

# Mathematical Analysis of a Method to Compute Guided Waves in Integrated Optics

D. Gómez Pedreira, Patrick Joly

► **To cite this version:**

D. Gómez Pedreira, Patrick Joly. Mathematical Analysis of a Method to Compute Guided Waves in Integrated Optics. [Research Report] RR-3933, INRIA. 2000. <inria-00072719>

**HAL Id: inria-00072719**

**<https://hal.inria.fr/inria-00072719>**

Submitted on 24 May 2006

**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.

*Mathematical analysis of a method to compute  
guided waves in integrated optics*

D. Gómez Pedreira and P. Joly

**N° 3933**

Mai 2000

————— THÈME 4 —————



*Rapport  
de recherche*



## Mathematical analysis of a method to compute guided waves in integrated optics

D. Gómez Pedreira and P. Joly\*

Thème 4 — Simulation et optimisation  
de systèmes complexes  
Projet Ondes

Rapport de recherche n° 3933 — Mai 2000 — 75 pages

**Abstract:** In this article, we propose a new method to solve an eigenvalue problem (posed in  $\mathbb{R}^2$ ) arising from the computation of guided modes in integrated optics electromagnetic waveguides under the weak guiding assumption.

We consider an open stratified waveguide translationally invariant in the infinite propagation direction. Its cross-section is also supposed to be an unbounded and stratified medium where an appropriate perturbation of the refraction index has been introduced to ensure the existence of guided modes.

The method presented here appears as a combination of analytical methods which take into account the unbounded and stratified character of the propagation medium and numerical computations which can be reduced to a neighborhood of the perturbation. In this report, we give a complete description of the method, present its main mathematical properties and achieve the convergence analysis with respect to the various approximation parameters,

**Key-words:** Guided modes, electromagnetic waveguides, integrated optics, spectral analysis, transparent boundary conditions, error estimates

*(Résumé : tsvp)*

\* dolores@zmat.usc.es, Patrick.Joly@inria.fr

# Analyse mathématique d'une méthode numérique pour le calcul d'ondes guidées en optique intégrée

**Résumé :** Dans ce rapport, nous proposons une nouvelle méthode pour résoudre un problème de valeurs propres en milieu non borné pour un opérateur différentiel du second ordre. Ce type de problème apparaît pour le calcul d'ondes électromagnétiques guidées en optique intégrée dans le cadre de l'hypothèse du faible guidage.

Nous considérons un guide ouvert invariant par translation dans une direction et occupant tout l'espace  $\mathbb{R}^3$ . La section transverse du guide apparaît comme une perturbation locale d'un milieu stratifié et on fait l'hypothèse que la distribution de l'indice de réfraction est telle que des modes guidés existent.

La méthode présentée ici apparaît comme une combinaison de méthodes analytiques pour la prise en compte du milieu non perturbé et de méthodes numériques pour la prise en compte de la perturbation. Nous décrivons cette méthode en détail, présentons ses principales propriétés mathématiques et menons l'analyse de la convergence par rapport aux divers paramètres d'approximation.

**Mots-clé :** Modes guidés, ondes électromagnétiques, optique intégrée, théorie spectrale, conditions aux limites transparentes, estimations d'erreur.

# Contents

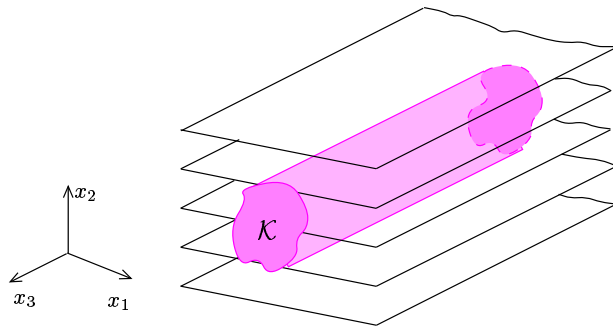
1	Introduction . . . . .	4
1.1	Mathematical setting . . . . .	7
1.2	Mathematical framework . . . . .	10
2	Presentation of the method . . . . .	15
2.1	Some notation . . . . .	15
2.2	A new formulation of the problem . . . . .	17
2.3	Decomposition of the operator $\tilde{S}$ . . . . .	19
2.4	Application to the computation of the operator $\tilde{S}$ . . . . .	20
3	Mathematical study of the problems $(\mathcal{P}_e)$ , $(\mathcal{P}_i)$ and $(\mathcal{P}_p)$ . . . . .	21
3.1	Study of the exterior problem $(\mathcal{P}_e)$ . . . . .	21
3.2	Study of the problem $(\mathcal{P}_i)$ . . . . .	22
3.3	Study of the problem $(\mathcal{P}_p)$ . . . . .	26
4	The operator $S$ and its spectral properties . . . . .	32
4.1	Definition of the operator $S$ . . . . .	32
4.2	Decomposition of the operator $S$ . . . . .	34
4.3	Smoothing and compactness properties of the operator $S_p$ . . . . .	37
4.4	Spectral properties of the operators $S_i$ , $S_e$ and $S$ . . . . .	40
4.5	Reformulation of the problem $(\mathcal{P}_S)$ . Introduction of the operator $K$ . . . . .	43
5	Numerical approximation . . . . .	45
5.1	Introduction of the parameter $R$ . The truncation of the domain $\Gamma$ . . . . .	45
5.2	Introduction of the parameter $N$ . The series truncation . . . . .	47
5.3	Description of the global numerical method . . . . .	49
6	Analysis of the error in the numerical approximation . . . . .	50
6.1	Analysis of the truncation related to $R$ . . . . .	51
6.2	Analysis of the truncation related to $N$ . . . . .	57

## 1 Introduction

An open electromagnetic waveguide consists, basically, of a dielectric structure which allows the electromagnetic energy to be confined inside, at least for certain frequencies depending on the geometry of the medium and on the refraction index  $n$  of the materials of which it is composed.

In the past few years, with the rapid growth of integrated optics communication techniques, the need for rigorous mathematical models describing in detail the guiding properties of open electromagnetic waveguides has also increased, because the design criteria which result from approximation methods often do not have the desired accuracy. The advantage of these devices versus the traditional optics components are their lighter weight and smaller size, lower cost and more stability. That is why the interest of the study.

In this work we shall be interested in open *stratified* electromagnetic waveguides. In this case, the device is composed by parallel layers each of them having a characteristic refraction index. The guide is supposed to be invariant with respect both to the geometry and to its physical characteristics under any translation along one privileged space direction, let us say  $x_3$ , which coincides with the propagation direction of the waves. In particular the refraction index will depend only on the two transverse coordinates  $x = (x_1, x_2)$  of the cross section of the guide.



**Figure 1.1:** Sketch of an open stratified waveguide and the choice of the coordinate system

If the materials composing the different layers are chosen with a suitable refraction index, then the energy of the wave is vertically confined in the material with the largest index. Nevertheless, to confine the light laterally—which is interesting for the design of these devices—it is necessary to introduce a compact perturbation  $\mathcal{K}$  where the refraction index depends on the two transverse coordinates, while in the

rest of the cross section depends only on the  $x_2$  coordinate. It is this perturbation which, if chosen appropriately, confines the waves inside a layer in a neighborhood of the perturbation. The dimensions of this perturbation are very small in comparison with the stratified medium and this is one of the reasons which leads to suppose, from the mathematical point of view, that the cross section of the guide is an unbounded domain.

For the sake of simplicity, we are going to consider a three-layer stratified medium where the perturbation of the refraction index involves only the central layer but, in fact, the method we outline in this work can be extended to an arbitrary number of parallel dielectric homogeneous layers, even if the perturbation is not completely embedded in one of them (cf. [13]).

Our aim will be to compute the *guided modes* (or *guided waves*) supported by such device. These are electromagnetic waves of finite transverse energy which can propagate without attenuation along the waveguide. The main difficulty to solve the problem comes from the unboundedness of the domain (the cross section of the guide) which makes the existence of guided modes not a priori obvious.

For the cases where the refractive index profile corresponds to an optical fiber or to a perturbed dioptr, theoretical studies of the problem have been carried out by Bonnet [3], Bamberger and Bonnet [2], Djellouli [7] or Gmati [12].

In our case, the originality of the study lies in the stratified character of the reference medium, which, at the same time, provides an additional difficulty.

In the case of a non perturbed stratified medium, the main references are the monographs by Wilcox [33] (for a scalar propagation model in acoustics) and Weder [32], which deals with a vectorial problem in electromagnetism.

In the case where the stratified medium is a perturbed medium, the more complete reference to our knowledge is the work by Bonnet *et al.* [4] who derive existence conditions for guided modes and bounds for the number of such guided modes. The main tool is the Min-Max principle, which allows characterizing the point spectrum below the essential spectrum. The analysis made in [4] shows that, contrary to what happens for optical fibers, the number of guided modes can remain bounded as the wave number tends to infinity.

From the numerical point of view, nothing has been done completely, to our knowledge, for the stratified perturbed case we are interested in, except an approximate method by Mahé [22].

There are some previous works for computing guided modes in the area of optical fibers but none of them can be applied to our situation due to the stratification of the media around the perturbation.



The idea in all of these works is to formulate a problem equivalent to the initial one but set in a bounded domain by introducing an artificial boundary which encloses the perturbation and to impose a transparent boundary condition:

- A first approach (cf. Johnson and Nedelec [16] or Jami and Lenoir [19] for scattering problems) consists in coupling finite elements and integral equations, taking into account an integral representation of the solution in the outside domain. The method requires the knowledge of the Green's function of the problem. This function is known analytically in a homogeneous media but is difficult to obtain in general stratified media (see Gmati [12] for a two layered medium). We can also refer Urbach [30] for a variant in volume of such an approach.
- Bonnet [3] for the scalar case and, more recently, Joly and Poirier [17] for the vectorial one propose, in the case of optical fibers, a “localized finite element method” (see also Picq [25], Masmoudi [23], Lenoir and Tounsi [20] or Givoli and Keller [11] for analogous works) consisting in reducing the domain of computation to a circle containing the perturbation. The solution of the exterior problem can be obtained by means of a Fourier expansion using separation of variables in polar coordinates. This allows to write an explicit “Dirichlet to Neumann” transparent condition. In our case there is no a separation of variables in polar coordinates because of the stratification of the medium.
- For a stratified reference medium as the one we are interested in, Mahé [22] propose an approximate method consisting in the introduction of two artificial horizontal boundaries where a Dirichlet or Neumann condition is applied. In a second step, the computational domain is bounded by considering two vertical boundaries enclosing the perturbation where exact transparent boundary conditions are imposed. The main drawback one can find in this method is that it is not exact: to get convergence to the true solution it is necessary to make the two horizontal boundaries go to infinity, which is very expensive numerically. It can also be expected that the accuracy of the results depends on frequency in the sense that this method will give very good results at high frequencies when the energy of the mode is very well confined in the vertical directions.

The method we propose in this work could be considered as an improvement of the Mahé's method in the sense that we give an approximation of the guided modes which converges to the exact ones for limit values of the approximation parameters

involved. The construction of the whole method extensively exploits the stratified nature of the reference medium. In particular, outside the perturbation, the two space variables  $x_1$  and  $x_2$  will be treated in a different way. Globally, it will appear as a mixed method combining three different techniques: Fourier transform, Fourier series and mixed finite elements.

The outline of this article is as follows: We begin with the statement of the problem in Section 1. We show how the study of guided modes under the weak guidance assumption amounts to the spectral analysis of a one-parameter dependent selfadjoint operator with non compact resolvent. The essential spectrum of this operator and bounds for its eigenvalues are determined.

In Section 2 we reformulate the problem and give a first description of the method we propose to compute the guided modes. We essentially insist on the ideas and do not consider the mathematical aspects that will be treated in detail in Section 3. The main point of the method involves the introduction of a selfadjoint operator  $S$  that we shall decompose, in view of numerical treatment, as the sum of three other operators. The main spectral properties of these operators will be the aim of Section 4.

In Section 5 we present the numerical approximation of the problem and finally, in Section 6, we derive the main theoretical results which justify this numerical approximation.

## 1.1 Mathematical setting

The propagation of electromagnetic waves is described by the Maxwell's equations

$$\begin{aligned} \mu \frac{\partial \mathbf{H}}{\partial t} + \operatorname{curl} \mathbf{E} &= 0 \\ \varepsilon \frac{\partial \mathbf{E}}{\partial t} - \operatorname{curl} \mathbf{H} &= 0 \end{aligned} \tag{1.1}$$

where  $\mathbf{E}(x, x_3, t)$  is the electric field,  $\mathbf{H}(x, x_3, t)$  is the magnetic field,  $t > 0$  denotes time and  $\varepsilon$  and  $\mu$  denote, respectively, the dielectric permittivity and the magnetic permeability of the material, specific characteristics which determine its electromagnetic behaviour. It is classical in guided optics to assume that  $\mu$  is constant and equal to  $\mu_0$  (the permeability of the vacuum), and then  $\varepsilon$  is related with the refraction index by the formula  $\varepsilon = \varepsilon_0 n^2$  ( $\varepsilon_0$  being the permittivity of the vacuum).

We will be interested in the computation of specific solutions of the Maxwell's equations called *guided modes*. By definition, a *guided mode* (or *guided wave*) is

solution of (1.1) in the form

$$\begin{aligned}\mathbf{E}(x, x_3, t) &= \mathbb{E}(x) e^{i(\omega t - \beta x_3)}, \\ \mathbf{H}(x, x_3, t) &= \mathbb{H}(x) e^{i(\omega t - \beta x_3)},\end{aligned}\tag{1.2}$$

where

- $\omega > 0$  is the *pulsation* of the wave,
- $\beta > 0$  is the *propagation constant* of the mode,
- $\mathbb{E}(x) = (E_1(x), E_2(x), E_3(x))$  and  $\mathbb{H}(x) = (H_1(x), H_2(x), H_3(x))$  are 2D vector fields describing the distribution of the electromagnetic field in each cross section and must satisfy

$$\int_{\mathbb{R}^2} (\epsilon |\mathbb{E}|^2 + \mu |\mathbb{H}|^2) dx < \infty.\tag{1.3}$$

The expression (1.2) represents an harmonic plane wave propagating without any attenuation ( $\omega$  and  $\beta$  are real numbers) in the  $x_3$ -direction, with *phase velocity*  $\omega/\beta$ . Such a solution is periodic in the direction  $x_3$  and the period  $\lambda = \frac{2\pi}{\beta}$  is called the *wavelength*.

The square integrability condition (1.3), which determines if a mode is guided or not, physically means that the energy of the mode remains practically confined in some bounded region of the cross section.

Our goal will be to compute numerically the frequencies for which guided modes can propagate in an open stratified waveguide. This means that the refraction index distribution  $n(x) : \mathbb{R}^2 \rightarrow \mathbb{R}_+^*$ , in addition to the usual hypotheses

$$\left| \begin{array}{l} n \in L^\infty(\mathbb{R}^2), \\ \inf_{(x_1, x_2) \in \mathbb{R}^2} n(x_1, x_2) > 0, \end{array} \right.\tag{1.4}$$

is such that there exists a compact set  $\mathcal{K} \in \mathbb{R}^2$  such that

$$n(x) = \bar{n}(x_2), \quad \forall x = (x_1, x_2) \notin \mathcal{K},$$

where  $\bar{n} \in L^\infty(\mathbb{R})$  is a positive function depending only on the  $x_2$  variable which represents the refraction index of the reference medium associated to the guide. We

suppose  $\bar{n}$  is given by

$$\bar{n}(x_2) = \begin{cases} n_\infty^- & \text{if } x_2 < 0 \\ n_0 & \text{if } 0 < x_2 < L \\ n_\infty^+ & \text{if } x_2 > L \end{cases}$$

The values  $n_\infty^-$  and  $n_\infty^+$  play symmetric roles, so we can assume without loss of generality that

$$n_\infty^- \leq n_\infty^+. \quad (1.5)$$

We also introduce:

$$\left\{ \begin{array}{l} n_+ = \sup_{x \in \mathbb{R}^2} n(x), \\ n_- = \inf_{x \in \mathbb{R}^2} n(x), \\ \bar{n}_+ = \sup_{x_2 \in \mathbb{R}} \bar{n}(x_2). \end{array} \right. \quad (1.6)$$

### The vectorial eigenvalue problem

By substituting the particular form of the guided modes into Maxwell's system and eliminating for instance  $\mathbb{H}$ , the problem to solve can be formulated as a vectorial eigenvalue problem in terms of the electric field

$$\text{curl}_\beta(\mu^{-1} \text{curl}_\beta(\mathbb{E})) = \omega^2 \varepsilon \mathbb{E}. \quad (1.7)$$

In (1.7)  $\text{curl}_\beta$  denotes the differential operator obtained from the curl by replacing  $\partial/\partial x_3$  by multiplication by  $-i\beta$ , i.e, if  $\vec{z} = (z_1, z_2, z_3)$ ,

$$\text{curl}_\beta \vec{z} = \left( \frac{\partial z_3}{\partial x_2} + i\beta z_2, -i\beta z_1 - \frac{\partial z_3}{\partial x_1}, \frac{\partial z_2}{\partial x_1} - \frac{\partial z_1}{\partial x_2} \right)^t.$$

### The simplified scalar model

We shall consider a particular case of problem (1.7) which arises under the assumption of *weak guidance* (that is, large wavenumber and weak variations of the refraction index). In this case (see, for instance, Snyder and Love [29], Bonnet [3],

Vassallo [31]), the third components of both magnetic and electric field can be neglected ( $E_3 \simeq 0 \simeq H_3$ ) and all the transverse components satisfy the same scalar equation for zero-order approximation of the Maxwell system, in such a way that searching of guided waves can be reduced to

$$(\mathcal{P}) \begin{cases} \text{Find } \omega > 0, \beta > 0 \text{ and } u \in L^2(\mathbb{R}^2) \text{ (} u \neq 0 \text{) such that} \\ -\Delta u + \beta^2 u = \omega^2 n^2 u, \\ \text{where } u(x) \equiv E_1(x), E_2(x), H_1(x) \text{ or } H_2(x). \end{cases} \quad (1.8)$$

For a given  $\beta > 0$ , this is a family of two dimensional scalar eigenvalue problems parameterized by  $\beta$ . The unknown  $\omega^2$  plays the role of eigenvalue of the operator  $\mathcal{A}_\beta = n^{-2}(-\Delta + \beta^2)$  and  $u$  is its associated eigenvector.

At this point, the first question is about the existence of guided modes or, in an equivalent manner, the existence of eigenvalues of the operator  $\mathcal{A}_\beta$ , which is not *a priori* obvious since  $\mathcal{A}_\beta$  has not compact resolvent.

This question is not the objective of our work. It has been exhaustively investigated by Bonnet *et al.* in [4] (see also [13] for complementary results). By using modern mathematical techniques related to spectral theory, they derive existence conditions for guided modes and bounds for the number of such guided modes which are related to conditions on the refraction index (not all perturbations lead to existence of such guided modes).

**Remark 1.1** *It is important to notice that guided modes, even when they exist, do not exist for just any values of  $\omega$  and  $\beta$  but for  $\omega$  and  $\beta$  linked by  $\omega = \omega(\beta)$ , the dispersion relation. The corresponding curves in the  $(\beta, \omega)$ -plane are by definition the dispersion curves.*

Thus we will be interested in the computation of the guided modes assuming they exist.

## 1.2 Mathematical framework

In this section, our aim is to provide a rigorous mathematical framework for the eigenvalue problem (1.8). We begin by giving a variational formulation and, after that, we rigorously define the operator  $\mathcal{A}_\beta$  as an unbounded operator on  $L^2(\mathbb{R}^2)$  associated with a symmetric bilinear form. We also establish the main properties of this bilinear form and of the operator  $\mathcal{A}_\beta$ . The main results are Theorems 1.1 and

1.2 which determine the essential spectrum of  $\mathcal{A}_\beta$  and give lower and upper bounds for the eigenvalues belonging to the discrete spectrum.

We shall consider that the Hilbert space  $L^2(\mathbb{R}^2)$  is endowed with the (weighted) inner product

$$(u, v) = \int_{\mathbb{R}^2} n^2 u v \, dx, \quad (1.9)$$

and denote by  $\|\cdot\|_{L^2(\mathbb{R}^2)}$  the corresponding Hilbert space norm.  $\|\cdot\|_{0,\mathbb{R}^2}$  will be denote the usual  $L^2$  norm. Then we can define on  $H^1(\mathbb{R}^2) \times H^1(\mathbb{R}^2)$  a symmetric bilinear form  $a(\beta; \cdot, \cdot)$  by

$$a(\beta; u, v) = \int_{\mathbb{R}^2} (\nabla u \cdot \nabla v + \beta^2 u v) \, dx. \quad (1.10)$$

By Green's formula, the *variational problem* associated to problem  $(\mathcal{P})$  states

$$\begin{cases} \text{For } \beta > 0, \text{ find } \omega > 0 \text{ and } u \in H^1(\mathbb{R}^2) \text{ (} u \neq 0 \text{) such that} \\ a(\beta; u, v) = \omega^2 (u, v), \quad \forall v \in H^1(\mathbb{R}^2). \end{cases} \quad (1.11)$$

To formulate our problem under a suitable form to use spectral theory, let us consider  $\mathcal{A}_\beta$  as an unbounded operator

$$\mathcal{A}_\beta : \mathcal{D}(\mathcal{A}_\beta) = H^2(\mathbb{R}^2) \subset L^2(\mathbb{R}^2) \longrightarrow L^2(\mathbb{R}^2)$$

given by

$$\mathcal{A}_\beta u = \frac{1}{n^2} (-\Delta u + \beta^2 u) \quad \forall u \in \mathcal{D}(\mathcal{A}_\beta).$$

Then a second formulation of our problem is

$$\begin{cases} \text{For } \beta > 0, \text{ find } \omega > 0 \text{ and } u \in \mathcal{D}(\mathcal{A}_\beta) \text{ (} u \neq 0 \text{) such that} \\ \mathcal{A}_\beta u = \omega^2 u, \end{cases} \quad (1.12)$$

which is the *spectral formulation* of the guided mode problem.

To study the eigenproblem (1.11) or (1.12) we need a rather detailed study of the bilinear form  $a(\beta; \cdot, \cdot)$  and of the operator  $\mathcal{A}_\beta$ .

**Lemma 1.1** *The bilinear form  $a(\beta; \cdot, \cdot)$  defined in (1.10) is continuous symmetric on  $H^1(\mathbb{R}^2) \times H^1(\mathbb{R}^2)$  and satisfies*

$$a(\beta; u, u) \geq \frac{1}{n_+^2} \left( \|\nabla u\|_{L^2(\mathbb{R}^2)}^2 + \beta^2 \|u\|_{L^2(\mathbb{R}^2)}^2 \right). \quad (1.13)$$

From previous classical characterizations of selfadjoint operators and the spectrum (cf. Schechter [28]) we deduce the following properties for the operator  $\mathcal{A}_\beta$ :

**Theorem 1.1** (i)  $\mathcal{A}_\beta$  is a selfadjoint and bounded from below operator.  
(ii) The spectrum  $\sigma(\mathcal{A}_\beta)$  of  $\mathcal{A}_\beta$  satisfies

$$\sigma(\mathcal{A}_\beta) \subset [\beta^2/n_+^2, +\infty).$$

As for any selfadjoint operator, the spectrum of  $\mathcal{A}_\beta$  splits into two different and complementary parts: a discrete set, the *discrete spectrum*  $\sigma_d(\mathcal{A}_\beta)$ , which is the set of isolated eigenvalues of finite multiplicity and a continuum, the *essential spectrum*  $\sigma_{ess}(\mathcal{A}_\beta)$ , which is the complementary set of the discrete spectrum in  $\sigma(\mathcal{A}_\beta)$ . We determine now the essential spectrum of  $\mathcal{A}_\beta$  while the discrete spectrum will be studied later.

### The essential spectrum of $\mathcal{A}_\beta$

The essential spectrum of  $\mathcal{A}_\beta$  can be characterized by a perturbation technique. We consider the operator  $\mathcal{A}_\beta$  as a perturbation of the operator  $\bar{\mathcal{A}}_\beta$

$$\bar{\mathcal{A}}_\beta : H^2(\mathbb{R}^2) \subset L^2(\mathbb{R}^2) \longrightarrow L^2(\mathbb{R}^2)$$

given by

$$\bar{\mathcal{A}}_\beta u = \frac{1}{\bar{n}^2} (-\Delta u + \beta^2 u) \tag{1.14}$$

associated with the non perturbed stratified medium obtained by replacing the function  $n(x)$  by  $\bar{n}(x_2)$ , the refraction index of the reference stratified medium.

