6.

**Theorem 13.** *Let  $\mu$  be the regularity of the eigenfunction  $\phi$  associated with  $\lambda$ , and  $\omega = \min(\mu, k)$ . Let  $\lambda_h$  be the discrete eigenvalue associated with Problem (16). The following a priori error estimate holds:  $|\lambda - \lambda_h| \lesssim h^{2\omega}$ .*

**Proof.** We apply the theory developed in [1]. The proof is decomposed as follows. We first show that there is no spectral pollution. Then, we derive the error estimate.

Let  $E : \underline{V} + \underline{V}_{h,NC}^k \rightarrow \underline{V} + \underline{V}_{h,NC}^k$  be the continuous spectral projector relative to  $\lambda$  defined by

$$E = \frac{1}{2\pi i} \int_{\Gamma} \left( z - T|_{\underline{V} + \underline{V}_{h,NC}^k} \right)^{-1} dz,$$

where  $\Gamma$  is a circle in the complex plane centred at  $\lambda$  which lies in  $\rho(T|_{\underline{V} + \underline{V}_{h,NC}^k})$  and encloses no other points of  $\sigma(T|_{\underline{V} + \underline{V}_{h,NC}^k})$ . The absence of spectral pollution relies on two properties. First, using interpolation results [9, Assumption 4.31] we have for all  $\phi \in E(\underline{V} + \underline{V}_{h,NC}^k)$ ,

$$\inf_{\psi_h \in \underline{V}_{h,NC}^k} \|\phi - \psi_h\|_{sip} \leq Ch^\mu,$$

where  $C$  is a constant independent of  $h$ . Second, we have for all  $\phi_h \in \underline{V}_{h,NC}^k$ ,

$$\begin{aligned} \|(T - T_h)\phi_h\|_{sip} &\leq Ch^\mu |T\phi_h|_{W^{2,p}(\mathcal{T}_h)}, \\ &\leq Ch^\mu \|T\phi_h\|_{(H^{1+\alpha_p}(\mathcal{R}))^{\tilde{N} \times G}}, \\ &\leq Ch^\mu \|\phi_h\|_{\underline{L}}, \\ &\leq Ch^\mu \|\phi_h\|_{sip}, \end{aligned}$$

where we used Theorem 10 in the second line and regularity results [17] in the third line. Applying [1, Theorem 3.7], we obtain that there is no spectral pollution.

Moreover, we apply [1, Theorem 3.14] to state the error estimate on the eigenvalue,

$$|\lambda - \lambda_h| \leq C \delta_h \delta_{*,h},$$

where

$$\begin{aligned} \delta_h &= \gamma_h + \left\| (T - T_h)|_{E(\underline{V} + \underline{V}_{h,\text{NC}}^k)} \right\|_{\text{sip}}, \\ \delta_{*,h} &= \gamma_{*,h} + \left\| (T_* - T_{*,h})|_{E(\underline{V} + \underline{V}_{h,\text{NC}}^k)} \right\|_{\text{sip}}, \end{aligned}$$

with

$$\begin{aligned} \gamma_h &= \delta(E(V + \underline{V}_{h,\text{NC}}^k), \underline{V}_{h,\text{NC}}^k), \\ \gamma_{*,h} &= \delta(E_*(V + \underline{V}_{h,\text{NC}}^k), \underline{V}_{h,\text{NC}}^k), \end{aligned}$$

where

$$\delta(Y, Z) = \sup_{y \in Y, \|y\|_{\text{sip}}=1} \left( \inf_{z \in Z} \|y - z\|_{\text{sip}} \right),$$

and  $E_* : \underline{V} + \underline{V}_{h,\text{NC}}^k \rightarrow \underline{V} + \underline{V}_{h,\text{NC}}^k$  is the continuous spectral projector of the adjoint operator  $T_*|_{\underline{V} + \underline{V}_{h,\text{NC}}^k}$  relative to  $\bar{\lambda}$ .

Using again elliptic regularity results [17] and Theorem 10, we obtain

$$\begin{aligned} \left\| (T - T_h)|_{E(\underline{V} + \underline{V}_{h,\text{NC}}^k)} \right\|_{\text{sip}} &\leq Ch^\mu, \\ \left\| (T_* - T_{*,h})|_{E(\underline{V} + \underline{V}_{h,\text{NC}}^k)} \right\|_{\text{sip}} &\leq Ch^\mu. \end{aligned}$$

Using elliptic regularity results, we get

$$\|\varphi\|_{(H^{1+\alpha_p}(\mathcal{R}))^{\hat{N} \times G}} \leq C \|\varphi\|_{\underline{L}} \leq C \|\varphi\|_{\underline{V}}.$$