By using partial Fourier transform with respect to the  $x_1$  variable, the spectral analysis of  $\bar{\mathcal{A}}_\beta$  can be easily reduced to the spectral analysis of a family of selfadjoint ordinary differential operators in the  $x_2$  variable, namely  $\bar{\mathcal{A}}_{\beta,k}$ , formally defined by ( $k \in \mathbb{R}$  is the dual variable of  $x_1$ )

$$\bar{\mathcal{A}}_{\beta,k} u = \frac{1}{\bar{n}^2} \left( -\frac{d^2 u}{dx_2^2} + (\beta^2 + k^2) u \right). \tag{1.15}$$

It is then straightforward to establish the following result

**Lemma 1.2** *The spectrum of the operator  $\bar{\mathcal{A}}_\beta$  is purely continuous and satisfies*

$$\sigma(\bar{\mathcal{A}}_\beta) = \sigma_{ess}(\bar{\mathcal{A}}_\beta) = [\sigma_e(\beta), +\infty), \quad (1.16)$$

where  $\sigma_e(\beta)$  is given by

$$\sigma_e(\beta) = \inf_{v \in H^1(\mathbb{R}) - \{0\}} \frac{\int_{\mathbb{R}} \left( \left| \frac{dv}{dx_2} \right|^2 + \beta^2 |v|^2 \right) dx_2}{\int_{\mathbb{R}} \bar{n}^2 |v|^2 dx_2}. \quad (1.17)$$

Then, by using compact perturbation techniques, it is not difficult to prove that the essential spectrum of  $\mathcal{A}_\beta$  is the same as the essential spectrum of  $\bar{\mathcal{A}}_\beta$ .

**Theorem 1.2** *The essential spectrum of the operator  $\mathcal{A}_\beta$  is given by*

$$\sigma_{ess}(\mathcal{A}_\beta) = [\sigma_e(\beta), +\infty)$$

where  $\sigma_e(\beta)$  has been defined in (1.17).

**Remark 1.2** *The lower bound  $\sigma_e(\beta)$  is nothing but the lower bound of the spectrum of the differential operator  $\bar{\mathcal{A}}_{\beta,0}$ . An immediate consequence of formula (1.17) is the inequality*

$$\sigma_e(\beta) \leq \frac{\beta^2}{n_\infty^{+2}}. \quad (1.18)$$

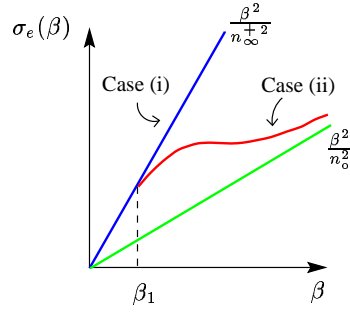
We can be more precise about  $\sigma_e(\beta)$  (see Fig. 1.2 and [13]):

- (i) When  $\bar{n}_+ = n_\infty^+$ , the spectrum of  $\bar{\mathcal{A}}_{\beta,0}$  is purely continuous and  $\sigma_e(\beta) = \beta^2/n_\infty^{+2}$ .
- (ii) If  $\bar{n}_+ > n_\infty^+$ , there exists a value  $\beta_1 > 0$ , called the first threshold of the reference stratified medium, such that

- for  $\beta < \beta_1$ ,  $\sigma_e(\beta) = \beta^2/n_\infty^{+2}$  (the spectrum of  $\bar{\mathcal{A}}_{\beta,0}$  is continuous),
- for  $\beta > \beta_1$ ,  $\sigma_e(\beta) < \beta^2/n_\infty^{+2}$ , and it is the smallest eigenvalue of  $\bar{\mathcal{A}}_{\beta,0}$ .

Moreover, one can prove the following behaviors





**Figure 1.2:** Cases for  $\sigma_\epsilon(\beta)$

$$\begin{aligned} \sigma_\epsilon(\beta) &\sim \beta^2/n_\infty^{+2} && \text{when } \beta \rightarrow 0 && (\text{see Proposition 2.8 in Mahé [22]}), \\ \sigma_\epsilon(\beta) &\sim \beta^2/n_0^2 && \text{when } \beta \rightarrow +\infty && (\text{see Lemma 1.9 in Mahé [22]}). \end{aligned}$$

In fact,  $\sigma_\epsilon(\beta)$  can be characterized as the solution of a simple transcendental equation and one shows that:

$$\sigma_\epsilon(\beta) = \beta^2/n_0^2 + r(\beta), \quad r(\beta) \rightarrow 0 \text{ exponentially when } \beta \rightarrow +\infty.$$

**The point spectrum of  $\mathcal{A}_\beta$**

The eigenvalues of  $\mathcal{A}_\beta$  can be divided into two categories (see Figure 1.2)

- (i) The ones which are strictly smaller than  $\sigma_\epsilon(\beta)$  which correspond to guided modes which propagate more slowly than any wave in the reference medium.
- (ii) The eigenvalues greater than  $\sigma_\epsilon(\beta)$  which are called embedded eigenvalues.



**Figure 1.3:** Sketch of the spectrum of the operator  $\mathcal{A}_\beta$ . The continuous line represents the points in the essential spectrum, the dotted line represents the points in the resolvent set and the cross signs  $\times$  represent the eigenvalues. There are two types of eigenvalues: those in the discrete spectrum (which are isolated points in the spectrum) and those in the essential spectrum.

The existence of both eigenvalues is not obvious and will depend on the nature of the perturbation  $n(x)$ . Eigenvalues in the discrete spectrum can be characterized

with the help of the min-max principle (we refer to Bonnet *et al.* [4] and G. Pedreira [13] for various existence results).

Although, in particular situations, eigenvalues embedded in the essential spectrum may exist (cf. Bonnet and Mahé [5]) there is a conjecture that the set of embedded eigenvalues is “generically” empty (with respect to the distribution  $n(x)$ ). Such modes are thus very unstable with respect to the physical imperfections of the guide.

That is why we shall be interested only in the eigenvalues corresponding to the discrete spectrum of  $\mathcal{A}_\beta$ . These eigenvalues continuously depend on  $n(x)$ . In what follows we shall call *guided mode* a mode corresponding to such an eigenvalue, and we shall look for eigenvalues  $\omega^2$  satisfying

$$\frac{\beta^2}{n_+^2} < \omega^2 < \sigma_e(\beta), \quad (1.19)$$

assuming that we are in a case where these eigenvalues may exist. More precisely, we are interested in computing the dispersion curves  $\beta \rightarrow \omega^2(\beta)$ , where  $\omega^2(\beta)$  is an eigenvalue of  $\mathcal{A}_\beta$ . Note that these functions are defined for  $\beta$  varying in an interval of the form  $]\beta_*, \beta^*[$ ,  $0 \leq \beta_* < \beta^* \leq +\infty$ , where  $\beta_*$  (respectively  $\beta^*$ ) is by definition a lower (respectively upper) threshold for the mode. These thresholds are such that

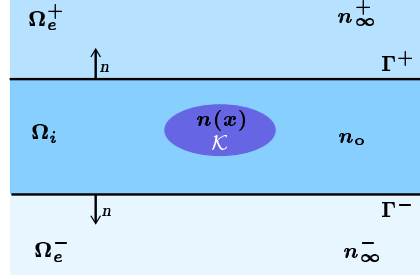
$$\lim_{\beta \rightarrow \beta_*} \omega(\beta)^2 = \sigma_e(\beta), \quad \lim_{\beta \rightarrow \beta^*} \omega(\beta)^2 = \sigma_e(\beta). \quad (1.20)$$

## 2 Presentation of the method

In this section, we make a sketch of the method proposed in this work for computing the guided modes. In order to lead the reader through the general framework in which development is taking place, we insist more on ideas than on rigorous mathematical aspects which will be developed in the following sections. Moreover, for the sake of simplicity of the exposition, we shall restrict ourselves to a simple model case that we present in the next paragraph.

### 2.1 Some notation

A cross section orthogonal to the  $x_3$  axis has been illustrated in Figure 2.1. We have denoted by  $\Omega_e^-$ ,  $\Omega_i$  and  $\Omega_e^+$  the different layers, with respective constant refractive index  $n_\infty^-$ ,  $n_0$  and  $n_\infty^+$ . The  $x_2$  axis is taken perpendicular to the interfaces between the layers, which coincide with the planes  $x_2 = 0$  and  $x_2 = L$ .



**Figure 2.4:** Cross section with three different layers.

More precisely, we shall use the following notations related to the geometry of the model domain:

- $\Gamma = \Gamma^- \cup \Gamma^+$ ,  $\Gamma^- = \{(x_1, 0), x_1 \in \mathbb{R}\}$ ,  $\Gamma^+ = \{(x_1, L), x_1 \in \mathbb{R}\}$ .
- $\Omega_e^- = \{(x_1, x_2) \in \mathbb{R}^2, x_2 < 0\}$ ,  $\Omega_e^+ = \{(x_1, x_2) \in \mathbb{R}^2, x_2 > L\}$ .
- $\Omega_e = \Omega_e^- \cup \Omega_e^+$ .
- $\Omega_i = \{(x_1, x_2) \in \mathbb{R}^2, 0 < x_2 < L\}$ .
- $\Omega = \Omega_i \cup \Omega_e = \mathbb{R}^2 \setminus \Gamma$ .
- $\mathcal{K}$  : the perturbed compact region,  $\mathcal{K} \subset \Omega_i$ .
- $n$  : the outer unit normal vector to  $\Gamma^-$  or  $\Gamma^+$ .

Finally,  $q$  being a function defined on  $\Omega$  with sufficient regularity,  $[q]_{\Gamma^+}$  and  $[q]_{\Gamma^-}$  will respectively denote the jump of function  $q$  across the boundaries  $\Gamma^+$  and  $\Gamma^-$  (we assume here that traces of  $q$  in  $\Gamma^+$  and  $\Gamma^-$  exist in some sense). More precisely

$$\begin{cases} [q]_{\Gamma^+} = (q|_{\Omega_i})|_{\Gamma^+} - (q|_{\Omega_e^+})|_{\Gamma^+}, \\ [q]_{\Gamma^-} = (q|_{\Omega_i})|_{\Gamma^-} - (q|_{\Omega_e^-})|_{\Gamma^-}. \end{cases} \quad (2.1)$$

## 2.2 A new formulation of the problem

According to inequality (1.19) we introduce the set

$$E = \{(\omega, \beta) \in \mathbb{R}^2 / \omega > 0, \beta > 0, \frac{\beta^2}{n_+^2} < \omega^2 < \sigma_e(\beta)\}. \quad (2.2)$$

Our goal is the resolution of the problem

$$(\mathcal{P}) \quad \begin{cases} \text{For } \beta > 0, \text{ find } \omega \text{ such that } (\omega, \beta) \in E, \text{ and } u \in L^2(\mathbb{R}^2) \text{ (} u \neq 0 \text{) satisfying} \\ -\Delta u + \beta^2 u = \omega^2 n^2 u. \end{cases} \quad (2.3)$$

The method we propose here consists in reducing the initial problem  $(\mathcal{P})$ , posed in  $\mathbb{R}^2$ , to another one posed only on the two artificial boundaries,  $\Gamma^+$  and  $\Gamma^-$  (see Fig. 2.1), and whose unknown  $\varphi$  will be the trace of  $u$  (the solution of  $(\mathcal{P})$ ) on  $\Gamma$ .

In the sequel, we adopt the notation  $H^{\frac{1}{2}}(\Gamma) = H^{\frac{1}{2}}(\Gamma^+) \times H^{\frac{1}{2}}(\Gamma^-)$ . We introduce an operator  $\tilde{S}(\omega, \beta)$  of Dirichlet-Neumann type depending on  $(\omega, \beta)$  defined as follows:

$$\begin{aligned} \tilde{S}(\omega, \beta) : H^{\frac{1}{2}}(\Gamma) &\longrightarrow H^{-\frac{1}{2}}(\Gamma), \\ \tilde{S}(\omega, \beta) \varphi &= \left[ \frac{\partial u(\varphi)}{\partial n} \right]_{\Gamma} = \left( \left[ \frac{\partial u(\varphi)}{\partial n} \right]_{\Gamma^+}, \left[ \frac{\partial u(\varphi)}{\partial n} \right]_{\Gamma^-} \right), \end{aligned}$$

where  $u(\varphi) \in H^1(\mathbb{R}^2)$  is the solution of the boundary value problem

$$(\mathcal{P}_{\varphi}) \quad \begin{cases} -\Delta u(\varphi) + (\beta^2 - n^2 \omega^2) u(\varphi) = 0 & \text{in } \Omega, \\ u(\varphi) = \varphi & \text{on } \Gamma. \end{cases} \quad (2.4)$$

The idea is the following: take a function  $\varphi$  defined on  $\Gamma$  and solve the boundary value problem  $(\mathcal{P}_{\varphi})$ , (which consists, in fact, of two decoupled problems: one outside the strip  $\Omega_i$  and another one inside). By construction, the function  $u(\varphi)$  is continuous. In order that  $u(\varphi)$  be a solution of problem  $(\mathcal{P})$ , it is enough to ensure the matching of the normal derivatives on the lines  $x_2 = 0$  and  $x_2 = L$ : this means that the the jump of the normal derivative of  $u(\varphi)$  across  $\Gamma$ , namely  $\tilde{S}(\omega, \beta)\varphi$ , must be equal to 0.

In this way, looking for guided modes is equivalent to solve the problem:

$(\mathcal{P}_{\tilde{S}})$  For  $\beta > 0$ , find  $\omega > 0$  with  $(\omega, \beta) \in E$ , and  $\varphi \neq 0, \varphi \in H^{\frac{1}{2}}(\Gamma)$  such that  $\tilde{S}(\omega, \beta) \varphi = 0$ .

In other words, for a given  $\beta$ , the values of  $\omega^2$  for which the operator  $\tilde{S}(\omega, \beta)$  is not injective, are the eigenvalues of the operator  $\mathcal{A}_\beta$  associated to problem  $(\mathcal{P})$ .

**Remark 2.1** *As we shall see in §4, the operator  $\tilde{S}(\omega, \beta)$  will not be defined for all pairs  $(\omega, \beta)$  in  $E$  but only for  $\omega \notin G_i(\beta)$ , where  $G_i(\beta)$  is a finite set of real numbers, to be defined later. This set can be assimilated to the famous irregular frequencies appearing in the solution of scattering problems by integral equations.*

**Remark 2.2** *Notice that if  $u$  is a solution of  $(\mathcal{P})$ ,  $u$  belongs to  $H^2(\mathbb{R}^2)$ , so that  $\varphi = u|_\Gamma$  belongs to  $H^{3/2}(\Gamma)$ . In fact, since  $\Gamma$  does not intersect  $\mathcal{K}$ , we even know that  $\varphi$  belongs to  $H^s(\Gamma)$ , for all  $s > 0$  (elliptic regularity).*

**Remark 2.3** *For theoretical purposes it will be useful (see §5) to introduce the restriction  $S(\omega, \beta)$  of the operator  $\tilde{S}(\omega, \beta)$  to  $H^1(\Gamma)$ , in such a way that  $S(\omega, \beta)$  will appear as an unbounded selfadjoint operator with domain  $\mathcal{D}(S) = H^1(\Gamma)$ . Then, the problem  $(\mathcal{P}_{\tilde{S}})$  will be reformulated as*

$(\mathcal{P}_S)$  For  $\beta > 0$ , find  $\omega > 0$  with  $(\omega, \beta) \in E$  such that 0 is an eigenvalue of  $S(\omega, \beta)$ .

Let us notice that:

- With the new formulation, we have replaced the problem originally posed in the whole domain  $\mathbb{R}^2$  by a new one posed in  $\Gamma$ .
- Doing so, we have added a new difficulty: instead of looking for  $\omega^2$  by directly solving an eigenvalue problem, now we look for values of  $\omega$  which make that the operator  $S$ , which depends on  $\omega$ , has 0 as an eigenvalue. This can be seen as an additional non linearity. This fact is rather general in numerical methods for solving open waveguide problems (see, for instance, Duterte [9]).

Of course, the efficiency of such method relies on an efficient way to evaluate numerically the operator  $\tilde{S}(\omega, \beta)$ . This is the aim of the next section.

### 2.3 Decomposition of the operator $\tilde{S}$

The construction of the operator  $\tilde{S}(\omega, \beta)$  requires, *a priori*, the resolution of the problem  $(\mathcal{P}_\varphi)$ . This problem cannot be directly solved because of the unboundedness of the domain. From both theoretical and numerical points of view, we shall represent  $\tilde{S}(\omega, \beta)$  via a decomposition into another three operators

$$\tilde{S}(\omega, \beta) = \tilde{S}_i(\omega, \beta) + \tilde{S}_p(\omega, \beta) - \tilde{S}_e(\omega, \beta) \quad (2.5)$$

that will be defined below. Let us write:

$$\left\{ \begin{array}{l} u(\varphi) = u_e(\varphi) \quad \text{in } \Omega_e, \\ u(\varphi) = u_i^p(\varphi) \quad \text{in } \Omega_i, \end{array} \right. \quad (2.6)$$

where

- $u_e = u_e(\varphi)$  is the unique solution in  $H^1(\Omega_e)$  of

$$(\mathcal{P}_e) \left\{ \begin{array}{l} -\Delta u_e + (\beta^2 - n_\infty^{\pm 2} \omega^2) u_e = 0 \quad \text{in } \Omega_e = \Omega_e^+ \cup \Omega_e^- \\ u_e = \varphi \quad \text{on } \Gamma \end{array} \right.$$

- $u_i^p = u_i^p(\varphi) = u_i(\varphi) + u_p(\varphi)$ , with  $u_i = u_i(\varphi)$  the unique solution of

$$(\mathcal{P}_i) \left\{ \begin{array}{l} -\Delta u_i + (\beta^2 - n_o^2 \omega^2) u_i = 0 \quad \text{in } \Omega_i \\ u_i = \varphi \quad \text{on } \Gamma \end{array} \right.$$

and  $u_p = u_p(\varphi)$  the unique solution, if  $\omega \notin G_i(\beta)$  (cf. Remark 2.1 and Theorem 3.5), of

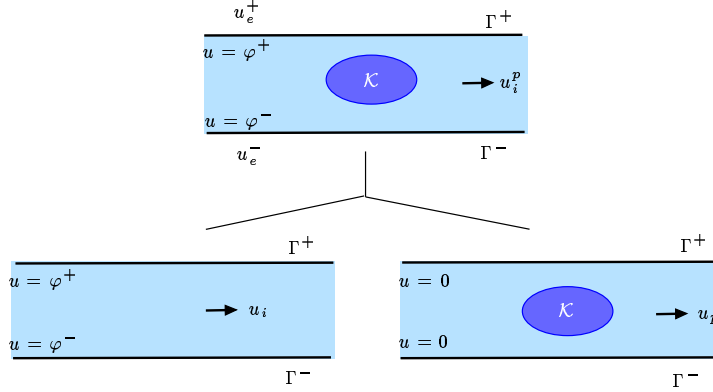
$$(\mathcal{P}_p) \left\{ \begin{array}{l} -\Delta u_p + (\beta^2 - n^2 \omega^2) u_p = (n^2 - n_o^2) \omega^2 u_i \quad \text{in } \Omega_i \\ u_p = 0 \quad \text{on } \Gamma \end{array} \right.$$

Then, we define the operator  $\tilde{S}_e(\omega, \beta)$  as

$$\tilde{S}_e(\omega, \beta) \varphi = \left( \frac{\partial u_e^+}{\partial n} \Big|_{\Gamma^+}, \frac{\partial u_e^-}{\partial n} \Big|_{\Gamma^-} \right), \quad (2.7)$$

and, similarly, the operators  $\tilde{S}_i(\omega, \beta)$  and  $\tilde{S}_p(\omega, \beta)$  as

$$\begin{aligned} \tilde{S}_i(\omega, \beta) \varphi &= \left( \frac{\partial u_i}{\partial n} \Big|_{\Gamma^+}, \frac{\partial u_i}{\partial n} \Big|_{\Gamma^-} \right), \\ \tilde{S}_p(\omega, \beta) \varphi &= \left( \frac{\partial u_p}{\partial n} \Big|_{\Gamma^+}, \frac{\partial u_p}{\partial n} \Big|_{\Gamma^-} \right). \end{aligned} \quad (2.8)$$



**Figure 2.5:** Computation of  $u_i^p = u_i + u_p$

Notice that  $u_e$  is the solution outside the strip, and  $u_i^p$  the solution in the perturbed layer;  $u_i$  is the solution inside in absence of perturbation and  $u_p$  the correction due to the perturbation (see Figure 2.3).

By construction

$$\left( \tilde{S}_i(\omega, \beta) + \tilde{S}_p(\omega, \beta) \right) \varphi = \left( \frac{\partial u_i^p}{\partial n} \Big|_{\Gamma^+}, \frac{\partial u_i^p}{\partial n} \Big|_{\Gamma^-} \right). \quad (2.9)$$

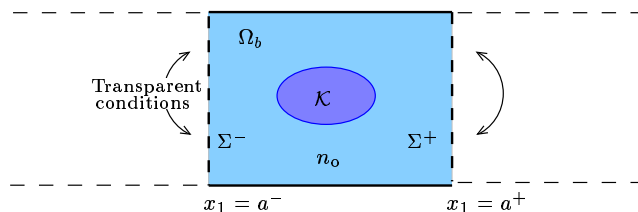
Then, by subtracting (2.7) from (2.9) we get formula (2.5).

## 2.4 Application to the computation of the operator $\tilde{S}$

For a given  $\varphi \in H^{\frac{1}{2}}(\Gamma)$ ,  $\tilde{S}(\omega, \beta)\varphi$  will be computed via the decomposition (2.5). This computation will involve a combination of analytical and numerical techniques. More precisely:

- Exploiting the fact that problems  $(\mathcal{P}_e)$  and  $(\mathcal{P}_i)$  are invariant under translation in the  $x_1$  direction, they can be solved explicitly using the partial Fourier transform in this direction. Therefore  $\tilde{S}_i(\omega, \beta)\varphi$  and  $\tilde{S}_e(\omega, \beta)\varphi$  are known explicitly via their Fourier transform.
- Since the boundary condition for problem  $(\mathcal{P}_p)$  is homogeneous and the right hand side has compact support in  $\Omega_i$ , we can reduce the solution of  $(\mathcal{P}_p)$  to a rectangular domain  $\Omega_b$  containing  $\mathcal{K}$ . This domain is delimited by boundaries

$\Gamma^+$  and  $\Gamma^-$  in the  $x_2$  direction and by two additional artificial boundaries, namely two “vertical” segments  $\Sigma^+$  and  $\Sigma^-$  (see Fig. 2.4) on which we apply appropriate transparent boundary conditions. This will be detailed in Section 3.3.



**Figure 2.6:** Bounded domain  $\Omega_b$  and transparent boundary conditions.

### 3 Mathematical study of the problems $(\mathcal{P}_e)$ , $(\mathcal{P}_i)$ and $(\mathcal{P}_p)$

In what follows, if  $\mathcal{O}$  be an open subset of  $\mathbb{R}^n$ , we shall denote

$$H^1(\Delta, \mathcal{O}) = \{v \in H^1(\mathcal{O}); \Delta v \in L^2(\mathcal{O})\} \tag{3.1}$$

endowed with the norm:

$$\|v\|_{H(\Delta, \mathcal{O})} = (\|v\|_{H^1(\mathcal{O})}^2 + \|\Delta v\|_{L^2(\mathcal{O})}^2)^{1/2}.$$

Besides, if  $\mathcal{O}$  has Lipschitz continuous boundary, we have the Green’s formula (cf. Dautray and Lions [6] or Girault and Raviart [10])

$$\left\langle \frac{\partial u}{\partial n} \Big|_{\partial \mathcal{O}}, \bar{v} \right\rangle_{\partial \mathcal{O}} = \int_{\mathcal{O}} \Delta u \bar{v} \, dx + \int_{\mathcal{O}} \nabla u \cdot \nabla \bar{v} \, dx, \quad \forall u \in H(\Delta, \mathcal{O}), v \in H^1(\mathcal{O}). \tag{3.2}$$

#### 3.1 Study of the exterior problem $(\mathcal{P}_e)$

First note that  $\omega^2 < \sigma_e(\beta)$ , coupled with inequality (1.18) and (1.5), implies that

$$\begin{cases} \beta^2 - n_{\infty}^{+2} \omega^2 > 0, \\ \beta^2 - n_{\infty}^{-2} \omega^2 > 0. \end{cases} \tag{3.3}$$

Then we have the following classical result, which is a direct consequence of the Lax-Milgram Lemma and of well-known estimates for Dirichlet’s non homogeneous problems for the Laplace operator (see, for instance, Girault and Raviart [10]).



**Theorem 3.1** *Problem  $(\mathcal{P}_e)$  has a unique solution  $u_e$  belonging to  $H^1(\Delta, \Omega_e)$ . Besides, we have the following estimate*

$$\|u_e\|_{H^1(\Delta, \Omega_e)} \leq C \|\varphi\|_{H^{\frac{1}{2}}(\Gamma)}. \quad (3.4)$$

As the refraction indices  $n_\infty^+$  and  $n_\infty^-$  are constant in  $\Omega_e^+$  and  $\Omega_e^-$ ,  $u_e$  can be calculated explicitly by using a Fourier transform in the  $x_1$  direction:

$$u_e(x_1, x_2) \xrightarrow{\mathcal{F}_{x_1}} \widehat{u}_e(k, x_2)$$

( $k \in \mathbb{R}$  denotes the dual variable of  $x_1$ ). A simple computation gives

$$\widehat{u}_e(k, x_2) = \begin{cases} \widehat{\varphi}^+(k) e^{-\xi_\infty^+(k)(x_2-L)} & \text{if } x_2 > L, \\ \widehat{\varphi}^-(k) e^{\xi_\infty^-(k)x_2} & \text{if } x_2 < 0, \end{cases} \quad (3.5)$$

where

$$\xi_\infty^+(k) = (k^2 + \beta^2 - \omega^2 n_\infty^{+2})^{1/2} > 0, \quad (3.6)$$

$$\xi_\infty^-(k) = (k^2 + \beta^2 - \omega^2 n_\infty^{-2})^{1/2} > 0. \quad (3.7)$$

Then, it is immediate to obtain an expression for the operator  $\widetilde{\mathcal{S}}_e(\omega, \beta)$  in the Fourier domain.

**Theorem 3.2** *For any  $\varphi \in H^{\frac{1}{2}}(\Gamma)$ , we have*

$$\left[ \widetilde{\mathcal{S}}_e(\omega, \beta) \varphi \right] (k) = M_e(k) \widehat{\varphi}(k), \quad (3.8)$$

where

$$\widehat{\varphi}(k) = \begin{pmatrix} \widehat{\varphi}^+(k) \\ \widehat{\varphi}^-(k) \end{pmatrix} \quad M_e(k) = \begin{pmatrix} -\xi_\infty^+(k) & 0 \\ 0 & -\xi_\infty^-(k) \end{pmatrix}. \quad (3.9)$$

### 3.2 Study of the problem $(\mathcal{P}_i)$

The existence and uniqueness of the solution is not as obvious as for  $(\mathcal{P}_e)$  since  $\beta^2 - n_o^2 \omega^2$  is not necessarily positive. However, the coerciveness of the problem is preserved thanks to the

**Lemma 3.1** *The following property holds*

$$\sigma_e(\beta) \leq \frac{1}{n_o^2} \left( \frac{\pi^2}{L^2} + \beta^2 \right). \quad (3.10)$$

*Proof.* We start from the formula

$$\inf_{\substack{u \in H_0^1(0, L) \\ u \neq 0}} \frac{\int_0^L \left( \left| \frac{du}{dx_2} \right|^2 + \beta^2 |u|^2 \right) dx}{\int_0^L n_o^2 |u|^2 dx} = \frac{1}{n_o^2} \left( \frac{\pi^2}{L^2} + \beta^2 \right). \quad (3.11)$$

For any  $\psi \in H_0^1(0, L)$  we consider its extension by 0 to  $\mathbb{R}$ , i.e., the function of  $H^1(\mathbb{R})$

$$\Psi(z) = \begin{cases} \psi(z) & \text{if } z \in (0, L), \\ 0 & \text{otherwise.} \end{cases}$$

We have

$$\int_{\mathbb{R}} \left( \left| \frac{d\Psi}{dx_2} \right|^2 + \beta^2 |\Psi|^2 \right) dx = \int_0^L \left( \left| \frac{d\psi}{dx_2} \right|^2 + \beta^2 |\psi|^2 \right) dx$$

and since  $\bar{n} = n_o$  in  $(0, L)$

$$\int_{\mathbb{R}} n^2 |\Psi|^2 dx = \int_0^L n_o^2 |\psi|^2 dx.$$

It is easy to conclude thanks to formula (1.17).  $\square$

**Theorem 3.3** *If  $\omega^2 < \sigma_e(\beta)$ , the bilinear form*

$$a_i(\omega, \beta; u, v) = \int_{\Omega_i} \nabla u \cdot \nabla v dx + (\beta^2 - n_o^2 \omega^2) \int_{\Omega_i} uv dx$$

*is coercive in  $H_0^1(\Omega_i)$ . As a consequence, problem  $(\mathcal{P}_i)$  has a unique solution  $u_i$  belonging to  $H^1(\Delta, \Omega_i)$ . Besides*

$$\|u_i\|_{H^1(\Delta, \Omega_i)} \leq C \|\varphi\|_{H^{\frac{1}{2}}(\Gamma)}. \quad (3.12)$$

*Proof.* Let  $u \in H_0^1(\Omega_i)$ . We can write

$$a_i(\omega, \beta; u, u) = \varepsilon \int_{\Omega_i} |\nabla u|^2 dx + (1 - \varepsilon) \int_{\Omega_i} |\nabla u|^2 dx + (\beta^2 - n_o^2 \omega^2) \int_{\Omega_i} |u|^2 dx,$$

where  $0 < \varepsilon < 1$  has to be properly chosen. From

$$\int_0^L \left| \frac{du}{dx_2} \right|^2 dx_2 \geq \frac{\pi^2}{L^2} \int_0^L |u|^2 dx_2 \quad (3.13)$$

we deduce

$$\int_{\Omega_i} |\nabla u|^2 dx \geq \frac{\pi^2}{L^2} \int_{\Omega_i} |u|^2 dx,$$

and then

$$a_i(\omega, \beta; u, v) \geq \varepsilon \int_{\Omega_i} |\nabla u|^2 dx + \left[ (1 - \varepsilon) \left( \frac{\pi^2}{L^2} + \beta^2 \right) - n_o^2 \omega^2 \right] \int_{\Omega_i} |u|^2 dx.$$

We choose now  $\varepsilon$  such that:

$$\left( \frac{\pi^2}{L^2} + \beta^2 \right) (1 - \varepsilon) > n_o^2 \omega^2$$

which is possible thanks to inequality (3.10) because  $\omega^2 < \sigma_e(\beta)$ . The coerciveness of  $a_i(\omega, \beta; \cdot, \cdot)$  deduces easily. Then the existence and uniqueness results follows immediately from Lax-Milgram's Lemma. Inequality (3.12) can be deduced in the same way as (3.4).  $\square$

In the same way as in the previous section,  $u_i$  can be computed by using Fourier transform in the  $x_1$  variable. Indeed, it is easy to check that

$$\widehat{u}_i(k, x_2) = \frac{\sinh(\xi_o(k) x_2)}{\sinh(\xi_o(k) L)} \widehat{\varphi}^+(k) + \frac{\sinh(\xi_o(k) (L - x_2))}{\sinh(\xi_o(k) L)} \widehat{\varphi}^-(k), \quad (3.14)$$

where

$$\xi_o(k) = \begin{cases} (k^2 + \beta^2 - n_o^2 \omega^2)^{1/2} & \text{if } k^2 + \beta^2 - n_o^2 \omega^2 \geq 0, \\ i(n_o^2 \omega^2 - k^2 - \beta^2)^{1/2} & \text{if } k^2 + \beta^2 - n_o^2 \omega^2 < 0. \end{cases} \quad (3.15)$$

**Remark 3.1** *The importance of the condition  $\omega^2 < \sigma_e(\beta)$  appears in formula (3.14) by expressing that  $\sinh(\xi_o(k)L)$  must be different from zero. Indeed, the equality could*

take place only for  $k^2 + \beta^2 - n_o^2 \omega^2 < 0$ , which occurs for some  $k$  when  $\beta < n_o \omega$ . This corresponds to the equation

$$\sin(\mu_o(k)L) = 0 \tag{3.16}$$

with  $\mu_o(k) = (n_o^2 \omega^2 - k^2 - \beta^2)^{1/2}$ . However, the condition  $\omega^2 < \sigma_e(\beta)$  combined with (3.10) implies

$$0 < \mu_o(k)L < \pi, \quad \forall k \in \mathbb{R},$$

so that equation (3.16) never holds.

We deduce from (3.14) an expression for the operator  $S_i$  in the Fourier domain.

**Theorem 3.4** For any  $\varphi \in H^{\frac{1}{2}}(\Gamma)$  we have

$$\left[ \widetilde{S}_i(\omega, \beta) \varphi \right] (k) = M_i(k) \widehat{\varphi}(k) \tag{3.17}$$

where  $\widehat{\varphi}(k)$  denotes the same vector as in (3.9) and

$$M_i(k) = \begin{pmatrix} \xi_o(k) \coth(\xi_o(k)L) & -\frac{\xi_o(k)}{\sinh(\xi_o(k)L)} \\ -\frac{\xi_o(k)}{\sinh(\xi_o(k)L)} & \xi_o(k) \coth(\xi_o(k)L) \end{pmatrix}. \tag{3.18}$$

We conclude this paragraph with a technical lemma which will be useful in the sequel.  $\langle \cdot, \cdot \rangle_\Gamma$  will denote in what follows the duality pairing between  $H^{-\frac{1}{2}}(\Gamma)$  and  $H^{\frac{1}{2}}(\Gamma)$ .

**Lemma 3.2** Let  $\varphi \in H^{\frac{1}{2}}(\Gamma)$ . The solution  $u_i$  of the problem  $(\mathcal{P}_i)$  satisfies

$$\|u_i\|_{H^{\frac{1}{2}}(\Omega_i)} \leq \|\varphi\|_{L^2(\Gamma)}. \tag{3.19}$$

*Proof.* The result could be proved using the explicit expression of  $u_i$ . We prefer to use a more direct transposition technique. Let  $G : H^2(\Omega_i) \cap H_0^1(\Omega_i) \rightarrow L^2(\Omega_i)$  the operator given by

$$Gz = -\Delta z + (\beta^2 - n_o^2 \omega^2)z.$$

Due to standard regularity results,  $G$  is an isomorphism. Multiplying problem  $(\mathcal{P}_i)$  by  $z \in H^2(\Omega_i) \cap H_0^1(\Omega_i)$ , integrating over  $\Omega_i$  and using Green's formula twice, we get

$$\int_{\Omega_i} u_i Gz = - \left\langle \varphi, \frac{\partial z}{\partial n} \right\rangle_{\Gamma}.$$

Choosing  $z = G^{-1}u_i$ , we get

$$\|u_i\|_{L^2(\Omega_i)}^2 = - \left\langle \varphi, \frac{\partial z}{\partial n} \right\rangle_{\Gamma} \leq \|\varphi\|_{H^{-\frac{1}{2}}(\Gamma)} \left\| \frac{\partial z}{\partial n} \right\|_{H^{\frac{1}{2}}(\Gamma)}.$$

Moreover, by trace theorem

$$\left\| \frac{\partial z}{\partial n} \right\|_{H^{\frac{1}{2}}(\Gamma)} < C \|z\|_{H^2(\Omega_i)} < C \|u_i\|_{L^2(\Omega_i)}.$$

Therefore

$$\|u_i\|_{L^2(\Omega_i)} \leq C \|\varphi\|_{H^{-\frac{1}{2}}(\Gamma)}. \quad (3.20)$$

On the other hand, we already know that

$$\|u_i\|_{H^1(\Omega_i)} \leq C \|\varphi\|_{H^{\frac{1}{2}}(\Gamma)}. \quad (3.21)$$

One concludes by interpolation (cf. Lions and Magenes [21]).  $\square$

### 3.3 Study of the problem $(\mathcal{P}_p)$

The existence and uniqueness of solution is not a priori obvious since the problem is no longer coercive in the general situation. In fact, we can distinguish two cases depending on the maximum of the refraction index: If the maximum of the refraction index is equal to the refraction index  $n_o$  of the central layer, then the problem is coercive and thus well-posed. Otherwise the problem is well-posed except for a certain number of singular frequencies  $\omega$  whose number increases with  $\beta$ , and which play the same role as the irregular frequencies which appear in integral equations for solving scattering problems.

To understand this, we introduce the operator  $\mathcal{A}_{p,\beta}$  of domain  $H^2(\Omega_i) \cap H_0^1(\Omega_i)$  defined by

$$\mathcal{A}_{p,\beta} u = \frac{1}{n^2} (-\Delta u + \beta^2 u).$$

The problem  $(\mathcal{P}_p)$  can be rewritten:

Find  $u_p \in H^2(\Omega_i)$  such that  $\mathcal{A}_{p,\beta}u_p - \omega^2 u_p = g$ , with  $g = \frac{(n^2 - n_o^2)\omega^2 u_i}{n^2} \in L^2(\Omega_i)$ .

Therefore, the problem  $(\mathcal{P}_p)$  is well-posed if and only if  $\omega^2$  does not belong to the spectrum of  $\mathcal{A}_{p,\beta}$ . Then Theorem 3.5, the main result of this section, follows immediately from Lemma 3.1 and from the following result, that we shall admit here and whose standard proof can be found in [13] for instance.

**Lemma 3.3**  $\mathcal{A}_{p,\beta}$  is a positive selfadjoint operator whose essential spectrum satisfies

$$\sigma_{ess}(\mathcal{A}_{p,\beta}) = \left[ \frac{1}{n_o^2} \left( \frac{\pi^2}{L^2} + \beta^2 \right), +\infty \right).$$

If  $n_+ = n_o$ , the discrete spectrum  $\sigma_d(\mathcal{A}_{p,\beta})$  of  $\mathcal{A}_{p,\beta}$  is empty. If  $n_+ > n_o$ , there exists an increasing sequence  $\beta_m^p$  tending to  $+\infty$  such that if  $\beta_m^p < \beta < \beta_{m+1}^p$ , then  $\sigma_d(\mathcal{A}_{p,\beta})$  consists of  $m$  eigenvalues, denoted  $G_i(\beta)$ , counted according to algebraic multiplicities.

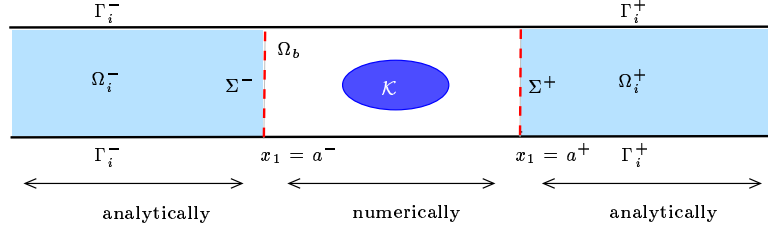
In what follows, we set  $G_i(\beta) = \emptyset$  if  $n_+ = n_o$ .

**Theorem 3.5** If  $\omega^2 < \sigma_e(\beta)$ , the problem  $(\mathcal{P}_p)$  is well-posed if and only if  $\omega^2 \notin G_i(\beta)$ . The solution  $u_p$ , when it exists, belongs to  $H^2(\Omega_i)$  and

$$\|u_p\|_{H^2(\Omega_i)} \leq C \|u_i\|_{L^2(\Omega_i)}. \tag{3.22}$$

The function  $u_p$  appears to be as the solution of a boundary value problem in the strip  $\Omega_i$  whose right hand side has compact support included in  $\mathcal{K}$ . As  $n$  depends on  $(x_1, x_2)$ , it is not possible to compute  $u_p$  analytically. But, since  $u_p$  verifies an homogeneous Dirichlet condition, one can compute it by using the so-called *Localized Finite Element Method* (according to the terminology of [20]). This consists in rewriting the problem  $(\mathcal{P}_p)$  as an equivalent one (in a certain sense) set in a bounded domain, with the help of transparent boundary conditions.

The idea is to introduce two fictitious vertical boundaries,  $\Sigma^+$  and  $\Sigma^-$ , to split the entire region  $\Omega_i$  into a bounded subdomain, namely the domain  $\Omega_b$ , containing the perturbation and another one,  $\Omega_i \setminus \Omega_b$ , where the refraction index is constant. Then two separate problems are considered: a problem in the subdomain  $\Omega_b$ , and another one in the subdomain  $\Omega_i \setminus \Omega_b$ , and we shall say that the whole problem in



**Figure 3.7:** Computation of  $u_p$ .

$\Omega_i$  is solved if the continuity across the artificial boundaries of  $u_p$  and its normal derivatives is ensured.

Therefore, the computation of  $u_p$  will be made by a mixed procedure: a numerical one inside  $\Omega_b$  and an analytical one outside (see Fig. 3.3).

Before giving more details, let us make precise some notations. We shall set

- $\Sigma^+ = \{(a^+, x_2) \in \mathbb{R}^2; x_2 \in [0, L]\}$ ,       $\Sigma^- = \{(a^-, x_2) \in \mathbb{R}^2; x_2 \in [0, L]\}$ ,
- $\Omega_i^+ = \{(x_1, x_2) \in \Omega_i; x_1 > a^+\}$ ,       $\Omega_i^- = \{(x_1, x_2) \in \Omega_i; x_1 < a^-\}$ ,
- $\Gamma_i^+ = \{(x_1, x_2) \in \Gamma; x_1 > a^+\}$ ,       $\Gamma_i^- = \{(x_1, x_2) \in \Gamma; x_1 < a^-\}$ .

In the sequel, it will be useful to identify  $\Sigma^+$  and  $\Sigma^-$  with the segment  $\Sigma = [0, L]$ . We shall denote by  $\{w_k\}$  the orthonormal basis of  $L^2(\Sigma)$  (or  $L^2(\Sigma^+)$  or  $L^2(\Sigma^-)$ ) given by

$$w_k(x_2) = \sqrt{\frac{2}{L}} \sin\left(\frac{k\pi x_2}{L}\right) \in H_0^1(\Sigma), \quad k \geq 1. \quad (3.23)$$

These are nothing but the eigenfunctions of the operator  $-d^2/dx_2^2$  in  $\Sigma$  with homogeneous boundary conditions. If  $\varphi$  is any function in  $H^{-1}(\Sigma)$  we shall set

$$\varphi_k = \langle \varphi, w_k \rangle_{\Sigma}$$

where  $\langle \cdot, \cdot \rangle_{\Sigma}$  denotes the duality product between  $H_0^1(\Sigma)$  and  $H^{-1}(\Sigma)$ . Of course if  $\varphi$  belongs to  $L^2(\Sigma)$ ,  $\varphi_k = (\varphi, w_k)_{L^2(\Sigma)}$ , and the  $\varphi_k$  are the expansion coefficients of  $\varphi$  in the basis  $\{w_k\}$ .

Let us recall that the Hilbert space  $H_{00}^{\frac{1}{2}}(\Sigma)$  is characterized by

$$H_{00}^{\frac{1}{2}}(\Sigma) = \{\varphi \in L^2(\Sigma) / \|\varphi\|_{H_{00}^{\frac{1}{2}}(\Sigma)}^2 = \sum_{k=1}^{\infty} |\varphi_k|^2 (1+k^2)^{1/2} < +\infty\} \quad (3.24)$$

while its topological dual, namely  $H_{00}^{\frac{1}{2}}(\Sigma)'$ , is characterized by

$$H_{00}^{\frac{1}{2}}(\Sigma)' = \{\varphi \in H^{-1}(\Sigma) / \|\varphi\|_{H_{00}^{\frac{1}{2}}(\Sigma)'}^2 = \sum_{k=1}^{\infty} |\varphi_k|^2 (1+k^2)^{-1/2} < +\infty\}. \quad (3.25)$$

The idea is first to compute the solution  $u_p$  in the exterior domain  $\Omega_i^+$  and  $\Omega_i^-$  assuming that its traces  $g^+$  and  $g^-$  on  $\Sigma^+$  and  $\Sigma^-$  are given. This leads us to consider the following problems, ( $g^+$  and  $g^-$  being given respectively in  $H_{00}^{\frac{1}{2}}(\Sigma^+)$  and  $H_{00}^{\frac{1}{2}}(\Sigma^-)$ )

$$\begin{aligned} (\mathcal{P}_i^+) \left\{ \begin{array}{l} -\Delta u_p + (\beta^2 - n_o^2 \omega^2) u_p = 0 \quad \text{in } \Omega_i^+, \\ u_p = 0 \quad \text{on } \Gamma_i^+, \\ u_p = g^+ \quad \text{on } \Sigma^+, \end{array} \right. \quad u_p \in H^1(\Omega_i^+), \\ \\ (\mathcal{P}_i^-) \left\{ \begin{array}{l} -\Delta u_p + (\beta^2 - n_o^2 \omega^2) u_p = 0 \quad \text{in } \Omega_i^-, \\ u_p = 0 \quad \text{on } \Gamma_i^-, \\ u_p = g^- \quad \text{on } \Sigma^-. \end{array} \right. \quad u_p \in H^1(\Omega_i^-), \end{aligned}$$

The following lemma concerns the definition of two operators  $T_+(\omega, \beta)$  and  $T_-(\omega, \beta)$  related to the transparent boundary conditions.

**Lemma 3.4** *Problem  $(\mathcal{P}_i^+)$  (respectively  $(\mathcal{P}_i^-)$ ) has a unique solution  $u^+$  (respectively  $u^-$ ) belonging to  $H^1(\Omega_i^+)$  (respectively  $H^1(\Omega_i^-)$ ) given by:*

$$\begin{aligned} u_p^+(x_1, x_2) &= \sum_{k=1}^{\infty} g_k^+ w_k(x_2) e^{-\xi_k(x_1 - a^+)} \quad \text{if } x_1 > a^+, \\ u_p^-(x_1, x_2) &= \sum_{k=1}^{\infty} g_k^- w_k(x_2) e^{-\xi_k(a^- - x_1)} \quad \text{if } x_1 < a^-, \end{aligned} \quad (3.26)$$

where  $\xi_k$  is given by

$$\xi_k = (k^2 \pi^2 / L^2 + \beta^2 - n_o^2 \omega^2)^{1/2}, \quad k \in \mathbb{N}. \quad (3.27)$$

*Proof.* The proof is straightforward. Formula (3.26) follows immediately from a classical technique of separation of variables.  $\square$



Using the trace Theorem 1.5.3.10 in Grisvard [15] applied to  $H^1(\Delta, \Omega_i^+)$  and  $H^1(\Delta, \Omega_i^-)$ , we can define the two operators  $T_+$  and  $T_-$  :

$$\begin{aligned} T_+(\omega, \beta) : H_{00}^{\frac{1}{2}}(\Sigma^+) &\longrightarrow H_{00}^{\frac{1}{2}}(\Sigma^+)' \\ g^+ &\longmapsto \frac{\partial u_p^+}{\partial \nu} \Big|_{\Sigma^+} = \frac{\partial u_p^+}{\partial x_1} \Big|_{\Sigma^+} \end{aligned} \quad (3.28)$$

$$\begin{aligned} T_-(\omega, \beta) : H_{00}^{\frac{1}{2}}(\Sigma^-) &\longrightarrow H_{00}^{\frac{1}{2}}(\Sigma^-)' \\ g^- &\longmapsto \frac{\partial u_p^-}{\partial \nu} \Big|_{\Sigma^-} = -\frac{\partial u_p^-}{\partial x_1} \Big|_{\Sigma^-} \end{aligned} \quad (3.29)$$

where  $u_p^+$  and  $u_p^-$  are the respective solutions of  $(\mathcal{P}_i^+)$  and  $(\mathcal{P}_i^-)$ .

Through the identification  $\Sigma^+ \equiv \Sigma^- \equiv \Sigma = [0, L]$ , we can identify  $T_+(\omega, \beta)$  and  $T_-(\omega, \beta)$  to a single operator

$$T(\omega, \beta) \in \mathcal{L}(H_{00}^{\frac{1}{2}}(\Sigma), H_{00}^{\frac{1}{2}}(\Sigma)')$$

Using Lemma 3.4 it is easy to prove the following theorem.

**Theorem 3.6** *The operator  $T(\omega, \beta)$  diagonalizes in the basis  $\{\mathbf{w}_k\}$  of  $H_{00}^{\frac{1}{2}}(\Sigma)$  and has the following expression:*

$$[T(\omega, \beta)\varphi](x_2) = \sum_{k=1}^{\infty} \xi_k(\omega, \beta) \mathbf{w}_k(x_2) \varphi_k \quad (3.30)$$

where  $\xi_k$  is defined by (3.27).

Let us denote  $u_p^b$  the restriction of  $u_p$  to the domain  $\Omega_b$ . By construction of  $T_+(\omega, \beta)$  and  $T_-(\omega, \beta)$  the continuity of  $u_p$  and of its normal derivatives across the two interfaces  $\Sigma^+$  and  $\Sigma^-$  shows that  $u_p^b$  satisfies the following boundary conditions:

$$\begin{aligned} \frac{\partial u_p^b}{\partial \nu} + T_+ u_p^b &= 0 \quad \text{on } \Sigma^+, \\ \frac{\partial u_p^b}{\partial \nu} + T_- u_p^b &= 0 \quad \text{on } \Sigma^-. \end{aligned}$$

This naturally lead us to introduce the following problem posed in the bounded domain  $\Omega_b$

$$(\mathcal{P}_p^\Sigma) \left\{ \begin{array}{ll} -\Delta u_p^b + (\beta^2 - n^2 \omega^2) u_p^b = (n^2 - n_0^2) \omega^2 u_i & \text{in } \Omega_b, \quad u_p^b \in H^1(\Omega_b) \\ u_p^b = 0 & \text{on } \Gamma_b, \\ \frac{\partial u_p^b}{\partial \nu} = -T_\pm u_p^b & \text{on } \Sigma^\pm, \end{array} \right. \quad (3.31)$$

where  $u_i$  is the solution of  $(\mathcal{P}_i)$  and we have set  $\Gamma_b = \Gamma \cap \partial\Omega_b$ . The problem  $(\mathcal{P}_p^\Sigma)$  is equivalent to the problem  $(\mathcal{P}_p)$  in the following sense (we omit the proof which is trivial):

**Theorem 3.7** (i) Let  $u_p$  be a solution of the problem  $(\mathcal{P}_p)$ . Then  $u_p^b := u_p|_{\Omega_b}$  is a solution of the problem  $(\mathcal{P}_p^\Sigma)$ .