Applying Theorem 10, we infer that

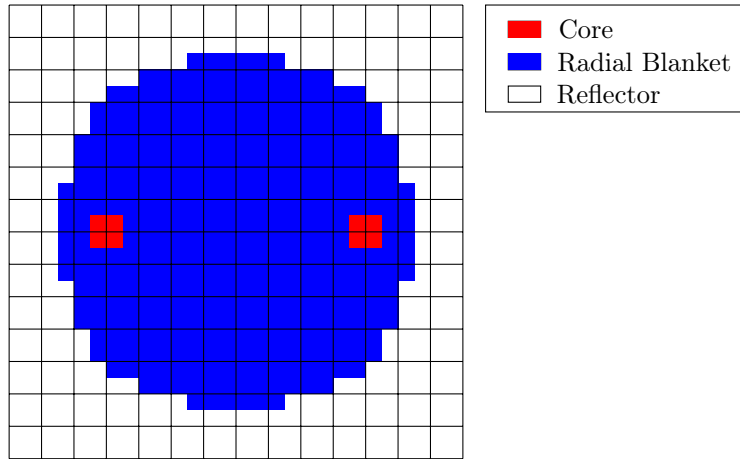
$$\begin{aligned} \gamma_h &\leq Ch^\mu, \\ \gamma_{*,h} &\leq Ch^\mu. \end{aligned}$$

This concludes the proof. □

### 7. Numerical Results

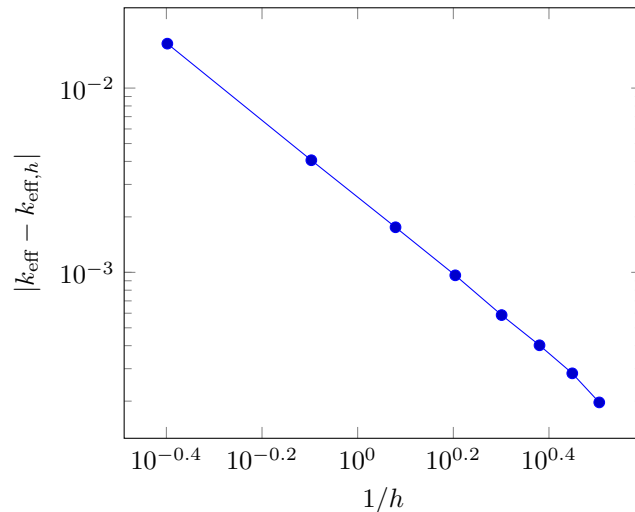
We consider the test case Model 2, case 1 from the benchmark of Takeda and Ikeda [20]. The geometry of the core is three-dimensional and the domain is  $\{(x, y, z) \in \mathbb{R}^3, 0 \leq x \leq 140 \text{ cm}; 0 \leq y \leq 140 \text{ cm}; 0 \leq z \leq 150 \text{ cm}\}$ . This test is defined with 4 energy groups, isotropic scattering and vacuum boundary conditions. Figure 1 represents the cross-sectional geometry on the plane  $z = 75 \text{ cm}$ .

Since the scattering is isotropic, the  $SP_3$  formulation can easily be reformulated as a multi-group diffusion problem with 8 energy groups and an isotropic albedo boundary condition [3]. We then made the computations with the PRIAM solver from the code CRONOS2 [14] for the conforming case and with the MINARET solver [15] from the APOLLO3<sup>®</sup> code [19] for the SIPG discretization.



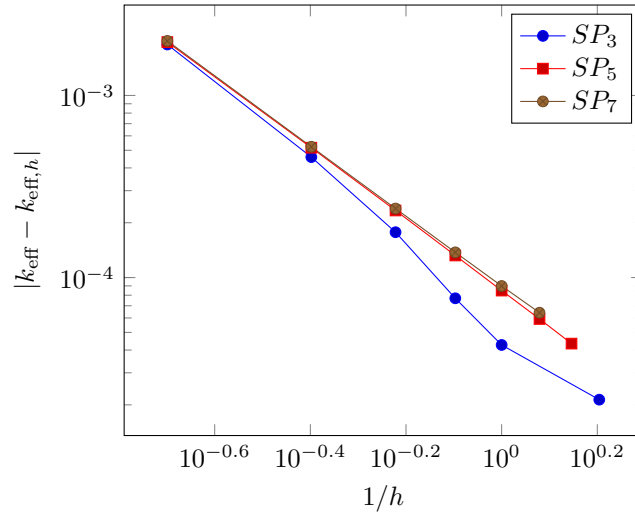
**Figure 1.** Cross-sectional view of the core ( $z = 75$  cm).

In Figure 2, we consider the convergence of the fundamental mode where we used the  $SP_3$  formulation with  $Q^1$  finite elements and a regular cartesian mesh of size  $h$ . The approximated order of convergence is 2.22.