(ii) Conversely, let  $u_p^b$  be a solution of  $(\mathcal{P}_p^\Sigma)$ . Then  $u_p^b$  can be extended to  $\Omega_i$  in a unique way yielding a solution  $u_p$  of  $(\mathcal{P}_p)$  as follows :

$$u_p(x_1, x_2) = \begin{cases} u_p^b(x_1, x_2) & \text{if } (x_1, x_2) \in \Omega_b, \\ \sum_{k=1}^{\infty} (u_p^b)_k^+ w_k(x_2) e^{-\xi_k(x_1 - a^+)} & \text{if } x_1 > a^+, \\ \sum_{k=1}^{\infty} (u_p^b)_k^- w_k(x_2) e^{-\xi_k(a^- - x_1)} & \text{if } x_1 < a^-, \end{cases} \quad (3.32)$$

where  $(u_p^b)_k^+$  and  $(u_p^b)_k^-$  are the expansion coefficients of the traces of  $u_p^b$  on  $\Sigma^+$  and  $\Sigma^-$  in the basis  $\{w_k\}$ .

**Remark 3.2** Toward this end, we shall denote  $f = (n^2 - n_0^2) \omega^2 u_i$ , where  $u_i$  is the solution of  $(\mathcal{P}_i)$ .

From Theorem 3.7 we deduce that the problem  $(\mathcal{P}_p^\Sigma)$  has unique solution  $u_p \in H^1(\Delta, \Omega_b)$  which satisfies

$$\|u_p\|_{H^1(\Omega_b)} \leq C(\omega, \beta) \|f\|_{L^2(\Omega_b)}, \quad (3.33)$$

$$\|\Delta u_p\|_{L^2(\Omega_b)} \leq C(\omega, \beta) \|f\|_{L^2(\Omega_b)}. \quad (3.34)$$

This theorem clearly shows that the computation of  $u_p$  has to be done in a mixed way, in the following sense:

- One computes  $u_p$  inside  $\Omega_b$  by solving the boundary value problem  $(\mathcal{P}_p^\Sigma)$ .
- Knowing  $u_p$  inside  $\Omega_b$  we compute it analytically in the exterior domain  $\Omega_i^+$  and  $\Omega_i^-$  via formula (3.32).

## 4 The operator $S$ and its spectral properties

In this section, we show that the operator  $\tilde{S}(\omega, \beta)$  (as well as  $\tilde{S}_e(\omega, \beta)$ ,  $\tilde{S}_i(\omega, \beta)$ ,  $\tilde{S}_p(\omega, \beta)$ ) can be reduced to an appropriate domain to be a selfadjoint operator, namely  $S(\omega, \beta)$ , in  $L^2(\Gamma)$ . Then we make the spectral analysis of  $S(\omega, \beta)$ .

### 4.1 Definition of the operator $S$

Let us introduce the bilinear form

$$s(\omega, \beta; \varphi, \psi) := \langle \tilde{S}(\omega, \beta)\varphi, \psi \rangle_\Gamma = \left\langle \left[ \frac{\partial u(\varphi)}{\partial \mathbf{n}} \right]_\Gamma, \psi \right\rangle_\Gamma \quad (4.1)$$

where  $u(\varphi)$  is the solution of  $(\mathcal{P}_\varphi)$  (see (2.4)). Thanks to Green's formula, it is easy to show that

$$s(\omega, \beta; \varphi, \psi) = \int_{\mathbb{R}^2} \nabla u(\varphi) \cdot \nabla u(\psi) + \int_{\mathbb{R}^2} (\beta^2 - n^2 \omega^2) u(\varphi) u(\psi) \quad \forall \varphi, \psi \in H^{\frac{1}{2}}(\Gamma), \quad (4.2)$$

which allows us to establish the following result

**Lemma 4.1** *The bilinear form  $s(\omega, \beta; \cdot, \cdot)$  is symmetric and satisfies the Gårding's inequality (where the constants  $C_1 > 0$  and  $C_2 \geq 0$  may depend on  $(\omega, \beta)$ )*

$$s(\omega, \beta; \varphi, \varphi) \geq C_1 \|\varphi\|_{H^{\frac{1}{2}}(\Gamma)}^2 - C_2 \|\varphi\|_{L^2(\Gamma)}^2, \quad \forall \varphi \in H^{\frac{1}{2}}(\Gamma).$$

*Proof.* The symmetry of  $s(\omega, \beta; \cdot, \cdot)$  is obvious from formula (4.2). Moreover we have

$$\begin{aligned} s(\omega, \beta; \varphi, \varphi) &\geq \min\{1, \beta^2\} \|u(\varphi)\|_{H^1(\mathbb{R}^2)}^2 - \omega^2 n_+^2 \|u(\varphi)\|_{0, \mathbb{R}^2}^2 \\ &\geq C_1 \|\varphi\|_{H^{\frac{1}{2}}(\Gamma)}^2 - \omega^2 n_+^2 \|u(\varphi)\|_{0, \mathbb{R}^2}^2 \end{aligned}$$

by trace theorem for  $u(\varphi)$ . To conclude, it is sufficient to prove that

$$\|u(\varphi)\|_{0,\mathbb{R}^2} \leq C \|\varphi\|_{L^2(\Gamma)}. \quad (4.3)$$

For this, we shall employ a duality argument. Let  $v(\varphi)$  be the solution of the problem (which exists because  $\omega^2 \notin G_i(\beta)$ )

$$\begin{cases} -\Delta v(\varphi) + (\beta^2 - n^2 \omega^2) v(\varphi) = u(\varphi) & \text{in } \Omega, \\ v(\varphi) = 0 & \text{on } \Gamma. \end{cases} \quad (4.4)$$

Since  $u(\varphi) \in L^2(\mathbb{R}^2)$ , then  $v(\varphi) \in H^2(\Omega)$  (elliptic regularity)

$$\|v(\varphi)\|_{H^2(\Omega)} \leq C \|u(\varphi)\|_{0,\mathbb{R}^2}. \quad (4.5)$$

By trace theorem, the function  $\left[ \frac{\partial v(\varphi)}{\partial n} \right]_{\Gamma}$  belongs to  $H^{\frac{1}{2}}(\Gamma)$  and using (4.5)

$$\left\| \left[ \frac{\partial v(\varphi)}{\partial n} \right] \right\|_{L^2(\Gamma)} \leq C \|u(\varphi)\|_{0,\mathbb{R}^2}. \quad (4.6)$$

Let us multiply (4.4) by  $u(\varphi)$ , integrate over  $\Omega$  and apply Green's formula in  $\Omega_i$  and  $\Omega_e$ . We get

$$\int_{\mathbb{R}^2} |u(\varphi)|^2 dx = \int_{\mathbb{R}^2} \nabla v(\varphi) \cdot \nabla u(\varphi) dx + \int_{\mathbb{R}^2} (\beta^2 - n^2 \omega^2) v(\varphi) u(\varphi) dx + \left\langle \left[ \frac{\partial v(\varphi)}{\partial n} \right], \varphi \right\rangle_{\Gamma}.$$

We apply once again Green's formula in  $\Omega$  to the right hand side. Taking into account the equation (2.4) satisfied by  $u(\varphi)$  and the fact that  $v(\varphi)$  vanishes on  $\Gamma$ , we obtain

$$\int_{\mathbb{R}^2} |u(\varphi)|^2 dx = \left\langle \left[ \frac{\partial v(\varphi)}{\partial n} \right], \varphi \right\rangle_{\Gamma}.$$

Then

$$\|u(\varphi)\|_{0,\mathbb{R}^2}^2 \leq \left\| \left[ \frac{\partial v(\varphi)}{\partial n} \right] \right\|_{L^2(\Gamma)} \|\varphi\|_{L^2(\Gamma)},$$

which together with (4.6) finishes the proof.  $\square$

As a consequence of previous lemma we have the

**Theorem 4.1** *The bilinear form  $s(\omega, \beta; \cdot, \cdot)$  defines a bounded from below selfadjoint operator  $S(\omega, \beta)$  in  $L^2(\Gamma)$  with domain*

$$\mathcal{D}(S(\omega, \beta)) = \{\varphi \in H^{\frac{1}{2}}(\Gamma) / \tilde{S}(\omega, \beta)\varphi \in L^2(\Gamma)\}, \quad (4.7)$$

that is to say,

$$(S(\omega, \beta) \phi, \psi)_{L^2(\Gamma)} = s(\omega, \beta; \phi, \psi), \quad \forall \phi, \psi \in \mathcal{D}(S(\omega, \beta)) \times H^{\frac{1}{2}}(\Gamma). \quad (4.8)$$

Of course,  $S(\omega, \beta)$  is also the restriction of  $\tilde{S}(\omega, \beta)$  to  $\mathcal{D}(S(\omega, \beta))$ . We shall prove later that  $\mathcal{D}(S(\omega, \beta)) = H^1(\Gamma)$ . This will use a decomposition of the operator  $S$  analogous to (2.5).

## 4.2 Decomposition of the operator $S$

Let us introduce the three bilinear forms on  $H^{\frac{1}{2}}(\Gamma)$  associated respectively to the operators  $\tilde{S}_e(\omega, \beta)$ ,  $\tilde{S}_i(\omega, \beta)$  and  $\tilde{S}_p(\omega, \beta)$ :

$$\begin{aligned} s_e(\omega, \beta; \varphi, \psi) &:= \langle \tilde{S}_e \varphi, \psi \rangle_{\Gamma} := \left\langle \frac{\partial u_e(\varphi)}{\partial n}, \psi \right\rangle_{\Gamma}, \\ s_i(\omega, \beta; \varphi, \psi) &:= \langle \tilde{S}_i \varphi, \psi \rangle_{\Gamma} := \left\langle \frac{\partial u_i(\varphi)}{\partial n}, \psi \right\rangle_{\Gamma}, \\ s_p(\omega, \beta; \varphi, \psi) &:= \langle \tilde{S}_p \varphi, \psi \rangle_{\Gamma} := \left\langle \frac{\partial u_p(\varphi)}{\partial n}, \psi \right\rangle_{\Gamma}, \end{aligned} \quad (4.9)$$

where  $u_e(\varphi)$ ,  $u_i(\varphi)$  and  $u_p(\varphi)$  are the respective solutions of problems  $(\mathcal{P}_e)$ ,  $(\mathcal{P}_i)$  and  $(\mathcal{P}_p)$  introduced in Section 3. Using Theorem 3.2 and Theorem 3.4 and Parseval's identity, we deduce the following expressions for  $s_e(\omega, \beta; \cdot, \cdot)$  and  $s_i(\omega, \beta; \cdot, \cdot)$

$$\begin{aligned} s_e(\omega, \beta; \varphi, \psi) &= \int_{\mathbb{R}} (M_e(k) \hat{\varphi}(k) \hat{\psi}(k)) dk, \\ s_i(\omega, \beta; \varphi, \psi) &= \int_{\mathbb{R}} (M_i(k) \hat{\varphi}(k) \hat{\psi}(k)) dk. \end{aligned} \quad (4.10)$$

Then, using standard arguments based on the fact that the hermitian matrices  $M_e(k)$  and  $M_i(k)$  depend analytically on  $k$  and satisfy

$$M_i(k) = |k|I + O\left(\frac{1}{|k|}\right), \quad M_e(k) = -|k|I + O\left(\frac{1}{|k|}\right), \quad |k| \longrightarrow +\infty \quad (4.11)$$

it is no difficult to prove the following result (we do not give the details but the same type of arguments will be used in Section 4.4, lemmas 4.4 and 4.5).

**Lemma 4.2** *The bilinear forms  $s_e(\omega, \beta; \cdot, \cdot)$ , and  $s_i(\omega, \beta; \cdot, \cdot)$  define two selfadjoint operators in  $L^2(\Gamma)$ ,  $S_e(\omega, \beta)$  and  $S_i(\omega, \beta)$ , with domain  $H^1(\Gamma)$  (these are restrictions to  $H^1(\Gamma)$  of  $\tilde{S}_e(\omega, \beta)$  and  $\tilde{S}_i(\omega, \beta)$ ). The spectrum of these operators is purely continuous.*

Contrary to the operators  $S_i$  and  $S_e$  which are unbounded in  $L^2(\Gamma)$ , we next prove that the operator  $\tilde{S}_p(\omega, \beta)$  can be uniquely extended as a bounded selfadjoint operator in  $L^2(\Gamma)$ . This will be a consequence of the following result about the bilinear form  $s_p(\omega, \beta; \cdot, \cdot)$ .

**Lemma 4.3** *The bilinear form  $s_p(\omega, \beta; \cdot, \cdot)$  associated with the operator  $\tilde{S}_p(\omega, \beta)$  satisfies:*

$$|s_p(\omega, \beta; \varphi, \psi)| \leq C \|\varphi\|_{L^2(\Gamma)} \|\psi\|_{L^2(\Gamma)} \quad \forall \varphi, \psi \in H^{\frac{1}{2}}(\Gamma). \quad (4.12)$$

*Proof.* Let us recall that by definition, for any  $\psi \in H^{\frac{1}{2}}(\Gamma)$ ,  $u_i(\psi)$  is the solution of the problem

$$\begin{cases} -\Delta u_i(\psi) + (\beta^2 - n_o^2 \omega^2) u_i(\psi) = 0 & \text{in } \Omega_i, \\ u_i(\psi) = \psi & \text{on } \Gamma. \end{cases} \quad (4.13)$$

Moreover, classical regularity results imply that

$$\|u_p(\psi)\|_{H^2(\Omega_i)} \leq C \|u_i(\psi)\|_{L^2(\Omega_i)}, \quad (4.14)$$

where, for any  $\varphi \in H^{\frac{1}{2}}(\Gamma)$ ,  $u_p(\varphi)$  denotes the solution of

$$\begin{cases} -\Delta u_p(\varphi) + (\beta^2 - n^2 \omega^2) u_p(\varphi) = (n^2 - n_o^2) \omega^2 u_i(\varphi) & \text{in } \Omega_i, \\ u_p(\varphi) = 0 & \text{on } \Gamma. \end{cases} \quad (4.15)$$

Let us multiply the first equation in (4.15) by  $u_i(\psi)$  and integrate over  $\Omega_i$ . Thanks to Green's formula we get

$$\begin{aligned} \left\langle \frac{\partial u_p(\varphi)}{\partial \mathbf{n}}, \psi \right\rangle_{\Gamma} + \int_{\Omega_i} (n^2 - n_o^2) \omega^2 u_i(\varphi) u_i(\psi) dx = \\ \int_{\Omega_i} \nabla u_p(\varphi) \cdot \nabla u_i(\psi) + \int_{\Omega_i} (\beta^2 - n^2 \omega^2) u_p(\varphi) u_i(\psi) dx. \end{aligned} \quad (4.16)$$

In the same way, from (4.13) we obtain

$$\int_{\Omega_i} \nabla u_p(\varphi) \cdot \nabla u_i(\psi) dx = - \int_{\Omega_i} (\beta^2 - n_o^2 \omega^2) u_p(\varphi) u_i(\psi) dx. \quad (4.17)$$

By substituting (4.17) into (4.16) we deduce, by using (4.9),

$$s_p(\omega, \beta; \varphi, \psi) = \int_{\Omega_i} (n_0^2 - n^2) \omega^2 u_p(\varphi) u_i(\psi) dx - \int_{\Omega_i} (n^2 - n_0^2) \omega^2 u_i(\varphi) u_i(\psi) dx.$$

Therefore, we obtain the estimate

$$\begin{aligned} |s_p(\omega, \beta; \varphi, \psi)| &\leq (n_0^2 - n_-^2) \omega^2 \int_{\Omega_i} |u_p(\varphi) u_i(\psi)| dx + (n_+^2 - n_0^2) \omega^2 \int_{\Omega_i} |u_i(\varphi) u_i(\psi)| dx \\ &\leq C_1 \|u_p(\varphi)\|_{L^2(\Omega_i)} \|u_i(\psi)\|_{L^2(\Omega_i)} + C_2 \|u_i(\varphi)\|_{L^2(\Omega_i)} \|u_i(\psi)\|_{L^2(\Omega_i)} \end{aligned}$$

and by using (4.14) with  $\psi = \varphi$

$$|s_p(\omega, \beta; \varphi, \psi)| \leq C \|u_i(\varphi)\|_{L^2(\Omega_i)} \|u_i(\psi)\|_{L^2(\Omega_i)}.$$

We conclude thanks to Lemma 3.2.  $\square$

As a corollary, we have:

**Theorem 4.2** *The operator  $S_p(\omega, \beta) : L^2(\Gamma) \longrightarrow L^2(\Gamma)$  defined by*

$$(S_p(\omega, \beta)\varphi, \psi) = s_p(\omega, \beta; \varphi, \psi)$$

*extends the operator  $\tilde{S}_p(\omega, \beta)$  to the domain  $L^2(\Gamma)$ .*

Since we have the identity

$$s(\omega, \beta; \phi, \psi) = s_i(\omega, \beta; \phi, \psi) - s_e(\omega, \beta; \phi, \psi) + s_p(\omega, \beta; \phi, \psi) \quad \forall \phi, \psi \in H^{\frac{1}{2}}(\Gamma) \quad (4.18)$$

we deduce immediately that

$$S(\omega, \beta) = S_i(\omega, \beta) - S_e(\omega, \beta) + S_p(\omega, \beta),$$

as an equality between unbounded operators in  $L^2(\Gamma)$ . This allows us to characterize the domain of the operator  $S(\omega, \beta)$ .

**Theorem 4.3** *The domain of the operator  $S(\omega, \beta)$  coincides with the set  $H^1(\Gamma)$ .*

*Proof.* We know that  $\tilde{S}(\omega, \beta) \in \mathcal{L}(\mathbf{H}^{\frac{1}{2}}(\Gamma), \mathbf{H}^{-\frac{1}{2}}(\Gamma))$ . Using regularity results for problem  $(\mathcal{P}_\varphi)$  (cf. Remark 2.2), we deduce  $\tilde{S}(\omega, \beta) \in \mathcal{L}(\mathbf{H}^{\frac{3}{2}}(\Gamma), \mathbf{H}^{\frac{1}{2}}(\Gamma))$ . Then, by interpolation (cf. Lions and Magenes [21]),

$$\tilde{S}(\omega, \beta) \in \mathcal{L}(\mathbf{H}^1(\Gamma), \mathbf{L}^2(\Gamma)),$$

which means that  $\mathbf{H}^1(\Gamma) \subset \mathcal{D}(S)$ .

To prove the reverse inclusion, let  $\phi \in \mathcal{D}(S(\omega, \beta))$ , which means that there exists a constant  $C(\phi)$  such that

$$|s(\omega, \beta; \phi, \psi)| \leq C(\phi) \|\psi\|_{\mathbf{L}^2(\Gamma)} \quad (4.19)$$

Using decomposition (4.18), Lemma 4.3 and identity (4.19) allows us to state, with another constant  $C(\phi)$

$$|(s_i - s_e)(\omega, \beta; \phi, \psi)| \leq C(\phi) \|\psi\|_{\mathbf{L}^2(\Gamma)}. \quad (4.20)$$

Equivalently, thanks to (4.10),

$$\int_{\mathbb{R}} (M_i - M_e)(k) \hat{\phi}(k) \hat{\psi}(k) dk \leq C(\phi) \|\psi\|_{\mathbf{L}^2(\Gamma)}. \quad (4.21)$$

From (4.11), we deduce that

$$M_i(k) - M_e(k) = 2|k|I + Q(k), \quad |Q(k)| \leq C, \forall k \in \mathbb{R}. \quad (4.22)$$

Therefore, with another constant  $C(\phi)$

$$\int_{\mathbb{R}} |k|(\hat{\phi}(k), \hat{\psi}(k)) dk \leq C(\phi) \|\psi\|_{\mathbf{L}^2(\Gamma)} \quad \forall \psi \in \mathbf{L}^2(\Gamma).$$

Since  $\mathbf{H}^{\frac{1}{2}}(\Gamma)$  is dense in  $\mathbf{L}^2(\Gamma)$ , we deduce that  $|k|\hat{\phi}(k)$  belongs to  $\mathbf{L}^2(\mathbb{R})$ , i.e.,  $\phi$  belongs to  $\mathbf{H}^1(\Gamma)$ .  $\square$

### 4.3 Smoothing and compactness properties of the operator $S_p$

In this section, we shall prove that the operator  $S_p(\omega, \beta)$  is a compact operator which maps the space  $\mathbf{L}^2$  into smooth rapidly decaying functions.

**Theorem 4.4** *The operator  $S_p(\omega, \beta)$  is a compact operator of  $\mathbf{L}^2(\Gamma)$ .*



*Proof.* It is enough to consider the chain

$$\begin{array}{ccccccc}
L^2(\Gamma) & \longrightarrow & L^2(\Omega_b) & \longrightarrow & H^2(\Omega_i) & \longrightarrow & H^{\frac{1}{2}}(\Gamma) \hookrightarrow L^2(\Gamma) \\
\varphi & \longmapsto & u_i(\varphi) & \longmapsto & u_p(\varphi) & \longmapsto & \frac{\partial u_p(\varphi)}{\partial n} = S_p \varphi|_{\Gamma} \\
& & \text{compact} & & \text{continuous} & & \text{continuous.}
\end{array}$$

where  $u_i(\varphi)$  is the solution of problem  $(\mathcal{P}_i)$  and  $u_p(\varphi)$  is the solution of  $(\mathcal{P}_p)$ .

Notice that, according to inequality (3.19), it makes sense to apply the operator  $S_p(\omega, \beta)$  to functions  $\varphi \in L^2(\Gamma)$ . Since the injection of  $H^{1/2}(\Omega_b)$  into  $L^2(\Omega_b)$  is compact ( $\Omega_b$  is a bounded domain, cf. Section 2.4), the compactness of the mapping  $\varphi \longrightarrow u_i(\varphi)$  follows immediately from the chain

$$\begin{array}{ccccccc}
L^2(\Gamma) & \longrightarrow & H^{\frac{1}{2}}(\Omega_i) & \longrightarrow & H^{\frac{1}{2}}(\Omega_b) & \longrightarrow & L^2(\Omega_b) \\
\varphi & \longmapsto & u_i(\varphi) & \longmapsto & u_i(\varphi)|_{\Omega_b} & \longmapsto & u_i(\varphi) \\
& & \text{continuous} & & \text{continuous} & & \text{compact}
\end{array}$$

where we have used Lemma 3.2. The continuity of the other ones follows from classical regularity results and a trace theorem.  $\square$

We give below a more precise result. Let  $m$  be a positive integer and  $\alpha > 0$ . We set

$$H^{m,\alpha}(\Gamma) = H^{m,\alpha}(\Gamma^+) \times H^{m,\alpha}(\Gamma^-),$$

where  $H^{m,\alpha}(\Gamma^+)$  and  $H^{m,\alpha}(\Gamma^-)$  are identified to  $H^{m,\alpha}(\mathbb{R})$ , defined by:

$$H^{m,\alpha}(\mathbb{R}) = \left\{ \varphi \in H^m(\mathbb{R}) \ / \ \sum_{j=0}^m \int_{\mathbb{R}} |D^j \varphi|^2 e^{\alpha|x|} dx < +\infty \right\} \quad (4.23)$$

which is a Hilbert space equipped with the norm

$$\|\varphi\|_{m,\alpha}^2 = \sum_{j=0}^m \int_{\mathbb{R}} |D^j \varphi|^2 e^{\alpha|x|} dx.$$

We have the following

**Theorem 4.5** *For all  $m \in \mathbb{N}$ ,  $S_p(\omega, \beta)$  can be extended as a bounded selfadjoint operator satisfying*

$$S_p \in \mathcal{L}(L^2(\Gamma), H^{m,\alpha}(\Gamma))$$

where  $\alpha < \alpha^* = 2 \left( \frac{\pi^2}{L^2} + \beta^2 - n_o^2 \omega^2 \right)^{1/2} > 0$ , ( $\alpha^* > 0$  deduces immediately from (1.19) and (3.10)).

*Proof.* We only sketch the proof and refer the reader to [13] for more details. Let  $\epsilon > 0$ , we set

$$\Omega_i^\epsilon(\Gamma) = \{x \in \Omega_i / d(x, \Gamma) < \epsilon\}.$$

For  $\epsilon$  small enough,

$$\mathcal{K} \cap \Omega_i^\epsilon(\Gamma) = \emptyset.$$

The first step of the proof consists in proving that  $u_p(\varphi)$  belongs to  $H^m(\Omega_i^\epsilon(\Gamma))$  with

$$\|u_p(\varphi)\|_{H^m(\Omega_i^\epsilon(\Gamma))} \leq C(m, \epsilon) \|\varphi\|_{L^2(\Gamma)}. \quad (4.24)$$

This is obvious for  $m = 2$  since we already know that  $u_p(\varphi)$  belongs to  $H^2(\Omega_i)$ . For greater  $m$ , it suffices to remark that in  $\Omega_i^\epsilon(\Gamma)$ ,  $u_p(\varphi)$  satisfies the homogeneous equation with constant coefficients:

$$-\Delta u_p(\varphi) + (\beta^2 - n_0^2 \omega^2) u_p(\varphi) = 0.$$

Then it suffices to use a localization procedure and a bootstrap argument to derive (4.24) (the fact that the boundary  $\Gamma$  is smooth and that  $u_p(\varphi)$  satisfies homogeneous Dirichlet conditions on  $\Gamma$  is also used in an essential way). The use of a trace theorem, allows us to deduce that  $S_p$  maps continuously  $L^2(\Gamma)$  into  $H^s(\Gamma)$ , for all  $s > 0$

$$\|u_p(\varphi)\|_{H^s(\Gamma)} \leq C(s, \epsilon) \|\varphi\|_{L^2(\Gamma)}. \quad (4.25)$$

The second step of the proof consists in obtaining weighted estimates of the type (this one concerns the trace on  $\Gamma^-$ , i.e, the line  $x_2 = 0$ )

$$\int_{a^+}^{\infty} \left| \frac{\partial^j}{\partial x_1^j} \left( \frac{\partial u_p}{\partial x_2} \right) (x_1, 0) \right|^2 e^{\alpha x_1} dx_1 < C \|\varphi\|_{L^2(\Gamma)}^2 \quad \forall j \in \mathbb{N}. \quad (4.26)$$

This part of the proof uses the expression of  $u_p$  for  $x_1 > a^+$  given in Theorem 3.7. We immediately get the following estimates for the trace of the normal derivative of  $u_p$  on  $x_2 = 0$ :

$$\left| \frac{\partial^j}{\partial x_1^j} \left( \frac{\partial u_p}{\partial x_2} \right) (x_1, 0) \right| \leq \sqrt{\frac{2}{L}} \frac{\pi}{L} \sum_{k=1}^{\infty} |u_k^+| k \xi_k^j e^{-\xi_k(x_1 - a^+)}, \quad \forall x_1 > a^+.$$

Since (see (3.27))  $\frac{\alpha^*}{2} \leq \xi_k \leq C(k^2 + 1)^{\frac{1}{2}}$ , we get

$$\left| \frac{\partial^j}{\partial x_1^j} \left( \frac{\partial u_p}{\partial x_2} \right) (x_1, 0) \right| \leq C e^{-\frac{\alpha^*}{2}(x_1 - a^+)} \sum_{k=1}^{\infty} |u_k^+| k^{j+1}.$$

But for  $m > \frac{j+2}{2}$ , we have

$$\sum_{k=1}^{\infty} |u_k^+| k^{j+1} \leq \left( \sum_{k=1}^{\infty} \frac{1}{k^2} \right)^{\frac{1}{2}} \left( \sum_{k=1}^{\infty} k^{2(2+j)} |u_k^+|^2 \right)^{\frac{1}{2}} \leq C \sum_{k=1}^{\infty} (k^{2m} u_k^+)^2. \quad (4.27)$$

Using the well-known result

$$\sum_{k=1}^{\infty} (k^{2m} u_k^+)^2 \leq C \|u_p\|_{\mathbb{H}^{2m}(\Sigma^+)}^2 \quad \forall m \in \mathbb{N}, \quad (4.28)$$

we deduce

$$\left| \frac{\partial^j}{\partial x_1^j} \left( \frac{\partial u_p}{\partial x_2} \right) (x_1, 0) \right|^2 \leq C e^{-\alpha^* (x_1 - a^+)} \|u_p\|_{\mathbb{H}^{2m}(\Sigma^+)}^2$$

and consequently, for  $\alpha < \alpha^*$  and  $m > \frac{j+2}{2}$ , we have

$$\int_{a^+}^{+\infty} \left| \frac{\partial^j}{\partial x_1^j} \left( \frac{\partial u_p}{\partial x_2} \right) (x_1, 0) \right|^2 e^{\alpha x_1} dx_1 \leq C \frac{e^{\alpha a^+}}{\alpha^* - \alpha} \|u_p\|_{\mathbb{H}^{2m}(\Sigma^+)}^2. \quad (4.29)$$

To conclude it is sufficient to notice that, in the same way one proves (4.25), one can show

$$\|u_p(\varphi)\|_{\mathbb{H}^s(\Sigma^+)} \leq C(s, \epsilon) \|\varphi\|_{L^2(\Gamma)}, \quad \forall s \in \mathbb{N}$$

because  $\Sigma^+ \cap \mathcal{K}$  is empty.  $\square$

**Remark 4.1** *The compactness of  $S_p$  also follows from Theorem 4.5 and the fact that the embedding*

$$\mathbb{H}^{m,\alpha}(\Gamma) \hookrightarrow L^2(\Gamma)$$

*is compact for  $m \geq 1$  and  $\alpha > 0$ .*

#### 4.4 Spectral properties of the operators $S_i$ , $S_e$ and $S$

In this section, we treat the operator  $S(\omega, \beta)$  as a compact perturbation of the operator  $(S_i - S_e)(\omega, \beta)$ . We first give two properties of this operator.

**Lemma 4.4** *The operator  $(S_i - S_e)(\omega, \beta)$  is an isomorphism from  $\mathbb{H}^1(\Gamma)$  into  $L^2(\Gamma)$ .*

*Proof.* It suffices to prove that

- (i) For all  $k \in \mathbb{R}$ ,  $(M_i - M_e)(k)$  is invertible.
- (ii) There exists  $C > 0$  such that  $|(M_i - M_e)^{-1}(k)| \leq C(1 + |k|^2)^{-1/2}$ .

To prove (i), we could use an explicit expression of  $(M_i - M_e)(k)$  but this leads to tedious computations. We prefer a more indirect proof based on the fact that (cf. Section 3)

$$(M_i - M_e)(k)\hat{\phi} = \left\{ - \left( \frac{d\hat{u}_e}{dx_2} - \frac{d\hat{u}_i}{dx_2} \right) (L), \left( \frac{d\hat{u}_e}{dx_2} - \frac{d\hat{u}_i}{dx_2} \right) (0) \right\}$$

where  $\hat{\phi} = (\hat{\phi}^+, \hat{\phi}^-) \in \ker[(M_i - M_e)(k)]$  and  $\hat{u}_i$  and  $\hat{u}_e$  are the respective solutions of

$$\left\{ \begin{array}{l} -\frac{d^2\hat{u}_i}{dx_2^2} + (k^2 + \beta^2 - n_o^2\omega^2)\hat{u}_i = 0 \quad \text{in } (0, L), \\ \hat{u}_i(L) = \hat{\phi}^+, \\ \hat{u}_i(0) = \hat{\phi}^-, \end{array} \right.$$

and

$$\left\{ \begin{array}{l} -\frac{d^2\hat{u}_e}{dx_2^2} + (k^2 + \beta^2 - n_\infty^{\pm 2}\omega^2)\hat{u}_e = 0 \quad \text{in } \mathbb{R} \setminus [0, L], \\ \hat{u}_e(L) = \hat{\phi}^+, \\ \hat{u}_e(0) = \hat{\phi}^-. \end{array} \right.$$

The equality  $[(M_i - M_e)(k)]\hat{\phi} = 0$  means that the function  $\hat{u}$  defined as:

$$\left\{ \begin{array}{l} \hat{u} = \hat{u}_i \quad \text{in } [0, L], \\ \hat{u} = \hat{u}_e \quad \text{in } \mathbb{R} \setminus [0, L], \end{array} \right.$$

satisfies

$$\left\{ \begin{array}{l} \hat{u} \in H^2((0, L) \cup (\mathbb{R} \setminus [0, L])), \\ [\hat{u}]_{x_2=0} = [\hat{u}]_{x_2=L} = 0, \\ \left[ \frac{d\hat{u}}{dx_2} \right]_{x_2=0} = \left[ \frac{d\hat{u}}{dx_2} \right]_{x_2=L} = 0, \end{array} \right.$$

which implies that  $\widehat{u} \in H^2(\mathbb{R})$  and satisfies

$$-\frac{d^2 \widehat{u}}{dx_2^2} + (k^2 + \beta^2 - \bar{n}^2 \omega^2) \widehat{u} = 0 \quad \text{in } \mathbb{R}.$$

If  $\widehat{u}$  were not identically 0, this would mean that  $\omega^2 \in \sigma(\bar{\mathcal{A}}_{\beta,k})$  ( $\bar{\mathcal{A}}_{\beta,k}$  being the operator introduced in (1.15)). This is not possible because

$$\omega^2 < \sigma_e(\beta) = \inf \sigma(\bar{\mathcal{A}}_{\beta,0}) = \inf_{\xi} \inf \sigma(\bar{\mathcal{A}}_{\beta,\xi}) \leq \inf \sigma(\bar{\mathcal{A}}_{\beta,k})$$

Therefore,  $\widehat{u} = 0$  which implies  $\phi = 0$ . To prove (ii) it suffices to remark that from (4.11) we deduce that

$$(M_i - M_e)^{-1}(k) = |2k|^{-1}I + O(|k|^{-3}).$$

□

Our next result characterizes the spectrum of the operator  $(S_i - S_e)(\omega, \beta)$ . The main consequence we shall derive is that 0 does not belong to the essential spectrum of this operator.

**Lemma 4.5** *The spectrum of  $(S_i - S_e)(\omega, \beta)$  is purely continuous and of the form:*

$$\sigma(S_i - S_e) = \sigma_{ess}(S_i - S_e) = [\sigma_*(\omega, \beta), +\infty)$$

with  $\sigma_*(\omega, \beta) > 0$ .

*Proof.* The fact that the spectrum of  $S_i - S_e$  is continuous derives from the fact that the matrix  $M_i - M_e$  depends analytically on  $k$ . Moreover,  $\sigma(S_i - S_e)$  is the set

$$\{\lambda_1(k), k \in \mathbb{R}\} \cup \{\lambda_2(k), k \in \mathbb{R}\}$$

where  $\lambda_1(k)$  and  $\lambda_2(k)$  are the two eigenvalues of  $(M_i - M_e)(k)$  which are continuous functions of  $k$  satisfying (cf. (4.22))

$$\lambda_i(k) \simeq 2|k|, \quad k \longrightarrow +\infty, \quad i = 1, 2.$$

As a consequence, the spectrum of  $S_i - S_e$  is an interval of the form  $[\sigma_*, +\infty)$ . The fact that  $\sigma_*(\omega, \beta) > 0$  is a consequence of Lemma 4.4 which indicates that 0 cannot belong to the spectrum of  $(S_i - S_e)(\omega, \beta)$ . □

As a consequence of the previous properties and the decomposition of the operator  $S$  we deduce the

**Lemma 4.6**  $\sigma_{ess}(S) = \sigma_{ess}(S_i - S_e) = [\sigma_*(\omega, \beta), +\infty)$ .

*Proof.* Since  $S_p(\omega, \beta)$  is a compact symmetric operator of  $L^2(\Gamma)$ ,  $S(\omega, \beta)$  is a compact perturbation of the selfadjoint operator  $(S_i - S_e)(\omega, \beta)$  and thus, as a consequence of the Weyl's theorem (cf. [27]), they have the same essential spectrum.  $\square$

This permits us to reformulate the characterization  $(\mathcal{P}_S)$  of our guided mode problem as follows

**Theorem 4.6** *There exists a guided mode associated to a pair  $(\omega, \beta)$  if and only if 0 is an isolated eigenvalue of finite multiplicity of the operator  $S(\omega, \beta)$ .*

A priori it seems we have not made a big step from  $(\mathcal{P}_S)$  to Theorem 4.6. However, Theorem 4.6 establishes an important property from the numerical point of view since isolated eigenvalues are those which can be easily distinguish from the others after discretization procedures.

#### 4.5 Reformulation of the problem $(\mathcal{P}_S)$ . Introduction of the operator $K$

One of the difficulties appearing in the numerical approximation of  $(\mathcal{P}_S)$  lies in the fact that the unknown function  $\varphi$  is defined on an unbounded domain, namely the boundary  $\Gamma$ . A possible approach would consist in a spectral approximation (in the physical space or in the Fourier domain) using appropriate bases of  $L^2(\mathbb{R})$  such as Hermite polynomials. We have preferred to study an alternative approach based on a truncation method which will be developed in the next section. To justify it, we need to work in the space  $L^2(\Gamma)$  instead of  $\mathcal{D}(S) = H^1(\Gamma)$ . This is our first motivation to reformulate our problem by introducing a new operator

$$K(\omega, \beta) = (S_i - S_e)^{-1} S_p, \tag{4.30}$$

which is well defined as a linear operator in  $L^2(\Gamma)$  since  $S_i - S_e$  is an isomorphism from  $H^1(\Gamma)$  into  $L^2(\Gamma)$ . The relationship between  $S(\omega, \beta)$  and  $K(\omega, \beta)$  is simply

$$S(\omega, \beta) = (S_i - S_e)\{I + K(\omega, \beta)\}, \tag{4.31}$$

so that it is obvious that the problem  $(\mathcal{P}_S)$  is equivalent to:

$(\mathcal{P}_K)$  For a given  $\beta$ , find  $\omega > 0$  with  $(\omega, \beta) \in E \setminus G_i(\beta)$ , such that  $-1$  be an eigenvalue of  $K(\omega, \beta)$

or, equivalently, there exists  $\varphi \in L^2(\Gamma)$ ,  $\varphi \neq 0$  such that

$$(I + K(\omega, \beta)) \varphi = 0. \quad (4.32)$$

Problem  $(\mathcal{P}_K)$  is the one we shall work with from the numerical point of view. The algorithm for searching guided modes will be:

1. Compute the eigenvalues  $\lambda(\omega, \beta)$  of  $K(\omega, \beta)$ .
2. Solve the equations  $\lambda(\omega, \beta) = -1$  in the plane  $(\omega, \beta)$ .

Let us summarize the most interesting properties of the operator  $K(\omega, \beta)$  in the following

**Theorem 4.7** *The operator  $K(\omega, \beta)$  is a compact operator from  $L^2(\Gamma)$  into  $L^2(\Gamma)$ . Except for 0, its spectrum is made of a sequence of non zero real eigenvalues having zero as unique accumulation point.*

*Proof.* The compactness of  $K$  is a consequence of the compactness of  $S_p$ . If we except 0, the general theory of compact operators says us that the rest of the spectrum is purely made of finite multiplicity eigenvalues having 0 as unique possible accumulation point. The fact that the spectrum is real comes from the fact that  $K$  can be written as the product of two selfadjoint operators. Indeed, let  $\lambda$  such an eigenvalue and  $\varphi \in L^2(\Gamma)$  the corresponding eigenfunction. Note that necessarily  $\varphi = \lambda^{-1}K\varphi \in H^1(\Gamma)$  ( $K$  maps  $L^2(\Gamma)$  into  $H^1(\Gamma)$  because of  $(S_i - S_e)^{-1}$ ) and that:

$$(I + K)\varphi = (\lambda + 1)\varphi.$$

Multiplying this equality by  $(S_i - S_e)$  we obtain

$$S\varphi = (1 + \lambda)(S_i - S_e)\varphi$$

which implies

$$(S\varphi, \varphi)_{L^2(\Gamma)} = (1 + \lambda)((S_i - S_e)\varphi, \varphi)_{L^2(\Gamma)}$$

which allows us to conclude because  $(S_i - S_e)\varphi, \varphi)_{L^2(\Gamma)}$  is real strictly positive (cf. Lemma 4.5).  $\square$

**Remark 4.2** *Note that  $K$  is not selfadjoint because  $S_p$  and  $(S_i - S_e)^{-1}$  do not commute.*

Apart from the interest of working in the space  $L^2(\Gamma)$ , the advantage is that  $K(\omega, \beta)$  no longer has continuous spectrum and that the theory of the approximation of the spectrum of compact operators, that we shall use in Section 6, is well-known (see Babuška and Osborn [1]).

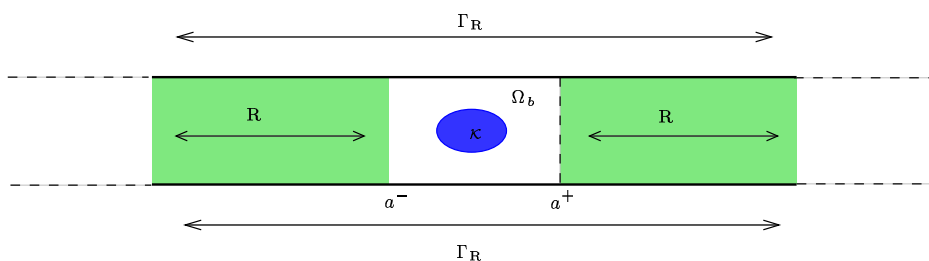
## 5 Numerical approximation

### 5.1 Introduction of the parameter $R$ . The truncation of the domain $\Gamma$

We propose a truncation procedure which leads to work with functions defined on the bounded domain (see Figure 5.8)

$$\Gamma_R = \Gamma \cap \{(x_1, x_2) / a^- - R < x_1 < a^+ + R\},$$

where  $R > 0$  is an approximation parameter devoted to tend to  $+\infty$ .



**Figure 5.8:** Truncation method.

We shall set  $\tilde{\Gamma}_R = \Gamma \setminus \Gamma_R$  and use the orthogonal decomposition

$$L^2(\Gamma) = L^2(\Gamma_R) \oplus L^2(\tilde{\Gamma}_R). \tag{5.1}$$

where  $L^2(\Gamma_R)$  (respectively  $L^2(\tilde{\Gamma}_R)$ ) denotes the subspace of functions in  $L^2(\Gamma)$  whose support is included in  $\Gamma_R$  (respectively  $\tilde{\Gamma}_R$ ).

We also introduce the orthogonal projector  $\Pi_R$  on  $L^2(\Gamma_R)$  defined by

$$\begin{aligned} \Pi_R : L^2(\Gamma) &\longrightarrow L^2(\Gamma) \\ \varphi &\longrightarrow \Pi_R \varphi = \chi_R \varphi \end{aligned} \tag{5.2}$$

where  $\chi_R$  denotes the characteristic function of  $\Gamma_R$ .



Let us start from equation (4.32). The idea would be to write an equation for  $\Pi_R \varphi$ , the “restriction” of  $\varphi$  to  $\Gamma_R$ . Such an equation does not exist but we are going to see that we can write an approximate equation, up to exponentially small errors, whose unknown  $\varphi_R$  will be an approximation of  $\Pi_R \varphi$ .

Let us decompose  $\varphi$  as

$$\varphi = \Pi_R \varphi + (\varphi - \Pi_R \varphi).$$

Applying the operator  $\Pi_R$  to (4.32), we obtain

$$\Pi_R \varphi + \Pi_R K \Pi_R \varphi + \Pi_R K (I - \Pi_R) \varphi = 0$$

that we can rewrite as

$$\begin{aligned} \Pi_R \varphi &+ \Pi_R (S_i - S_e)^{-1} \Pi_R S_p \Pi_R \varphi \\ &+ \{ \Pi_R (S_i - S_e)^{-1} (I - \Pi_R) S_p \Pi_R \\ &+ \Pi_R (S_i - S_e)^{-1} S_p (I - \Pi_R) \} \varphi = 0. \end{aligned}$$

The following estimates will be proved in the next section ( $\gamma = \gamma(\omega, \beta) > 0$ )

$$\|S_p (I - \Pi_R)\|_{\mathcal{L}(L^2(\Gamma))} \leq C e^{-\gamma R}, \quad (5.3)$$

$$\|(I - \Pi_R) S_p\|_{\mathcal{L}(L^2(\Gamma))} \leq C e^{-\gamma R}. \quad (5.4)$$

This allows us to approximate the equation (4.32) by

$$\{I + K_R\} \varphi_R = 0, \quad \varphi_R \in L^2(\Gamma_R) \quad (5.5)$$

where

$$K_R = \Pi_R (S_i - S_e)^{-1} \Pi_R S_p \Pi_R, \quad (5.6)$$

and where  $\varphi_R$  will be an approximation of  $\Pi_R \varphi$ .

The most interesting properties of  $K_R$  are the following:

- (i) According to the orthogonal decomposition (5.1), the operator  $K_R$  has the following block structure

$$K_R = \begin{pmatrix} \mathbb{K}_R & 0 \\ 0 & 0 \end{pmatrix}, \quad \mathbb{K}_R \in \mathcal{L}(L^2(\Gamma_R)). \quad (5.7)$$

This shows that  $K_R$  has an infinity dimensional kernel, namely  $L^2(\tilde{\Gamma}_R)$  and that

$$\sigma(K_R) = \{0\} \cup \sigma(\mathbb{K}_R).$$

Moreover, identifying  $\varphi_R$  with its restriction to  $\Gamma_R$  equation (5.5) can also be rewritten as

$$\{I + \mathbb{K}_R\}\varphi_R = 0, \quad \varphi_R \in L^2(\Gamma_R). \quad (5.8)$$

In practice, we shall only have to work with the operator  $\mathbb{K}_R$ , that is to say, in the bounded domain  $\Gamma_R$ .

- (ii)  $K_R$  appears to be the product of two selfadjoint and compact operators, namely  $\Pi_R S_p \Pi_R$  and  $\Pi_R (S_i - S_e)^{-1} \Pi_R$ . As  $S_p$  is compact (cf. Theorem 4.4),  $K_R$  is compact too. Moreover,  $\Pi_R (S_i - S_e)^{-1} \Pi_R$  is strictly positive when restricted to  $L^2(\Gamma_R)$ . This shows that the spectrum of  $K_R$  is made of 0 plus a sequence of real eigenvalues tending to 0. From the numerical point of view, looking for the eigenvalues of  $K_R$  will lead, after discretization, to a classical symmetric generalized eigenvalue problem.

## 5.2 Introduction of the parameter $N$ . The series truncation

The second difficulty one meets in the numerical approximation of our problem and, more specifically, of problem  $(\mathcal{P}_p^\Sigma)$  is linked to the fact that the operator  $T(\omega, \beta)$  described in Theorem 3.6 is not well suited because it involves a series. In order to solve numerically the problem  $(\mathcal{P}_p^\Sigma)$  we are led to truncate it at rank  $N$  (see [20], [18] for similar questions and developments).

We shall denote  $T_N(\omega, \beta) \in \mathcal{L}(H_{00}^{\frac{1}{2}}(\Sigma), H_{00}^{\frac{1}{2}}(\Sigma)')$  the *truncated operator* defined as

$$[T_N(\omega, \beta)\varphi](x_2) = \sum_{k=1}^N \xi_k(\omega, \beta) \mathbf{w}_k(x_2) \varphi_k. \quad (5.9)$$

In the sequel, the following problem will be referred to as the *semidiscrete problem*:

$$(\mathcal{P}_p^\Sigma)_N \left\{ \begin{array}{ll} -\Delta u_p^N + (\beta^2 - n^2 \omega^2) u_p^N = (n^2 - n_o^2) \omega^2 u_i & \text{in } \Omega_b, \quad u_p^N \in H^1(\Omega_b) \\ u_p^N = 0 & \text{on } \Gamma_b, \\ \frac{\partial u_p^N}{\partial \nu} = -T_N u_p^N & \text{on } \Sigma^\pm. \end{array} \right. \quad (5.10)$$

Thus, the solution  $u_p^b$  of  $(\mathcal{P}_p^\Sigma)$  will be approximated by the solution  $u_p^N$  of the semidiscrete problem  $(\mathcal{P}_p^\Sigma)_N$ . The existence of solution of this problem, which is not a trivial question, will be proved in Theorem 6.7 of Section 6.2.

To compute the operator  $S_p\varphi$ , we need to know  $u_p$  outside  $\Omega_b$ . Consistently with Lemma 3.4 (or Theorem 3.7),  $u_p^N$  will be extended in the exterior domain  $\Omega_{ext} = \Omega_i \setminus \Omega_b$  by:

$$\begin{aligned} u_p^N(x_1, x_2) &= \sum_{k=1}^N (u_k^N)^+ w_k(x_2) e^{-\xi_k(x_1 - a^+)} \quad \text{if } x_1 > a^+, \\ u_p^N(x_1, x_2) &= \sum_{k=1}^N (u_k^N)^- w_k(x_2) e^{-\xi_k(a^- - x_1)} \quad \text{if } x_1 < a^-. \end{aligned} \tag{5.11}$$

Remark that  $u_p^N$  is defined in such a way that

$$-\Delta u_p^N + (\beta^2 - n_0^2 \omega^2) u_p^N = 0 \quad \text{in } \Omega_{ext},$$

and such that its normal derivative across the two boundaries  $\Sigma^\pm$  is continuous

$$\left[ \frac{\partial u_p^N}{\partial \nu} \right] \Big|_{\Sigma^\pm}$$

However, even though  $u_p^N$  is  $H^2$  in both domains  $\Omega_b$  and  $\Omega_{ext}$ , it is not continuous across  $\Sigma^\pm$ , more precisely

$$u_p^N \notin H^1(\Omega_i).$$

Nevertheless, this is sufficient to define an approximation  $S_p^N$  of the operator  $S_p$ . Let us denote

$$\Gamma_{ext} = \Gamma \setminus \Gamma_b = \Gamma_{ext}^+ \cup \Gamma_{ext}^-$$

with

$$\begin{aligned} \Gamma_{ext}^+ &= \Gamma_{ext} \cap \Gamma^+, \\ \Gamma_{ext}^- &= \Gamma_{ext} \cap \Gamma^-. \end{aligned}$$

The operator  $S_p^N$  is defined as follows:

$$\begin{cases} (S_p^N \varphi)|_{\Gamma_b} = \frac{\partial u_p^N}{\partial n} \Big|_{\Gamma_b}, \\ (S_p^N \varphi)|_{\Gamma_{ext}} = \frac{\partial u_p^N}{\partial n} \Big|_{\Gamma_{ext}}. \end{cases} \tag{5.12}$$

More precisely, we shall prove in Section 6.2 that, if  $\omega \notin G_i(\beta)$ , (5.12) defines, at least for  $N$  large enough,  $S_p^N$  as a bounded operator in  $L^2(\Gamma)$ . Moreover, according to formula (5.11), we have the following explicit representations of  $(S_p^N \varphi)|_{\Gamma_{ext}^+}$  :

$$\left\{ \begin{array}{l} (S_p^N \varphi)|_{\Gamma_{ext}^+} = \sum_{k=1}^N (-1)^k \sqrt{\frac{2}{L}} \frac{k\pi}{L} (u_k^N)^+ e^{-\xi_k(x_1 - a^+)} \quad \text{if } x_1 > a^+, \\ (S_p^N \varphi)|_{\Gamma_{ext}^+} = \sum_{k=1}^N (-1)^k \sqrt{\frac{2}{L}} \frac{k\pi}{L} (u_k^N)^- e^{-\xi_k(a^- - x_1)} \quad \text{if } x_1 < a^-, \end{array} \right. \quad (5.13)$$

and similar expressions on  $\Gamma_{ext}^-$ . In (5.13),  $(u_k^N)^+$  and  $(u_k^N)^-$  denotes, as usual, the expansion coefficients of the trace of  $u_p^N$  on  $\Sigma^+$  and  $\Sigma^-$  in the basis  $\{w_k\}$ .