**Figure 2.** Error on the discrete eigenvalue for the  $SP_3$  formulation with  $Q^1$  finite elements

In Figure 3, we consider the convergence of the fundamental mode for different the  $SP_N$  formulations with discontinuous  $P^1$  finite elements and a prismatic mesh of size  $h$ . The approximated orders of convergence are given in Table 1.



**Figure 3.** Error on the discrete eigenvalue for the  $SP_3$  formulation with discontinuous linear finite elements

**Table 1.** Approximated order of convergence associated with Figure 3

$SP_3$	$SP_5$	$SP_7$
1.88	1.96	1.92

## 8. Conclusion

We did the numerical analysis of the approximation with an  $H^1$ -conforming finite element method of the neutron multigroup  $SP_N$  equations. We also studied the numerical analysis of the approximation with the Symmetric Interior Penalty Galerkin method of the neutron multigroup  $SP_N$  equations. We then illustrated numerically the convergence results on a benchmark representative of a nuclear core. Those results can be extended to a mixed finite element method, see [5] for the diffusion case with an  $H^1$ -conforming finite element method.

### Acknowledgements

The authors gratefully acknowledge P. Ciarlet for fruitful discussions.

### References

- [1] A. Alonso, A. D. Russo, "Spectral approximation of variationally-posed eigenvalue problems by nonconforming methods", *J. Comput. Appl. Math.* **223** (2009), no. 1, p. 177-197.
- [2] I. Babuška, J. E. Osborn, "Eigenvalue problems", in *Handbook of numerical analysis, vol. II*, Handbook of Numerical Analysis, vol. 2, North-Holland, 1991, p. 645-785.
- [3] A.-M. Baudron, J.-J. Lautard, "Simplified  $P_N$  transport core calculations in the Apollo3 system", International Conference on Mathematics and Computational Methods Applied to Nuclear Science and Engineering (M&C 2011).
- [4] H. Brezis, *Functional analysis, Sobolev spaces and partial differential equations*, Universitext, Springer, 2010.

- [5] P. Ciarlet, Jr., L. Giret, E. Jamelot, F. D. Kpadonou, “Numerical analysis of the mixed finite element method for the neutron diffusion eigenproblem with heterogeneous coefficients”, *ESAIM, Math. Model. Numer. Anal.* **52** (2018), no. 5, p. 2003-2035.
- [6] P. Ciarlet, Jr., E. Jamelot, F. D. Kpadonou, “Domain decomposition methods for the diffusion equation with low-regularity solution”, *Comput. Math. Appl.* **74** (2017), no. 10, p. 2369-2384.
- [7] P. G. Ciarlet, *Linear and nonlinear functional analysis with applications*, Society for Industrial and Applied Mathematics, 2013.
- [8] R. Dautray, J.-L. Lions, *Analyse mathématique et calcul numérique pour les sciences et les techniques*, Masson, 1985.
- [9] D. A. Di Pietro, A. Ern, *Mathematical aspects of discontinuous Galerkin methods*, Mathématiques & Applications, vol. 69, Springer, 2011.
- [10] J. J. Duderstadt, L. J. Hamilton, *Nuclear reactor analysis*, John Wiley & Sons, Inc., 1976.
- [11] A. Ern, J.-L. Guermond, *Theory and practice of finite elements*, Applied Mathematical Sciences, vol. 159, Springer, 2013.
- [12] E. M. Gelbard, “Application of spherical harmonics method to reactor problems”, 1960, Bettis Atomic Power Laboratory, West Mifflin, PA, Technical Report No. WAPD-BT-20.
- [13] L. Giret, “Non-conforming domain decomposition for the multigroup neutron SPN equation”, PhD Thesis, Paris Saclay, 2018.
- [14] J.-J. Lautard, S. Loubière, C. Fedon-Magnaud, “CRONOS: a modular computational system for neutronic core calculations”, 1992.
- [15] J.-J. Lautard, J.-Y. Moller, “Minaret, a deterministic neutron transport solver for nuclear core calculations”, International Conference on Mathematics and Computational Methods Applied to Nuclear Science and Engineering (M&C 2011).
- [16] G. I. Marchuk, V. I. Lebedev, *Numerical methods in the theory of neutron transport*, Harwood Academic Pub., 1986.
- [17] S. Nicaise, A.-M. Sändig, “General interface problems. I, II”, *Math. Methods Appl. Sci.* **17** (1994), no. 6, p. 395-429, 431-450.
- [18] J. E. Osborn, “Spectral approximation for compact operators”, *Math. Comp.* **29** (1975), no. 131, p. 712-725.
- [19] D. Schneider *et al.*, “APOLLO3<sup>®</sup>: CEA/DEN deterministic multi-purpose code for reactor physics analysis”, PHYSOR-2016, May 1-5 2016, Sun Valley, Idaho, USA.
- [20] T. Takeda, H. Ikeda, “3-D neutron transport benchmarks”, *Journal of Nuclear Science and Technology* **28** (1991), no. 7, p. 656-669.