This approximation of  $S_p$  leads naturally to introduce the approximation of  $K$ , namely  $K_N$  defined by

$$K_N(\omega, \beta) = (S_i - S_e)^{-1} S_p^N \in \mathcal{L}(L^2(\Gamma)).$$

**Remark 5.1** *It is easy to see that the operator  $K_N$  is compact. The compactness of  $K_N$  derives from the one of  $S_p^N$ , which can be easily proved as the compactness of  $S_p$  (see Theorem 4.4).*

### 5.3 Description of the global numerical method

In practice, we need to make approximations with respect to  $N$  and  $R$  and then to define the operator

$$K_R^N(\omega, \beta) = (\Pi_R (S_i - S_e)^{-1} \Pi_R) (\Pi_R S_p^N \Pi_R),$$

and consider the approximate problem

$$(\mathcal{P}_{K_R^N}) \quad \begin{array}{l} \text{For a given } \beta, \text{ find } \omega \in \mathbb{R} \text{ with } (\omega, \beta) \in E \setminus G_i(\beta), \\ \text{such that } -1 \text{ is an eigenvalue of } K_R^N(\omega, \beta). \end{array}$$

However, one still cannot handle the operator  $K_R^N$  numerically since the computation of  $S_p^N$  involves the resolution of the boundary value problem  $(\mathcal{P}_p^\Sigma)_N$ . This problem needs a numerical approximation, for instance, using a finite element method associated to a mesh of stepsize  $h$ . This approximation will produce a new

approximate operator  $S_p^{h,N}$  of  $S_p^N$ . We do not present in detail the step of the numerical method since it is quite standard.

In its turn, this leads to a new approximation of the operator  $K_R^N(\omega, \beta)$ , namely,

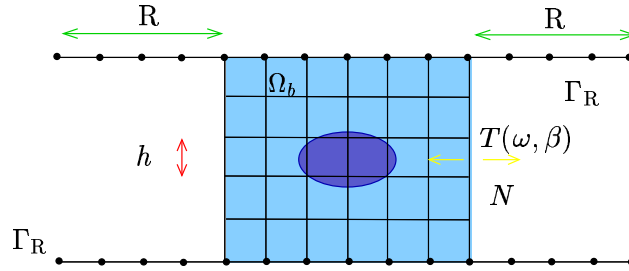
$$K_R^{h,N}(\omega, \beta) = (\Pi_R (S_i - S_e)^{-1} \Pi_R)_h (\Pi_R S_p^{h,N} \Pi_R),$$

and the numerical method we propose consists in solving the problem

$$(\mathcal{P}_{K_R^{h,N}}) \quad \text{For a given } \beta, \text{ find } \omega > 0 \text{ with } (\omega, \beta) \in E \setminus G_i(\beta),$$

$$\text{such that } -1 \text{ is an eigenvalue of } K_R^{h,N}(\omega, \beta).$$

We summarize the meaning of the different parameters in the Figure 5.9.



**Figure 5.9:** Reinterpretation of the numerical method.

Hence, after the numerical approximation, the problem of searching guided modes can be reduced to:

1. Compute the eigenvalues  $\lambda(\omega, \beta)$  of  $K_R^{h,N}(\omega, \beta)$ .
2. Solve the equations  $\lambda_R^{h,N}(\omega, \beta) = -1$  in the plane  $(\omega, \beta)$ .

## 6 Analysis of the error in the numerical approximation

In this section, we present an analysis of the numerical error due to the approximations related with the parameters  $R$  and  $N$ . More precisely, if  $\lambda(\omega, \beta)$ ,  $\lambda_N(\omega, \beta)$ ,  $\lambda_R(\omega, \beta)$  and  $\lambda_{N,R}(\omega, \beta)$  denote, respectively, the eigenvalues of  $K(\omega, \beta)$ ,  $K_N(\omega, \beta)$ ,  $K_R(\omega, \beta)$  and  $K_{N,R}^N(\omega, \beta)$ , we derive error estimates for

$$|\lambda(\omega, \beta) - \lambda_N(\omega, \beta)|, |\lambda(\omega, \beta) - \lambda_R(\omega, \beta)| \text{ and } |\lambda(\omega, \beta) - \lambda_{N,R}(\omega, \beta)|.$$

For simplicity, in what follows we shall use letter  $C$  to denote a positive generic constant, not necessarily the same at each occurrence, which is independent of the approximation parameters involved, unless we specify the contrary.

### 6.1 Analysis of the truncation related to $R$

We are going to prove that the nonzero eigenvalues of  $K_R(\omega, \beta)$  converge exponentially to the nonzero eigenvalues of  $K(\omega, \beta)$ . This would be immediate via Osborn's theory for the approximation of the spectrum of compact operators if the operator  $K_R(\omega, \beta)$  would converge to  $K(\omega, \beta)$  in the operator norm. This is not the case since  $(K(\omega, \beta) - K_R(\omega, \beta))\varphi = K(\omega, \beta)\varphi$  for any  $\varphi$  with support in  $\Gamma \setminus \Gamma_R$ .

That is why we introduce the intermediate operator

$$\tilde{K}_R(\omega, \beta) = (S_i - S_e)^{-1} \Pi_R S_p \Pi_R.$$

Notice that  $K_R = \Pi_R \tilde{K}_R$ . The idea of the proof is the following

- We prove that  $\tilde{K}_R(\omega, \beta)$  has the same nonzero eigenvalues as  $K_R(\omega, \beta)$  (Theorem 6.1).
- We prove that  $\|\tilde{K}_R(\omega, \beta) - K(\omega, \beta)\|$  converges to 0 (Theorem 6.2).
- We apply Osborn's theory to conclude (Theorem 6.3).

**Theorem 6.1** *Operators  $K_R(\omega, \beta)$  and  $\tilde{K}_R(\omega, \beta)$  have the same non zero eigenvalues. The corresponding eigenfunctions have their support included in  $\Gamma_R$ .*

*Proof.* Let  $\lambda \neq 0$  be an eigenvalue of  $\tilde{K}_R$ . Then

$$(S_i - S_e)^{-1} \Pi_R S_p \Pi_R \varphi = \lambda \varphi, \quad \varphi \neq 0. \tag{6.1}$$

Let us notice that  $\Pi_R \varphi \neq 0$ . Otherwise we would have  $\varphi = 0$ . By applying  $\Pi_R$  to (6.1) we obtain

$$K_R(\Pi_R \varphi) = \lambda(\Pi_R \varphi),$$

which proves that  $\lambda$  is an eigenvalue of  $K_R$ .

Now, let  $\lambda \neq 0$  be an eigenvalue of  $K_R$ . Then, there exists  $\varphi \in L^2(\Gamma)$ ,  $\varphi \neq 0$ , such that

$$K_R \varphi = \Pi_R (S_i - S_e)^{-1} \Pi_R S_p \Pi_R \varphi = \lambda \varphi.$$

Notice that since  $\lambda \neq 0$ , the support of  $\varphi$  is included in  $\Gamma_R$ . We set

$$\tilde{\varphi} = \varphi + \frac{1}{\lambda} (\mathbf{I} - \Pi_R) \tilde{K}_R \varphi.$$

Since the two addends are orthogonal,  $\|\tilde{\varphi}\|_{L^2} \geq \|\varphi\|_{L^2}$  and then  $\tilde{\varphi} \neq 0$ . We have

$$\begin{aligned} \tilde{K}_R \tilde{\varphi} &= \tilde{K}_R \varphi \quad (\text{since } \tilde{K}_R(\mathbf{I} - \Pi_R) = 0) \\ &= \Pi_R \tilde{K}_R \varphi + (\mathbf{I} - \Pi_R) \tilde{K}_R \varphi \\ &= K_R \varphi + \lambda \left( \frac{1}{\lambda} (\mathbf{I} - \Pi_R) \tilde{K}_R \varphi \right) \\ &= \lambda \left( \varphi + \frac{1}{\lambda} (\mathbf{I} - \Pi_R) \tilde{K}_R \varphi \right) \\ &= \lambda \tilde{\varphi}. \end{aligned}$$

Then,  $\lambda$  is a nonzero eigenvalue of  $\tilde{K}_R$ .  $\square$

The estimate of  $\|\tilde{K}_R(\omega, \beta) - K(\omega, \beta)\|$  will follow from the estimates (5.3) and (5.4) announced in Section 5.1 and which are the object of our next two lemmas.

**Lemma 6.1** *The following inequality holds*

$$\|(\mathbf{I} - \Pi_R) S_p\|_{\mathcal{L}(L^2(\Gamma))} \leq C e^{-\gamma R}$$

where  $\gamma = \gamma(\omega, \beta) = \xi_1(\omega, \beta) = \left( \frac{\pi^2}{L^2} + \beta^2 - n_o^2 \omega^2 \right)^{1/2} > 0$ .

*Proof.* Let  $\varphi \in L^2(\Gamma)$  and  $u_p$  be the solution of  $(\mathcal{P}_p)$  associated with  $\varphi$ . From the definition of the operators  $S_p$  and  $\Pi_R$  (cf. (2.8) and (5.2)), the proof reduces to get bounds, in terms of  $\|\varphi\|_{L^2(\Gamma)}$ , of the integrals

$$\begin{aligned} \int_{a^+ + R}^{+\infty} \left| \frac{\partial u_p}{\partial n}(x_1, 0) \right|^2 dx_1, & \quad \int_{a^+ + R}^{+\infty} \left| \frac{\partial u_p}{\partial n}(x_1, L) \right|^2 dx_1, \\ \int_{-\infty}^{a^- - R} \left| \frac{\partial u_p}{\partial n}(x_1, 0) \right|^2 dx_1, & \quad \int_{-\infty}^{a^- - R} \left| \frac{\partial u_p}{\partial n}(x_1, L) \right|^2 dx_1. \end{aligned}$$

Let us restrict ourselves to the first one (the other ones can be treated in a similar way). We have (cf. (3.26))

$$\begin{aligned} \int_{a^++R}^{+\infty} \left| \frac{\partial u_p}{\partial x_2}(x_1, 0) \right|^2 dx_1 &= \frac{2}{L} \int_{a^++R}^{+\infty} \left| \sum_{k=1}^{\infty} u_k^+ \frac{k\pi}{L} e^{-\xi_k(x_1-a^+)} \right|^2 dx_1 \\ &\leq \frac{2}{L} \left( \int_{a^++R}^{+\infty} e^{-2\xi_1(x_1-a^+)} dx_1 \right) \left( \sum_{k=1}^{\infty} |u_k^+| \frac{k\pi}{L} \right)^2 \\ &= \frac{2}{L} \frac{1}{2\xi_1} e^{-2\xi_1 R} \frac{\pi^2}{L^2} \left( \sum_{k=1}^{\infty} |u_k^+| k \right)^2. \end{aligned}$$

The series has been estimated in the proof of Theorem 4.5 (it is a particular case of (4.27) with  $j = 0$ ). Indeed, from formulas (4.27) and (4.28), we deduce that

$$\left( \sum_{k=1}^{\infty} |u_k^+| k \right)^2 \leq C \|u_p\|_{\mathbb{H}^2(\Sigma^+)}.$$

Using (3.19) and (3.22), we have

$$\left( \sum_{k=1}^{\infty} |u_k^+| k \right)^2 \leq C \|u_p\|_{\mathbb{H}^2(\Sigma^+)} \leq C \|u_i\|_{L^2(\Omega_i)} \leq C \|\varphi\|_{L^2(\Gamma)}, \quad (6.2)$$

so that we deduce

$$\int_{a^++R}^{+\infty} \left| \frac{\partial u_p}{\partial n}(x_1, 0) \right|^2 dx_1 \leq C e^{-2\xi_1 R} \|\varphi\|_{L^2(\Gamma)}.$$

□

**Lemma 6.2** *The following estimate holds*

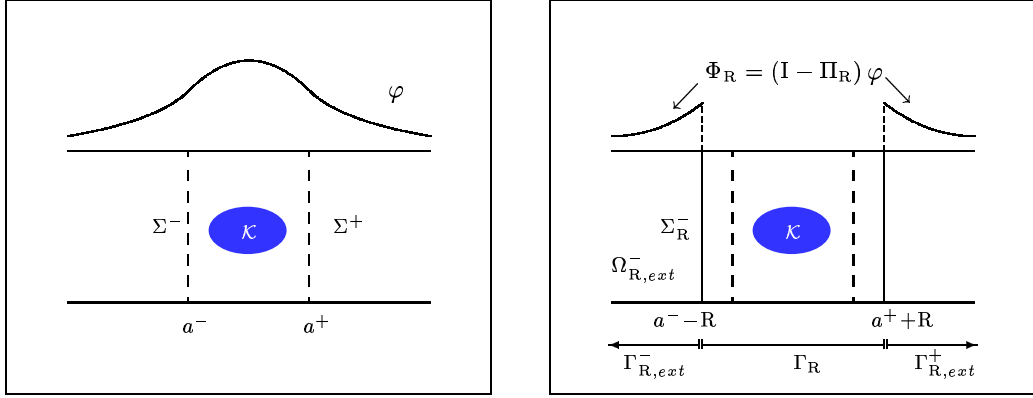
$$\|S_p(\mathbb{I} - \Pi_R)\|_{\mathcal{L}(L^2(\Gamma))} \leq C e^{-\gamma R}$$

where  $\gamma$  has been defined in Lemma 6.1.

*Proof.* For  $\varphi \in L^2(\Gamma)$  given, we set  $\Phi_R = (\mathbb{I} - \Pi_R)\varphi \in L^2(\Gamma)$ . By definition

$$S_p \Phi_R = \frac{\partial u_p}{\partial n} \Big|_{\Gamma}$$





**Figure 6.10:** Sketch of a function  $\varphi \in L^2(\Gamma)$  and the resulting  $\Phi = (I - \Pi_R)\varphi$ .

with  $u_p = u_p(\Phi_R)$ . By trace theorem we have

$$\|S_p \Phi_R\|_{L^2(\Gamma)} \leq C \left\| \frac{\partial u_p}{\partial n} \right\|_{H^{\frac{1}{2}}(\Gamma)} \leq C \|u_p\|_{H^2(\Omega_i)} \leq C \|u_i\|_{L^2(\Omega_b)}. \quad (6.3)$$

Therefore, it is enough to find a bound for  $\|u_i\|_{L^2(\Omega_b)}$ . In what follows, we shall denote

- $\Gamma_{R,ext}^- = \{(x_1, x_2), x_1 < a^- - R, x_2 = 0 \text{ or } x_2 = L\}$
- $\Gamma_{R,ext}^+ = \{(x_1, x_2), x_1 > a^+ + R, x_2 = 0 \text{ or } x_2 = L\}$
- $\Sigma_R^- = \{(a^- - R, x_2), x_2 \in (0, L)\}$
- $\Omega_{R,ext}^- = \{(x_1, x_2) \in \mathbb{R}^2, x_1 < a^- - R, 0 < x_2 < L\}$

The function  $\Phi_R$  may be considered as the sum of two functions: one of them is supported by  $\Gamma_{R,ext}^-$  and the other one is supported by  $\Gamma_{R,ext}^+$ . We shall restrict ourselves to the case  $\text{supp } \Phi_R \subset \Gamma_{R,ext}^-$ . The other one deduces in a similar way. Since we have homogeneous Dirichlet boundary conditions on  $\Gamma \setminus \Gamma_{R,ext}^-$ , if we set  $\psi = u_i|_{\Sigma_R^-}$ , we can write

$$u_i(x_1, x_2) = \sum_{k=1}^{\infty} \psi_k w_k(x_2) e^{-\xi_k(x_1 - a^- + R)} \quad \forall x_1 > a^- - R,$$

where  $\psi_k$  are the expansion coefficients of  $\psi$  in the basis  $\{w_k\}$  of  $L^2(0, L)$ . Then

$$\int_{\Omega_b} |u_i(x)|^2 dx_1 dx_2 = \sum_{k=1}^{\infty} \psi_k^2 \int_{a^-}^{a^+} e^{-2\xi_k(x_1 - a^- + R)} dx_1 \leq C \sum_{k=1}^{\infty} \psi_k^2 \frac{1}{2\xi_k} e^{-2\xi_k R}.$$

Hence, by using the asymptotic behaviour of  $\xi_k$  for large  $k$ , the definition of the  $\|\psi\|_{(\mathbb{H}_{00}^{1/2}(\Sigma_R^-))'}$  norm (see (3.24)) and the inequality  $\xi_k \leq \xi_1 = \gamma$ , we have

$$\|u_i\|_{L^2(\Omega_b)} \leq C e^{-\gamma R} \|\psi\|_{(\mathbb{H}_{00}^{1/2}(\Sigma_R^-))'}. \tag{6.4}$$

A trace theorem which can be derived from Theorem 1.5.3.4 of Grisvard [15], allows us to write

$$\|\psi\|_{(\mathbb{H}_{00}^{1/2}(\Sigma_R^-))'} \leq C \|u_i\|_{L^2(\Delta, \Omega_{R,ext}^-)}. \tag{6.5}$$

On the other hand, since  $\Phi_R \in L^2(\Gamma)$  we can immediately deduce, as a consequence of Lemma 3.2, the estimate

$$\|u_i\|_{L^2(\Delta, \Omega_i)} \leq C \|\Phi_R\|_{L^2(\Gamma)} \leq C \|\varphi\|_{L^2(\Gamma)}.$$

This inequality, together with (6.4) and (6.5), allows us to conclude.  $\square$

**Theorem 6.2** *There exist two strictly positive constants  $C$  and  $\gamma$ , depending on  $\omega$  and  $\beta$  such that*

$$\|\tilde{K}_R(\omega, \beta) - K(\omega, \beta)\|_{\mathcal{L}(L^2(\Gamma))} \leq C e^{-\gamma R}. \tag{6.6}$$

*Proof.* We simply write

$$\|K(\omega, \beta) - \tilde{K}_R(\omega, \beta)\| \leq \|(S_i - S_e)^{-1}\| \{ \|S_p(I - \Pi_R)\| + \|(I - \Pi_R)S_p\| \},$$

and use lemmas 6.2 and 6.1.  $\square$

Theorems 6.2 and 6.1 will be used now to deduce estimates for the rate of convergence of the nonzero eigenvalues of  $K(\omega, \beta)$  by those of  $K_R(\omega, \beta)$ , taking into account some classical results of spectral theory for compact operators which can be found in Dunford and Schwartz [8] or Babuška and Osborn [1].

**Theorem 6.3** *Let  $\lambda(\omega, \beta)$  be a nonzero eigenvalue of the operator  $K(\omega, \beta)$  with algebraic multiplicity  $m$  and assume its ascent is  $\eta$ . Then, there is  $R_0 > 0$  such that, for  $R > R_0$ , there exist  $\lambda_R^1(\omega, \beta), \lambda_R^2(\omega, \beta), \dots, \lambda_R^m(\omega, \beta)$  eigenvalues of  $\tilde{K}_R(\omega, \beta)$  that converge to  $\lambda(\omega, \beta)$ . Besides, we have the following estimate*

$$|\lambda(\omega, \beta) - \lambda_R^j(\omega, \beta)| \leq C e^{\frac{-\gamma R}{\eta}}, \quad \text{for } j = 1, \dots, m,$$

with  $\gamma = \gamma(\omega, \beta) > 0$ .

*Proof.* The proof follows immediately, taking into account the compactness of the operators  $K(\omega, \beta)$ ,  $\tilde{K}_R$ , Theorem 6 of Osborn [24] which states

$$|\lambda(\omega, \beta) - \lambda_R^j(\omega, \beta)| \leq C \|(K - \tilde{K}_R)|_{\mathcal{R}(E)}\|^{1/\eta},$$

and Theorem 6.2.  $\square$

**Remark 6.1** *An important property to be emphasized is the fact that there exists  $\gamma_* > 0$  independent of  $\omega$  and  $\beta$ , such that*

$$\forall(\omega, \beta) \text{ such that } \omega < \sigma_e(\beta), \quad \gamma(\omega, \beta) \geq \gamma_* > 0. \quad (6.7)$$

Indeed, for  $\omega < \sigma_e(\beta)$ , we have

$$\gamma(\omega, \beta) > \tilde{\gamma}(\beta) = (\pi^2/L^2 + \beta^2 - n_0^2 \sigma_e(\beta))^{\frac{1}{2}}.$$

From Lemma 3.1, we know that  $\tilde{\gamma}(\beta)$  is a strictly positive continuous function of  $\beta$  which satisfies  $\tilde{\gamma}(0) = \pi/L$  and furthermore (cf. Remark 1.2):

$$\tilde{\gamma}(\beta) = (\pi^2/L^2 + (1 - n_0^2/n_\infty^{+2}) \beta^2)^{\frac{1}{2}},$$

$$\tilde{\gamma}(\beta) \rightarrow \frac{\pi}{L} \quad \text{when } \beta \rightarrow +\infty,$$

which proves that  $\gamma_* = \inf \tilde{\gamma}(\beta) > 0$ .

*This result is very important in that our truncation process guarantees a uniform (in  $\beta$ ) exponential rate of convergence in  $R$  even though the eigenmode  $u(\beta)$  one looks for has not a uniform exponential decay in the variable  $x_1$ . Indeed, if  $\omega(\beta)$  is the pulsation of the eigenmode, one easily shows that this decay is given by*

$$C \exp(-[(\sigma_e(\beta) - \omega(\beta)^2)^{\frac{1}{2}} |x_1|])$$

where  $(\sigma_e(\beta) - \omega^2) > 0$  goes to 0 when  $\beta$  approaches a threshold (cf. (1.20)). Therefore, the accuracy of our method does not deteriorate at the vicinity of thresholds, which would be the case with a brutal truncation method.

## 6.2 Analysis of the truncation related to $N$

The aim of this section is to analyze the numerical error due to the truncation parameter  $N$  (cf. Section 5.2). The effect of the series truncation is that we approximate the eigenvalues  $\lambda(\omega, \beta)$  of the operator  $K(\omega, \beta)$  by the eigenvalues  $\lambda_N(\omega, \beta)$  of the operator  $K_N(\omega, \beta)$ . In order to get an error estimate on  $\lambda(\omega, \beta) - \lambda_N(\omega, \beta)$ , we need an estimate of

$$\|K(\omega, \beta) - K_N(\omega, \beta)\|_{\mathcal{L}(L^2(\Gamma))},$$

that is, an estimate of

$$\|S_p - S_p^N\|_{\mathcal{L}(L^2(\Gamma))}.$$

For this, we shall establish an error estimate of the type

$$\|u_p - u_p^N\|_{H^2(\Omega_b)} \leq \epsilon(N) \|\varphi\|_{L^2(\Gamma)},$$

where  $\epsilon(N)$  denotes some positive function which tends to 0 when  $N$  tends to  $+\infty$ . The estimate of  $u_p - u_p^N$  outside  $\Omega_b$  is then quasi-explicit using formulas (3.32) and (5.11). The main difficulty is thus reduced to the approximation of  $(\mathcal{P}_p^\Sigma)$  by  $(\mathcal{P}_p^\Sigma)_N$ . The analysis appears not so obvious and will be splitted in several steps:

1. *Existence and uniqueness of  $u_p^N$  solution of  $(\mathcal{P}_p^\Sigma)_N$ .* The idea is to prove an existence and stability result (this means that we want to get estimates on  $u_p^N$  which are independent of  $N$ ) by using the fact that the problem  $(\mathcal{P}_p^\Sigma)_N$  is close enough, at least for  $N$  large enough, to the problem  $(\mathcal{P}_p^\Sigma)$ , which is well posed since  $\omega \notin G_i(\beta)$ .

The main tools will be a perturbation method and a fixed point technique. A technical difficulty arises due to the fact that the difference  $T_N(\omega, \beta) - T(\omega, \beta)$  does not converge to 0 in the operator norm. To overcome this difficulty, and following ideas developed by Razafiarivelo in [26] for instance, we write two problems  $(\mathcal{P}_p^{\Sigma^\circ})$  and  $(\mathcal{P}_p^{\Sigma^\circ})_N^*$  which are “equivalent” to the original problems  $(\mathcal{P}_p^\Sigma)$  and  $(\mathcal{P}_p^\Sigma)_N$  (in a sense to be precised later) but posed in a subdomain  $\Omega_o$  of  $\Omega_b$ , namely the rectangle

$$\Omega_o = \{(x_1, x_2) \mid a^- + \ell < x_1 < a^+ + \ell, \quad 0 < x_2 < L\}.$$

The advantage with these new problems is that the operators  $T(\omega, \beta)$  and  $T_N(\omega, \beta)$  appearing in  $(\mathcal{P}_p^\Sigma)$  and  $(\mathcal{P}_p^\Sigma)_N$  are replaced, respectively, by  $T^*(\omega, \beta)$  and  $T_N^*(\omega, \beta)$  where the difference  $T^*(\omega, \beta) - T_N^*(\omega, \beta)$  converges exponentially to 0, which permits to make the fixed point procedure work. Then, we are able to prove the existence and uniqueness of  $u_p^N$  only for  $N$  large enough, which is rather classical in the approximation of boundary value problems which are not coercive but only compact perturbations of coercive problems.

2. *Derivation of the error estimates.* Following the existence proof, error estimates for  $u_p - u_p^N$  are then obtained in three steps: first in the smaller domain  $\Omega_o$  using a perturbation technique, then successively in the domains  $\Omega_b \setminus \Omega_o$  and  $\Omega_i \setminus \Omega_b$  using two different but not difficult techniques.

### Analysis of $(\mathcal{P}_p^\Sigma)_N$ and construction of $K_N$

We shall begin by introducing two auxiliary problems posed in two vertical bands  $C^+$  and  $C^-$ , contained in the domain  $\Omega_b$ , where the refraction index  $n$  is constant and equal to  $n_o$ . This allows making analytical computations in these domains.

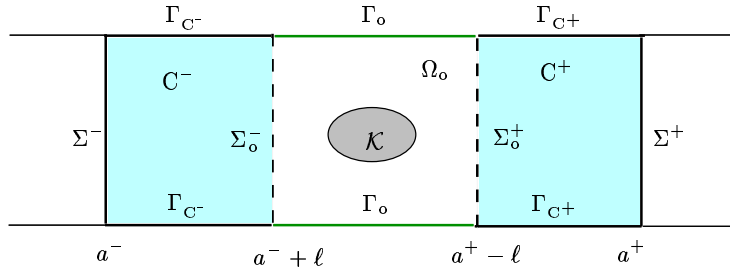


Figure 6.11: Sketch of the subdomains  $C^+$  and  $C^-$ .

Let  $\ell \in \mathbb{R}$ ,  $0 < \ell < (a^+ - a^-)/2$ . We shall denote (cf. Figure 6.11):

- $\Omega_o = \{(x_1, x_2) \in \mathbb{R}^2 / a^- + \ell < x_1 < a^+ - \ell, 0 < x_2 < L\}$ .
- $\Sigma_o^+ = \{(x_1, x_2) \in \mathbb{R}^2 / x_1 = a^+ - \ell, 0 < x_2 < L\}$ .
- $\Sigma_o^- = \{(x_1, x_2) \in \mathbb{R}^2 / x_1 = a^- + \ell, 0 < x_2 < L\}$ .
- $C^+ = \{(x_1, x_2) \in \mathbb{R}^2 / a^+ - \ell < x_1 < a^+, 0 < x_2 < L\}$ .
- $C^- = \{(x_1, x_2) \in \mathbb{R}^2 / a^- < x_1 < a^- + \ell, 0 < x_2 < L\}$ .
- $\Gamma_{C^+} = \Gamma \cap C^+$ ,
- $\Gamma_{C^-} = \Gamma \cap C^-$ .
- $\Gamma_o = \Gamma \cap \partial\Omega_o$ .

Let  $g^+$  and  $g^-$  be any functions belonging to  $H_{00}^{\frac{1}{2}}(\Sigma_o^+)$  and  $H_{00}^{\frac{1}{2}}(\Sigma_o^-)$ , respectively, and let  $T_N$  be the truncated operator introduced in (5.9). We consider the problems

$$(\mathcal{P}^{C^+}) \left\{ \begin{array}{l} -\Delta u + (\beta^2 - n_0^2 \omega^2) u = 0 \text{ in } C^+, \quad u \in H^1(C^+), \\ u = 0 \text{ on } \Gamma_{C^+}, \\ u = g^+ \text{ on } \Sigma_0^+, \\ \frac{\partial u}{\partial \nu} + T_N u = 0 \text{ on } \Sigma^+, \end{array} \right.$$

and

$$(\mathcal{P}^{C^-}) \left\{ \begin{array}{l} -\Delta u + (\beta^2 - n_0^2 \omega^2) u = 0 \text{ in } C^-, \quad u \in H^1(C^-), \\ u = 0 \text{ on } \Gamma_{C^-}, \\ u = g^- \text{ on } \Sigma_0^-, \\ \frac{\partial u}{\partial \nu} + T_N u = 0 \text{ on } \Sigma^-. \end{array} \right.$$

The following results will be useful for our purposes.

**Lemma 6.3** *For every  $g^+ \in H_{00}^{\frac{1}{2}}(\Sigma_0^+)$ , the problem  $(\mathcal{P}^{C^+})$  has a unique solution  $v^N \in H^1(C^+)$  given by:*

$$v^N(x_1, x_2) = \sum_{k=1}^N a_k e^{-\xi_k x_1} w_k(x_2) + \sum_{k=N+1}^{\infty} a_k \left( e^{-\xi_k x_1} + e^{\xi_k(x_1 - 2a^+)} \right) w_k(x_2) \quad (6.8)$$

where

$$a_k = \begin{cases} g_k^+ e^{\xi_k(a^+ - \ell)} & \text{if } k \leq N, \\ \frac{g_k^+ e^{\xi_k a^+}}{2 \cosh(\xi_k \ell)} & \text{if } k \geq N + 1, \end{cases} \quad (6.9)$$

and  $g_k^+$  denotes the expansion coefficients of  $g^+$  in the basis  $\{w_k\}$  of  $H_{00}^{\frac{1}{2}}(\Sigma_0^+)$ . Besides, there exists a constant  $C$  independent of  $N$  such that

$$\|v^N\|_{H(\Delta, C^+)} \leq C \|g\|_{H_{00}^{\frac{1}{2}}(\Sigma^+)}. \quad (6.10)$$

*Proof.* The proof is straightforward. Formula (6.8) follows immediately from a classical technique of separation of variables.  $\square$

By the symmetry of the domain, we have a similar result to Lemma 6.3 but in the subdomain  $C^-$ .

We now define transparent boundary conditions on  $\Sigma_o^+$  and  $\Sigma_o^-$  to formulate a problem posed in the domain  $\Omega_o$  whose solution is the restriction of the solution to problem  $(\mathcal{P}_p^\Sigma)_N$  to this domain. More precisely, as we did for the operators  $T_+$  and  $T_-$ , we define two new Dirichlet-Neumann operators as follows:

$$(T_N^+)^*(\omega, \beta) : \begin{array}{l} \mathbb{H}_{00}^{\frac{1}{2}}(\Sigma_o^+) \longrightarrow \mathbb{H}_{00}^{\frac{1}{2}}(\Sigma_o^+)' \\ g^+ \longmapsto \left. \frac{\partial v^N}{\partial \nu} \right|_{\Sigma_o^+} \end{array} \quad (6.11)$$

$$(T_N^-)^*(\omega, \beta) : \begin{array}{l} \mathbb{H}_{00}^{\frac{1}{2}}(\Sigma_o^-) \longrightarrow \mathbb{H}_{00}^{\frac{1}{2}}(\Sigma_o^-)' \\ g^- \longmapsto \left. \frac{\partial v^N}{\partial \nu} \right|_{\Sigma_o^-} \end{array} \quad (6.12)$$

$n$  being the solution of  $(\mathcal{P}^{C^+})$  or  $(\mathcal{P}^{C^-})$  associated with the Dirichlet data  $g^+$  or  $g^-$  (here,  $\nu$  denotes the outgoing normal vector to the boundaries  $\Sigma_o^+$  or  $\Sigma_o^-$  from  $C^+$  and  $C^-$  respectively.)

As we have done for the operator  $T$ , making the identifications  $\Sigma_o^+ \equiv \Sigma_o^- \equiv \Sigma_o = (0, L)$  we can identify  $(T_N^+)^*$  and  $(T_N^-)^*$  to a single operator  $T_N^*$  belonging to  $\mathcal{L}(\mathbb{H}_{00}^{\frac{1}{2}}(\Sigma_o), \mathbb{H}_{00}^{\frac{1}{2}}(\Sigma_o)')$ . Using Lemma 6.3 it is easy to prove the following theorem.

**Theorem 6.4** *The operator  $T_N^*$  has the following representation:*

$$[T_N^*(\omega, \beta) \varphi] = \sum_{k=1}^{\infty} t_k^* \mathbf{w}_k(x_2) \varphi_k \quad \forall \varphi \in \mathbb{H}_{00}^{\frac{1}{2}}(\Sigma_o), \quad (6.13)$$

where  $\varphi_k$  denotes the expansion coefficients of  $\varphi$  in the basis  $\{\mathbf{w}_k\}$  of  $\mathbb{H}_{00}^{\frac{1}{2}}(\Sigma_o)$  and

$$t_k^* = \begin{cases} \xi_k & \text{if } k \leq N, \\ \xi_k \tanh(\xi_k \ell) & \text{if } k > N. \end{cases} \quad (6.14)$$

This leads us naturally to introduce the problem:

$$(\mathcal{P}_p^{\Sigma_o})_N^* \left\{ \begin{array}{l} -\Delta u_{p,0}^N + (\beta^2 - n^2 \omega^2) u_{p,0}^N = (n^2 - n_o^2) \omega^2 u_i \quad \text{in } \Omega_o, \quad u_{p,0}^N \in H^1(\Omega_o), \\ u_{p,0}^N = 0 \quad \text{on } \Gamma_o, \\ \frac{\partial u_{p,0}^N}{\partial \nu} = -T_N^* u_{p,0}^N \quad \text{on } \Sigma_o^\pm. \end{array} \right. \quad (6.15)$$

Problems  $(\mathcal{P}_p^\Sigma)_N$  and  $(\mathcal{P}_p^{\Sigma_o})_N^*$  are equivalent in the sense of the following theorem (whose trivial proof is omitted)

**Theorem 6.5** *If  $u_p^N$  is a solution of  $(\mathcal{P}_p^\Sigma)_N$ , then  $u_p^N|_{\Omega_o}$  is a solution of  $(\mathcal{P}_p^{\Sigma_o})_N^*$ . Conversely, if  $u_{p,0}^N$  denotes a solution of  $(\mathcal{P}_p^{\Sigma_o})_N^*$ , then  $u_{p,0}^N$  can be extended to a solution of  $(\mathcal{P}_p^\Sigma)_N$  in a unique way by solving  $(\mathcal{P}^{C^+})$  (respectively  $(\mathcal{P}^{C^-})$ ) with  $g^+ = u_p^N|_{\Sigma_o^+}$  (respectively  $g^- = u_p^N|_{\Sigma_o^-}$ ).*

The main result of this section is Theorem 6.7 which ensures the existence of solution of the problem  $(\mathcal{P}_p^{\Sigma_o})_N^*$ , and consequently, taking into account the equivalence of both problems, the existence of solution of the initial semidiscretized problem  $(\mathcal{P}_p^\Sigma)_N$ . The proof of Theorem 6.7 will use Lemma 6.4 and Theorem 6.6 below. In the sequel, we shall denote  $\alpha = \pi\ell/L$ .

**Lemma 6.4** *Let  $T^* (\equiv T) \in \mathcal{L}(H_{00}^{\frac{1}{2}}(\Sigma_o), H_{00}^{\frac{1}{2}}(\Sigma_o)')$  be the operator given by (3.30), once  $\Sigma_o$  has been identified to  $[0, L]$ . There exists a constant  $C$  independent of  $N$  such that*

$$\|T^* - T_N^*\|_{\mathcal{L}(H_{00}^{\frac{1}{2}}(\Sigma_o), H_{00}^{\frac{1}{2}}(\Sigma_o)')} \leq C e^{-2\alpha N}. \quad (6.16)$$

*Proof.* Let  $\varphi \in H_{00}^{\frac{1}{2}}(\Sigma_o)$ . We have

$$\|(T^* - T_N^*)\varphi\|_{H_{00}^{\frac{1}{2}}(\Sigma_o)'}^2 = \sum_{k \geq N+1} \xi_k^2 (1 - \tanh(\xi_k \ell))^2 \frac{|\varphi_k|^2}{(1 + k^2)^{1/2}}.$$

Since  $\xi_k^2 \leq C(k^2 + 1)$  and  $1 - \tanh(z)$  is a decreasing function of  $z$ , it results that

$$\begin{aligned} \|(T^* - T_N^*)\varphi\|_{H_{00}^{\frac{1}{2}}(\Sigma_o)'}^2 &\leq C [1 - \tanh(\xi_{N+1} \ell)]^2 \sum_{k \geq N+1} (1 + k^2)^{1/2} |\varphi_k|^2 \\ &\leq C [1 - \tanh(\xi_{N+1} \ell)]^2 \|\varphi\|_{H_{00}^{\frac{1}{2}}(\Sigma_o)}^2. \end{aligned}$$



Thus, taking into account that  $\xi_{N+1} \ell \geq N\alpha$  and  $1 - \tanh(z) \leq e^{-z}$ , we obtain

$$\|T^* - T_N^*\|_{\mathcal{L}(\mathbb{H}_{00}^{\frac{1}{2}}(\Sigma_o), \mathbb{H}_{00}^{\frac{1}{2}}(\Sigma_o)')} \leq C e^{-2\alpha N}$$

with  $C$  a constant independent of  $N$ .  $\square$

The following theorem extends the inequalities (3.33) and (3.34) to any function  $f \in L^2(\Omega_i)$  with compact support in  $\Omega_b$ .

**Theorem 6.6** *Let  $f \in L^2(\Omega_i)$ , with  $\text{supp}(f) \subset \Omega_b$ ,  $g \in \mathbb{H}_{00}^{\frac{1}{2}}(\Sigma)'$ . Let  $T$  be the operator defined in (3.30). The problem*

$$\left\{ \begin{array}{l} \text{Find } u \in H^1(\Omega_b) \text{ such that} \\ -\Delta u + (\beta^2 - n^2 \omega^2)u = f \quad \text{in } \Omega_b \\ u = 0 \quad \text{on } \Gamma_b \\ \frac{\partial u}{\partial \nu} + Tu = g \quad \text{on } \Sigma^\pm \end{array} \right. \quad (6.17)$$

has a unique solution  $u \in H^1(\Delta, \Omega_b)$  which satisfies

$$\|u\|_{H^1(\Omega_b)} + \|\Delta u\|_{L^2(\Omega_b)} \leq C(\|f\|_{L^2(\Omega_b)} + \|g\|_{\mathbb{H}_{00}^{\frac{1}{2}}(\Sigma)'}). \quad (6.18)$$

The following theorem provides an existence and stability result for problem  $(\mathcal{P}_p^{\Sigma_o})_N^*$ .

**Theorem 6.7** *Problem  $(\mathcal{P}_p^{\Sigma_o})_N^*$  has a unique solution  $u_{p,o}^N$  for  $N$  large enough. This solution satisfies*

$$\|u_{p,o}^N\|_{H^1(\Omega_b)} \leq C \|f\|_{L^2(\Omega_i)}. \quad (6.19)$$

*Proof.* The boundary condition on  $\Sigma_0$  of  $(\mathcal{P}_p^{\Sigma_o})_N^*$  can be written as

$$\frac{\partial u_{p,o}^N}{\partial \nu} + T^* u_{p,o}^N = (T^* - T_N^*) u_{p,o}^N. \quad (6.20)$$

Taking this fact into account, the existence of solution of the problem  $(\mathcal{P}_p^{\Sigma_o})_N^*$  reduces to prove the contractivity of the following mapping:

$$\begin{array}{ccc} \mathcal{R}_N : H^1(\Omega_o) & \longrightarrow & H^1(\Delta, \Omega_o) \subset H^1(\Omega_o) \\ z & \longmapsto & \mathcal{R}_N z = v_N \end{array} \quad (6.21)$$

with  $v_N$  the unique solution of the problem:

$$\left\{ \begin{array}{l} -\Delta v_N + (\beta^2 - n^2 \omega^2) v_N = (n^2 - n_0^2) \omega^2 u_i \quad \text{in } \Omega_0, \\ v_N = 0 \quad \text{on } \Gamma_0, \\ \frac{\partial v_N}{\partial \nu} + T^* v_N = (T^* - T_N^*) z \quad \text{on } \Sigma_0^\pm. \end{array} \right. \quad (6.22)$$

Notice that  $\mathcal{R}_N$  is well defined because of Theorem 6.6.

Now we prove that  $\mathcal{R}_N$  is a contractive mapping. Let  $\zeta_N = \mathcal{R}_N z_1 - \mathcal{R}_N z_2$ . Then  $\zeta_N$  verifies:

$$\left\{ \begin{array}{l} -\Delta \zeta_N + (\beta^2 - n^2 \omega^2) \zeta_N = 0 \quad \text{in } \Omega_0, \\ \zeta_N = 0 \quad \text{on } \Gamma_0, \\ \frac{\partial \zeta_N}{\partial \nu} + T^* \zeta_N = (T^* - T_N^*) (z_1 - z_2) \quad \text{on } \Sigma_0^\pm, \end{array} \right. \quad (6.23)$$

and taking into account (6.16) and (6.18) we get

$$\begin{aligned} \|\zeta_N\|_{\mathbf{H}^1(\Delta, \Omega_0)} &= \|\mathcal{R}_N z_1 - \mathcal{R}_N z_2\|_{\mathbf{H}^1(\Delta, \Omega_0)} \\ &\leq C(\Omega_0) \|(T^* - T_N^*) (z_1 - z_2)\|_{\mathbf{H}_{00}^{\frac{1}{2}}(\Sigma_0)'} \\ &\leq C(\Omega_0) \|T^* - T_N^*\|_{\mathcal{L}(\mathbf{H}_{00}^{\frac{1}{2}}(\Sigma_0), \mathbf{H}_{00}^{\frac{1}{2}}(\Sigma_0)')} \|z_1 - z_2\|_{\mathbf{H}_{00}^{\frac{1}{2}}(\Sigma_0)} \\ &\leq C(\Omega_0) e^{-2\alpha N} \|z_1 - z_2\|_{\mathbf{H}^1(\Omega_0)}, \end{aligned}$$

from which it follows that  $\mathcal{R}_N$  is a contractive mapping for  $N$  large enough.

From boundary condition (6.20) and Theorem 6.6 it results that

$$\|u_{p,0}^N\|_{\mathbf{H}^1(\Omega_0)} \leq C(\|f\|_{\mathbf{L}^2(\Omega_i)} + \|(T^* - T_N^*) u_{p,0}^N\|_{\mathbf{H}_{00}^{\frac{1}{2}}(\Sigma_0^\pm)'}).$$

From Lemma 6.4 and trace theorem we deduce

$$\|u_{p,0}^N\|_{\mathbf{H}^1(\Omega_0)} \leq C\|f\|_{\mathbf{L}^2(\Omega_i)} + C e^{-2\alpha N} \|u_{p,0}^N\|_{\mathbf{H}^1(\Omega_0)}.$$

Since  $(1 - C e^{-2\alpha N}) \rightarrow 1$ , for  $N$  large enough we have  $(1 - C e^{-2\alpha N}) > 0$  and then

$$\|u_{p,0}^N\|_{\mathbf{H}^1(\Omega_0)} \leq \frac{C}{1 - C e^{-2\alpha N}} \|f\|_{\mathbf{L}^2(\Omega_i)},$$

which finishes the proof.  $\square$

Theorems 6.5 and 6.7 allow us to state the following

**Theorem 6.8** *The semidiscrete problem  $(\mathcal{P}_p^\Sigma)_N$  has a unique solution  $u_p^N$  for  $N$  large enough. This solution satisfies*

$$\|u_p^N\|_{H^1(\Omega_b)} \leq C \|f\|_{L^2(\Omega_i)}. \quad (6.24)$$

*Proof.* To get the estimate in  $\Omega_b$ , it suffices to use the estimate (6.10) with  $g^+ = g^+ = u_p^N|_{\Sigma_0^+}$  and to apply trace theorem to obtain an estimate in  $H^1(C^+)$ . One proceeds in the same manner to get the estimate in  $H^1(C^-)$ .  $\square$

In the sequel, we shall also need an  $H^2$  estimate for  $u_p^N$  in  $\Omega_b$ .

**Lemma 6.5** *Let  $u_p^N$  be the solution of the problem  $(\mathcal{P}_p^\Sigma)_N$ . Then  $u_p^N \in H^2(\Omega_b)$  and*

$$\|u_p^N\|_{H^2(\Omega_b)} \leq C N \|f\|_{L^2(\Omega_i)}. \quad (6.25)$$

*Proof.* Let us denote  $\delta = \beta^2 - n^2\omega^2 \in L^\infty$ . From (5.10), one has

$$\left\{ \begin{array}{l} -\Delta u_p^N + u_p^N = (1 - \delta) u_p^N + (n^2 - n_0^2) \omega^2 u_i \quad \text{in } \Omega_b, \\ \frac{\partial u_p^N}{\partial \nu} = -T_N u_p^N \quad \text{on } \Sigma^\pm, \\ u_p^N = 0 \quad \text{on } \Gamma_b. \end{array} \right.$$

It follows immediately that  $T_N u_p^N \in H_{00}^{\frac{1}{2}}(\Sigma)$  because it corresponds to a finite sum of  $C^\infty$  functions which vanish at the ends of  $\Sigma$ . Therefore, classical regularity results which can be found in [15], allow us to ensure that  $u_p^N \in H^2(\Omega_b)$  and

$$\|u_p^N\|_{H^2(\Omega_b)} \leq C (\|f + (1 - \delta) u_p^N\|_{L^2(\Omega_b)} + \|T_N u_p^N\|_{H_{00}^{\frac{1}{2}}(\Sigma)}), \quad (6.26)$$

and hence, using (6.24)

$$\|u_p^N\|_{H^2(\Omega_b)} \leq C (\|f\|_{L^2(\Omega_b)} + \|T_N u_p^N\|_{H_{00}^{\frac{1}{2}}(\Sigma)}). \quad (6.27)$$

On the other hand, by definition of  $T_N$ , we have

$$\|T_N u_p^N\|_{H_{00}^{\frac{1}{2}}(\Sigma)}^2 = \sum_{k \leq N} \xi_k^2 (u_p^N)_k^2 (1 + k^2)^{1/2}.$$

Then, since  $\xi_k^2 \sim C(1 + k^2)$  for large  $k$ , we get

$$\begin{aligned} \|T_N u_p^N\|_{\mathbb{H}_{00}^{\frac{1}{2}}(\Sigma)}^2 &\leq C \sum_{k \leq N} \xi_k^3 (u_p^N)_k^2 \leq C \xi_N^2 \|u_p^N\|_{\mathbb{H}_{00}^{\frac{1}{2}}(\Sigma)}^2 \\ &\leq C N^2 \|u_p^N\|_{\mathbb{H}^1(\Omega_b)}^2 \leq C N^2 \|f\|_{\mathbb{L}^2(\Omega_b)}^2. \end{aligned} \quad (6.28)$$

Then estimate (6.25) results from (6.27) and (6.28).  $\square$

It is now easy to get an  $\mathbb{H}^2$  estimate in  $\Omega_{ext}$  for the function  $u_p^N$  defined by (5.11).

**Lemma 6.6** *The function  $u_p^N$  defined by (5.11) satisfies:*

$$\|u_p^N\|_{\mathbb{H}^2(\Omega_{ext})} \leq C N \|f\|_{\mathbb{L}^2(\Omega_i)}. \quad (6.29)$$

Notice that

$$u_p^N \equiv (u_p^N|_{\Omega_b}, u_p^N|_{\Omega_{ext}}) \in \mathbb{H}^2(\Omega_b) \times \mathbb{H}^2(\Omega_{ext})$$

implies

$$\left( \frac{\partial u_p^N}{\partial n} \Big|_{\Gamma_b}, \frac{\partial u_p^N}{\partial n} \Big|_{\Gamma_{ext}} \right) \in \mathbb{H}^{\frac{1}{2}}(\Gamma_b) \times \mathbb{H}^{\frac{1}{2}}(\Gamma_{ext}),$$

so that the definition (5.12) permits to define  $S_p^N$  as a bounded operator in  $\mathbb{L}^2(\Gamma)$ . Proceeding as for the proof of Theorem 4.4, it is not difficult to prove that  $S_p^N$  is compact in  $\mathbb{L}^2(\Gamma)$  and hence the operator

$$K_N = (S_i - S_e)^{-1} S_p^N \in \mathcal{L}(\mathbb{L}^2(\Gamma))$$

is compact too.

### Error estimates

We derive error estimates for  $u_p - u_p^N$ , first in the domain  $\Omega_o$  (Lemma 6.7), then by extension in the domains  $\mathbb{C}^+$  and  $\mathbb{C}^-$  (Lemma 6.8) and finally in the exterior domain  $\Omega_{ext}$  (Lemma 6.9).

**Lemma 6.7** *Let  $u_p$  be the solution of the problem  $(\mathcal{P}_p^\Sigma)$  and  $u_p^N$  be the solution of  $(\mathcal{P}_p^\Sigma)_N$ . There exists a constant  $C$  independent of  $N$  such that*

$$\|u_p - u_p^N\|_{\mathbb{H}^1(\Delta, \Omega_o)} \leq C e^{-2\alpha N} \|f\|_{\mathbb{L}^2(\Omega_i)}. \quad (6.30)$$

*Proof.* It is immediate to check that  $u_p|_{\Omega_o}$  satisfies

$$(\mathcal{P}_p^{\Sigma_o}) \begin{cases} -\Delta u + (\beta^2 - n^2\omega^2)u = (n^2 - n_o^2)\omega^2 u_i & \text{in } \Omega_o, \\ u = 0 & \text{on } \Gamma_o, \\ \frac{\partial u}{\partial \nu} + T^* u = 0 & \text{on } \Sigma_o^\pm, \end{cases}$$

where we recall that  $T^*$  is the analogous operator to the operator  $T$  but defined on the boundaries  $\Sigma_o^+$  and  $\Sigma_o^-$ .

On the other hand,  $u_p^N|_{\Omega_o}$  satisfies  $(\mathcal{P}_p^{\Sigma_o})_N^*$ . Thus, the difference  $\tilde{u}^N = u_p - u_p^N$  satisfies:

$$\begin{cases} -\Delta \tilde{u}^N + (\beta^2 - n^2\omega^2)\tilde{u}^N = 0 & \text{in } \Omega_o, \\ \tilde{u}^N = 0 & \text{on } \Gamma_o, \\ \frac{\partial \tilde{u}^N}{\partial \nu} + T^* \tilde{u}^N = -(T^* - T_N^*)u_p^N & \text{on } \Sigma_o^\pm. \end{cases} \quad (6.31)$$

By Theorem 6.6 and Lemma 6.4

$$\begin{aligned} \|u_p - u_p^N\|_{H^1(\Delta, \Omega_o)} &\leq C \|(T^* - T_N^*)u_p^N\|_{H_{00}^{\frac{1}{2}}(\Sigma_o)^+} \\ &\leq C e^{-2\alpha N} \|u_p^N\|_{H^1(\Omega_o)} \end{aligned}$$

and taking into account (6.19) we get (6.30)  $\square$

**Lemma 6.8** *Let  $u_p$  (respectively  $u_p^N$ ) be the solution of  $(\mathcal{P}_p^\Sigma)$  (respectively of  $(\mathcal{P}_p^\Sigma)_N$ ). There exists a constant  $C$  independent of  $N$  such that:*

$$\|u_p - u_p^N\|_{H^1(\Delta, C^+)} \leq C e^{-\alpha N} \|f\|_{L^2(\Omega_i)}. \quad (6.32)$$

*Besides*

$$\left\| \frac{\partial u_p}{\partial \nu} - \frac{\partial u_p^N}{\partial \nu} \right\|_{H_{00}^{\frac{1}{2}}(\Sigma^+)^+} \leq C \|u_p - u_p^N\|_{H(\Delta, C^+)}. \quad (6.33)$$

*Proof.* We explain how to get the estimate in  $C^+$ . We denote by  $g$  and  $g^N$  the traces of  $u_p$  and  $u_p^N$  on  $\Sigma_o$ . The difference  $\tilde{u}^N = u_p - u_p^N$  satisfies:

$$\begin{cases} -\Delta \tilde{u}^N + (\beta^2 - n^2\omega^2)\tilde{u}^N = 0 & \text{in } C^+, \\ \tilde{u}^N = 0 & \text{on } \Gamma_{C^+}, \\ \tilde{u}^N = g - g^N & \text{on } \Sigma_o^+, \\ \frac{\partial \tilde{u}^N}{\partial \nu} + T^* \tilde{u}^N = -(T^* - T_N^*)u_p^N & \text{on } \Sigma^+. \end{cases} \quad (6.34)$$

We decompose  $\tilde{u}^N$  as

$$\tilde{u}^N = z^N + v^N$$

where the function  $z^N$  satisfies

$$\begin{cases} z^N \in H^1(C^+), & \|z^N\|_{H^1(C^+)} \leq C \|g - g^N\|_{H_{00}^{\frac{1}{2}}(\Sigma)}, \\ z^N(x_1, x_2) = 0 & \text{if } a^+ - \eta < x_1 < a^+, \quad \text{with } \eta > 0. \end{cases}$$

The existence of such a  $z^N$  is guaranteed by the trace theorem. Then, it is not difficult to check that  $v^N$  satisfies (it suffices to write the equation satisfied by  $v^N$ , to multiply by  $v^N$  and to integrate over  $C^+$ )

$$\begin{aligned} \int_{C^+} \{|\nabla v^N|^2 + (\beta^2 - n_0^2 \omega^2)|v^N|^2\} + \langle T_N v^N, v^N \rangle_{\Sigma_0^+} &= \\ = \int_{C^+} \{\nabla z^N \cdot \nabla v^N + (\beta^2 - n_0^2 \omega^2)z^N v^N\} - \langle (T_N - T)u_p, v^N \rangle_{\Sigma_0^+}. \end{aligned} \quad (6.35)$$

Since the operator is positive and

$$\int_{C^+} \{|\nabla v^N|^2 dx + (\beta^2 - n_0^2 \omega^2)|v^N|^2\} \geq C \|v^N\|_{H^1(C^+)}^2,$$

we easily get

$$\|v^N\|_{H^1(C^+)}^2 \leq C (\|g - g^N\| + \|(T_N - T)u_p\|). \quad (6.36)$$

Therefore, by the triangular inequality

$$\|\tilde{u}^N\|_{H^1(C^+)}^2 \leq C (\|g - g^N\|_{H_{00}^{\frac{1}{2}}(\Sigma_0^+)} + \|(T_N - T)u_p\|_{H_{00}^{\frac{1}{2}}(\Sigma^+)'}), \quad (6.37)$$

and from Lemma 6.7 and trace theorem, we know that

$$\|g - g^N\|_{H_{00}^{\frac{1}{2}}(\Sigma_0^+)} \leq C e^{-2\alpha N} \|f\|_{L^2(\Omega_i)}.$$

Moreover, if  $u_k$  are the coefficients of the expansion for the trace of  $u_p$  on  $\Sigma^+$  in the basis  $w_k$ , we have

$$\|(T_N - T)u_p\|_{H_{00}^{\frac{1}{2}}(\Sigma^+)'}^2 = \sum_{k>N} (1 + k^2)^{-1/2} |\xi_k|^2 |u_k|^2 \leq C \sum_{k>N} (1 + k^2)^{1/2} |u_k|^2.$$

In the region  $x_1 > a^+ - l$  in which  $n = n_o$ , we know that  $u_p$  is of the form

$$u_p(x_1, x_2) = \sum_{k=1}^{\infty} A_k^+ w_k(x_2) e^{-\xi_k(x_1 - a^+)},$$

therefore, if  $u_k^o$  are the coefficients of the expansion trace of  $u_p$  on  $\Sigma_o^+$  in the basis  $w_k$ , we know that  $u_k = u_k^o e^{-2\xi_k l}$ . Consequently, since  $\xi_{N+1}l \geq \alpha N$ ,

$$\begin{aligned} \|(T_N - T)u_p\|_{H_{00}^{\frac{1}{2}}(\Sigma^+)}^2 &\leq C \sum_{k>N} (1+k^2)^{1/2} e^{-2\xi_k l} |u_k^o|^2 \\ &\leq e^{-2\xi_{N+1}l} \|u_p\|_{H_{00}^{\frac{1}{2}}(\Sigma_o^+)}^2 \\ &\leq C e^{-2\alpha N} \|u_p\|_{H^1(\Omega_o)}^2. \end{aligned} \tag{6.38}$$

□

In the same way, it is possible to extend estimates (6.32) and (6.33) to the domain  $\Omega_{ext}$ . We do not give the details of the proof which is very similar to the previous one (see [13]).

**Lemma 6.9** *There exists a constant  $C$  independent of  $N$  such that*

$$\|u_p - u_p^N\|_{H^1(\Omega_{ext})} \leq C e^{-\alpha N} \|f\|_{L^2(\Omega_i)}. \tag{6.39}$$

In order to get a bound for  $\|S_p - S_p^N\|_{\mathcal{L}(L^2(\Gamma))}$ , we shall need  $H^2$  estimates instead of  $H^1$  ones. To obtain them, we shall use the following intermediate result

**Lemma 6.10** *The following properties hold:*

$$(i) \quad \|T_N(u_p - u_p^N)\|_{H_{00}^{\frac{1}{2}}(\Sigma)} \leq C N \|u_p - u_p^N\|_{H^1(\Omega_b)}. \tag{6.40}$$

$$(ii) \quad \|(T - T_N)u_p\|_{H_{00}^{\frac{1}{2}}(\Sigma)} \leq C e^{-\alpha N} \|f\|_{L^2(\Omega_i)}. \tag{6.41}$$

*Proof. (i)* If  $(u_p - u_p^N)_k$  denotes, as usual, the  $k^{\text{th}}$  Fourier coefficient of  $u_p|_\Sigma - u_p^N|_\Sigma$  in the basis  $\{w_k\}$ , we have

$$\begin{aligned} \|T_N(u_p - u_p^N)\|_{H_{00}^{\frac{1}{2}}(\Sigma)}^2 &= \sum_{k \leq N} |\xi_k|^2 (u_p - u_p^N)_k^2 (1 + k^2)^{1/2} \\ &\leq |\xi_N|^2 \sum_{k \leq N} (u_p - u_p^N)_k^2 (1 + k^2)^{1/2} \\ &\leq C(1 + N^2) \|u_p - u_p^N\|_{H_{00}^{\frac{1}{2}}(\Sigma_0)}^2 \\ &\leq C(1 + N^2) \|u_p - u_p^N\|_{H^1(\Omega_b)}^2. \end{aligned}$$

(ii) Using the notations (in particular  $u_k^0$ ) of the proof of Lemma 6.8, we can write

$$\begin{aligned} \|(T - T_N)u_p\|_{H_{00}^{\frac{1}{2}}(\Sigma)}^2 &= \sum_{k > N} |\xi_k|^2 (1 + k^2)^{1/2} (u_k)^2 \\ &\leq C \sum_{k > N} (1 + k^2)^{3/2} (u_k^0)^2 e^{-2\xi_k l} \leq C e^{-2\xi_{N+1} l} \|u_p\|_{H^{3/2}(\Sigma_0)}^2 \\ &\leq C e^{-2\alpha N} \|u_p\|_{H^2(\Omega_b)}^2 \leq C e^{-2\alpha N} \|f\|_{L^2(\Omega_i)}^2 \end{aligned}$$

which concludes the proof. We have used the fact that

$$\left( \sum_{k=1}^{\infty} (1 + k^2)^{3/2} (u_k^0)^2 \right)^{1/2}$$

is a norm equivalent to the norm of the interpolation space  $[H^2(\Sigma_0) \cap H_0^1(\Sigma_0), L^2(\Sigma_0)]_{1/4}$ , which coincides with the space  $H^{\frac{3}{2}}(\Sigma_0) \cap H_0^1(\Sigma_0)$  (cf. Grisvard [14], Theorem 8.1.1).  $\square$

Now we prove  $H^2$  estimates.

**Lemma 6.11** *There exist constants  $C$  independent of  $N$  such that*

$$\|u_p - u_p^N\|_{H^2(\Omega_b)} \leq C N e^{-\alpha N} \|f\|_{L^2(\Omega_i)}. \quad (6.42)$$

$$\|u_p - u_p^N\|_{H^2(\Omega_i)} \leq C N e^{-\alpha N} \|f\|_{L^2(\Omega_i)}. \quad (6.43)$$

*Proof.* The difference  $\tilde{u}^N = u_p - u_p^N$  satisfies:



$$\left\{ \begin{array}{ll} -\Delta \tilde{u}^N + (\beta^2 - n^2 \omega^2) \tilde{u}^N = 0 & \text{in } \Omega_b, \\ \tilde{u}^N = 0 & \text{on } \Gamma_b, \\ \frac{\partial \tilde{u}^N}{\partial \nu} = T_N u_p^N - T u_p & \text{on } \Sigma^\pm. \end{array} \right. \quad (6.44)$$

Classical regularity estimates yield

$$\|u_p - u_p^N\|_{\mathbb{H}^2(\Omega_b)} \leq C \left\{ \|T_N (u_p - u_p^N)\|_{\mathbb{H}_{00}^{\frac{1}{2}}(\Sigma)} + \|(T - T_N) u_p\|_{\mathbb{H}_{00}^{\frac{1}{2}}(\Sigma)} \right\},$$

and estimate (6.43) follows from Lemma 6.10.

In the domain  $\Omega_i^+$  ( $\Omega_i^-$  is treated in the same way),  $\tilde{u}^N$  satisfies

$$\left\{ \begin{array}{ll} -\Delta \tilde{u}^N + (\beta^2 - n^2 \omega^2) \tilde{u}^N = 0 & \text{in } \Omega_i^+, \\ \tilde{u}^N = 0 & \text{on } \Gamma_i^+, \\ \frac{\partial \tilde{u}^N}{\partial \nu} = T_N u_p^N - T u_p & \text{on } \Sigma^+, \end{array} \right.$$

and the same type of argument applies.  $\square$

We use this lemma to establish the main result of this section.

**Theorem 6.9** *There exists a constant  $C$  independent of  $N$  such that*

$$\|S_p - S_p^N\|_{\mathcal{L}(L^2(\Gamma))} \leq C N e^{-\alpha N}, \quad (6.45)$$

$$\|K - K_N\|_{\mathcal{L}(L^2(\Gamma))} \leq C N e^{-\alpha N}. \quad (6.46)$$

*Proof.* Estimate (6.46) is a direct consequence of (6.9), since

$$K - K_N = (S_i - S_e)^{-1} (S_p - S_p^N),$$

while (6.45) is a direct consequence of Lemma 6.11 through trace theorems in  $\mathbb{H}^2$ .  $\square$

A consequence of the previous result is obtained thanks to the approximation theory of compact spectral problems.

**Theorem 6.10** *Let  $\lambda(\omega, \beta)$  be a nonzero eigenvalue of the operator  $K(\omega, \beta)$  with algebraic multiplicity  $m$  and assume the ascent of  $\lambda - K$  is  $\eta$ . Then, there is  $N_0 > 0$  such that, for  $N > N_0$ , there exist  $\lambda_N^1(\omega, \beta), \lambda_N^2(\omega, \beta), \dots, \lambda_N^m(\omega, \beta)$  eigenvalues of  $K_N(\omega, \beta)$  converging to  $\lambda(\omega, \beta)$ . Besides, we have the following estimate*

$$|\lambda(\omega, \beta) - \lambda_N^j(\omega, \beta)| \leq C N e^{-\frac{\alpha N}{\eta}}, \quad \text{for } j = 1, \dots, m,$$

with  $\alpha > 0$ .

### The global error estimate with respect to $R$ and $N$

We are going to obtain error estimates for the eigenvalue approximation  $|\lambda - \lambda_{N,R}|$ . For this purpose, it is necessary to introduce the intermediate operator

$$\tilde{K}_R^N = (S_i - S_e)^{-1} \Pi_R S_p^N \Pi_R.$$

Notice that  $K_R^N = \Pi_R \tilde{K}_R^N$ . Then, using analogous arguments as those in Section 6.1, the idea is the following

- The operator  $\tilde{K}_R^N$  has the same non zero eigenvalues as  $K_R^N$ . We omit the proof of this step because it is analogous to the proof of Theorem 6.11.
- We prove that  $\|K - \tilde{K}_R^N\|$  converges to zero (Theorem 6.11).
- We apply the Osborn's theory to conclude (Theorem 6.12)

**Theorem 6.11** *There exists strictly positive constants  $C_1, C_2, \alpha$  and  $\gamma$  ( $C_1, C_2$  independent of  $N$ ) such that*

$$\|K - K_R^N\| \leq C_1 N e^{-\alpha N} + C_2 e^{-\gamma R}. \quad (6.47)$$

*Proof.* We simply write

$$\|K - \tilde{K}_R^N\| \leq \|K - \tilde{K}_R\| + \|\tilde{K}_R - \tilde{K}_R^N\|.$$

The norm  $\|K - \tilde{K}_R\|$  has been estimated in Theorem 6.2. For the other term, we take into account that

$$\begin{aligned} \|\tilde{K}_R - \tilde{K}_R^N\| &= \|(S_i - S_e)^{-1} \Pi_R S_p \Pi_R - (S_i - S_e)^{-1} \Pi_R S_p^N \Pi_R\| \\ &= \|(S_i - S_e)^{-1} \Pi_R (S_p - S_p^N)\| \\ &\leq C \|S_p - S_p^N\|. \end{aligned}$$

Then estimate (6.47) deduces immediately from (6.6) and (6.45).  $\square$

In the same way as for Theorem 6.2, we have the following consequence of the previous result

**Theorem 6.12** *Let  $\lambda(\omega, \beta)$  be a nonzero eigenvalue of the operator  $K(\omega, \beta)$  with algebraic multiplicity  $m$  and assume the ascent of  $\lambda - K$  is  $\eta$ . Then, there are  $N_0 > 0$  and  $R_0 > 0$  such that, for  $N > N_0$  and  $R > R_0$ , there exist  $\lambda_{N,R}^1(\omega, \beta), \lambda_{N,R}^2(\omega, \beta), \dots, \lambda_{N,R}^m(\omega, \beta)$  eigenvalues of  $K_{N,R}(\omega, \beta)$  converging to  $\lambda(\omega, \beta)$ . Besides, we have the following estimate*

$$|\lambda(\omega, \beta) - \lambda_{N,R}^j(\omega, \beta)| \leq C_1 N e^{\frac{-\alpha N}{\eta}} + C_2 e^{\frac{-\gamma R}{\eta}}, \quad \text{for } j = 1, \dots, m,$$

with  $\alpha > 0$  and  $\gamma > 0$ .

## Acknowledgments

The authors are specially grateful to Prof. Alfredo Bermúdez for his carefully reading of this paper and helpful suggestions.

# Bibliography

- [1] I. BABUŠKA and J. OSBORN. *Eigenvalue problems*, vol. 2 of *Handbook of Numerical Analysis*, pp. 641–787. North-Holland, Amsterdam, 1991.
- [2] A. BAMBERGER and A. S. BONNET. Mathematical analysis of the guided modes of an optical fiber. *SIAM J. Math. Anal.*, **21**(6), pp. 1487–1510, 1990.
- [3] A. S. BONNET BEN DHIA. Analyse mathématique de la propagation de modes guidés dans les fibres optiques. Technical Report 229, Ecole Nationale Supérieure de Techniques Avancées, 1989.
- [4] A. S. BONNET BEN DHIA, G. CALOZ, and F. MAHÉ. Guided modes of integrated optical guides. a mathematical study. *IMA J. of Appl. Math.*, **60**, pp. 225–261, 1998.
- [5] A. S. BONNET BEN DHIA and F. MAHÉ. A guided mode in the range of the radiation modes for a rib waveguide. *J. Opt.*, **28**, pp. 41–43, 1997.
- [6] R. DAUTRAY and J. L. LIONS. *Analyse mathématique et calcul numérique pour les sciences et les techniques*, vol. **1**. Masson, Paris, 1985.
- [7] R. DJELLOUILI. *Contributions à l'Analyse Mathématique et au calcul des modes guidés des fibres optiques*. PhD, Université Paris XI, 1988.
- [8] N. DUNFORD and J. T. SCHWARTZ. *Linear operators*, vol. **2**. Interscience Publishers, New York, 1963.
- [9] J. DUTERTE. *Analyse numérique d'ondes guidées par une perturbation géométrique cylindrique d'un demi-espace élastique homogène*. PhD thesis, Université Paris VI, 1995.

- 
- [10] V. GIRAULT and P. A. RAVIART. *Finite element method for Navier-Stokes equations*. Springer-Verlag, Berlin [etc.], 1986.
- [11] D. GIVOLI and J. B. KELLER. Exact nonreflecting boundary conditions. *J. Comput. Phys.*, 82(1),pp. 172–192, 1989.
- [12] N. GMATI. *Guidage et diffraction d'ondes en milieu non borné. Résolution numérique par une méthode de couplage entre éléments finis et représentation intégrale*. PhD, Université Paris VI, 1992.
- [13] M. D. GÓMEZ PEDREIRA. *A numerical method for the computation of guided waves in integrated optics*. PhD thesis, Universidad de Santiago de Compostela, 1999.
- [14] P. GRISVARD. Caractérisation de quelques espaces d'interpolation. *Arch. Rat. Mech. Anal.*, 25,pp. 40–63, 1967.
- [15] P. GRISVARD. *Elliptic problems in nonsmooth domains*. Pitman Advanced Publishing Program, 1985.
- [16] C. JOHNSON and J.C. NEDELEC. On the coupling of boundary integral and finite element methods. *Math. Comp.*, **35**,pp. 1036–1079, 1980.
- [17] P. JOLY and C. POIRIER. Mathematical analysis of electromagnetic open waveguides. *RAIRO Modél. Math. Anal. Numér.*, **29**(5),pp. 505–575, 1995.
- [18] P. JOLY and C. POIRIER. A numerical method for the computation of electromagnetic modes in optical fibres. *Math. Meth in Appl. Sciences*, 22,pp. 389–447, 1999.
- [19] M. LENOIR and A. JAMI. A variational formulation for exterior problems in linear hydrodynamics. *Comput. Methods Appl. Mech. and Engrg.*, **16**,pp. 341–359, 1978.
- [20] M. LENOIR and A. TOUNSI. The localized finite element method and its application to the two-dimensional sea-keeping problem. *SIAM J. Numer. Anal.*, **25**(4),pp. 729–752, 1988.
- [21] J. L. LIONS and E. MAGENES. *Problèmes aux limites non homogènes et applications*, vol. 1. Dunod, Paris, 1968.

- 
- [22] F. MAHÉ. *Etude Mathématique et numérique de la propagation d'ondes électromagnétiques dans les microguides de l'optique intégrée*. PhD, Université de Rennes I, 1993.
- [23] M. MASMOUDI. Numerical solution for exterior problems. *Numer. Math.*, 51(1),pp. 87–101, 1987.
- [24] J. E. OSBORN. Spectral approximation for compact operators. *Mathematics of computation*, **29**(131),pp. 712–725, 1975.
- [25] H. PICQ. *Détermination et calcul numérique de la première valeur propre d'opérateurs de Schrödinger dans le plan*. PhD thesis, Université de Nice, 1979.
- [26] J. RAZAFIARIVELO. *Optimisation de la forme de transitions entre guides électromagnétiques par une méthode intégrale d'éléments finis*. PhD thesis, Université de Paris VI, 1995.
- [27] M. REED and B. SIMON. *Methods of Modern Mathematical Physics*, vol. **4**. Academic Press, New York, 1978.
- [28] M. SCHECHTER. *Operators Methods in quantum Mechanics*. Elsevier Science Publishing Co., New York [etc.], 1981.
- [29] A. W. SNYDER and J. LOVE. *Optical waveguide theory*. Chapman and Hall, London [etc.], 1983.
- [30] H. P. URBACH. Analysis of the domain integral operator for anisotropic dielectric waveguides. *SIAM J. Math. Anal.*, 27(1),pp. 204–220, 1996.
- [31] C. VASSALLO. *Théorie des guides d'ondes électromagnétiques*. CNET-ENST, Paris, 1985.
- [32] R. WEDER. *Spectral and scattering theory for wave propagation in perturbed stratified media*, vol. **87** of *Applied Mathematical Sciences*. Springer-Verlag, Berlin [etc.], 1991.
- [33] C. H. WILCOX. *Sound propagation in stratified fluids*, vol. **50** of *Applied Mathematical Sciences*. Springer Verlag, 1984.



---

Unit e de recherche INRIA Lorraine, Technop le de Nancy-Brabois, Campus scientifique,  
615 rue du Jardin Botanique, BP 101, 54600 VILLERS L ES NANCY  
Unit e de recherche INRIA Rennes, Irisa, Campus universitaire de Beaulieu, 35042 RENNES Cedex  
Unit e de recherche INRIA Rh ne-Alpes, 655, avenue de l'Europe, 38330 MONTBONNOT ST MARTIN  
Unit e de recherche INRIA Rocquencourt, Domaine de Voluceau, Rocquencourt, BP 105, 78153 LE CHESNAY Cedex  
Unit e de recherche INRIA Sophia-Antipolis, 2004 route des Lucioles, BP 93, 06902 SOPHIA-ANTIPOLIS Cedex

---

 diteur  
INRIA, Domaine de Voluceau, Rocquencourt, BP 105, 78153 LE CHESNAY Cedex (France)  
<http://www.inria.fr>  
ISSN 0249-6399