# Numerical realization of Dirichlet-to-Neumann transparent boundary conditions for photonic crystal wave-guides

Dirk Klindworth[a,*], Kersten Schmidt[a], Sonia Fliss[b]

[a]*DFG Research Center MATHEON, TU Berlin, Berlin, Germany*
[b]*Laboratoire POEMS, UMR 7231 CNRS/ENSTA/INRIA, ENSTA ParisTech, Paris, France*

## Abstract

The computation of guided modes in photonic crystal wave-guides is a key issue in the process of designing devices in photonic communications. Existing methods, such as the super-cell method, provide an efficient computation of well-confined modes. However, if the modes are not well-confined, the modelling error of the super-cell method becomes prohibitive and advanced methods applying transparent boundary conditions for periodic media are needed. In this work we demonstrate the numerical realization of a recently proposed Dirichlet-to-Neumann approach and compare the results with those of the super-cell method. For the resulting non-linear eigenvalue problem we propose an iterative solution based on Newton's method and a direct solution using Chebyshev interpolation of the non-linear operator. Based on the Dirichlet-to-Neumann approach, we present a formula for the group velocity of guided modes that can serve as an objective function in the optimization of photonic crystal wave-guides.

*Keywords:* Dirichlet-to-Neumann map, photonic crystal wave-guide, super-cell method, high-order FEM, non-linear eigenvalue problem, Newton's method, Chebyshev interpolation, group velocity

## 1. Introduction

Photonic crystals (PhC) are nanostructures with a periodic refractive index [1]. The periodicity, which is in the order of the wavelength of light, is induced by periodically spaced holes in an otherwise homogeneous medium. A typical approximation, the so called *2D planar PhC*, of this three dimensional structure is obtained by assuming invariance along the direction of the holes. PhCs have been studied extensively due to their ability to tailor the propagation of light, see for example [2–12] and the references therein.

Of particular interest are *2D planar PhC wave-guides* which are obtained by introducing a line defect in a 2D planar PhC. Light is guided efficiently along the line defect of a 2D planar PhC wave-guide, while decaying exponentially in the PhC. For homogeneous line defects the existence of these so called *guided modes* was shown in [13], while the mathematical justification of this observation in full generality is still under investigation. An important feature of 2D planar PhC wave-guides is the possibility to tailor the dispersion of guided modes, and hence, obtaining, for example, slow light modes [14, 15], *i. e.* guided modes with a small group velocity. Slow light modes lead to a simultaneous enhancement of the light intensity and are thus relevant for the construction of devices in non-linear optics [16].

For 2D planar PhC wave-guides with infinite extend, that we will deal with in this work, a plane wave expansion [17] as used in [11] for the homogeneous exterior domain of finite PhC wave-guides is not appropriate since it cannot account for the periodicity of the infinite medium. For the computation of guided modes the super-cell approach [18, 19] has proven to be an efficient yet reliable method if the modes are well-confined, *i. e.* decay exponentially inside the PhCs with large decay rate. However, if the guided mode is not well-confined, the computational complexity of the super-cell method increases significantly. To overcome this problem, transparent boundary conditions for periodic media serve as an alternative that allow for an *exact* computation of guided modes without introducing a modelling error.

Transparent boundary conditions for periodic media using a Dirichlet-to-Neumann (DtN) approach were introduced in [20], see also [21–23]. Very recently, this approach was rigorously extended to the computation of guided modes in PhC wave-guides [24]. In this paper we want to elaborate on the numerical implementation of the approach in [24] for the exact computation of guided modes in PhC wave-guides, extend it by some important features, such as the computation of the group velocity of guided modes, and compare it to the results of the super-cell method.

---

*Corresponding author. Address: DFG Research Center MATHEON, TU Berlin, Secr. MA 6-4, Straße des 17. Juni 136, 10623 Berlin, Germany. Telephone: +49 (0)30 314-25192. Telefax: +49 (0)30 314-21110
*Email address:* klindworth@math.tu-berlin.de (Dirk Klindworth)

This paper is organized as follows: In Section 2 we describe the model problem, introduce function spaces and summarize results of the spectral theory of the model problem. Furthermore, we provide in this section a short summary of the super-cell method. In Section 3 we introduce a non-linear eigenvalue problem for the computation of guided modes in PhC wave-guides using DtN operators and explain the computation of these DtN operators via local cell problems and a so called *Ricatti equation*. We finish this section with a proof of the differentiability of these DtN operators. In Section 4 we derive a formula for the group velocity of guided modes, that can for example be used to efficiently search for slow light modes. The discretization of the problem using the *finite element method* (FEM) is explained in Section 5. With the help of this discretization we will also show in Section 5 how to solve the discrete form of the Ricatti equation. We continue with solution techniques for the non-linear eigenvalue problem in Section 6 before we present numerical results in Section 7. Finally, we give some concluding remarks in Section 8.

## 2. Model problem

### 2.1. The geometry of photonic crystal wave-guides

A general approach to describe the medium of a 2D PhC wave-guide in the above mentioned configuration is a piecewise definition of its permittivity $\varepsilon$, where the permittivity in the holes is equal to the vacuum permittivity and in the bulk it takes some constant value. Let us first consider two (infinite) PhCs in 2D which are characterized by their periodic permittivities $\varepsilon_{\mathrm{PhC}}^{\pm} : \mathbb{R}^2 \to \mathbb{R}^+ \setminus \{0\}$ that satisfy $\varepsilon_{\mathrm{PhC}}^{\pm}(\mathbf{x} + \mathbf{a}_i^{\pm}) = \varepsilon_{\mathrm{PhC}}^{\pm}(\mathbf{x})$ with the periodicity vectors $\mathbf{a}_i^{\pm} \in \mathbb{R}^2$, $i = 1, 2$, where we assume that $\mathbf{a}_1^+ = \mathbf{a}_1^-$ and w.l.o.g. $\mathbf{a}_1 := \mathbf{a}_1^{\pm} = a_1 (1, 0)^{\mathrm{T}}$, $a_1 > 0$. The periodicity vectors $\mathbf{a}_2^{\pm}$, however, do not need to be identical or parallel, and neither do they have to be orthogonal to $\mathbf{a}_1$. Here and in the sequel, the superscript "+" indicates quantities related to the PhC on top of the guide and the superscript "−" indicates quantities related to the PhC below the guide.

Note that for a PhC with square lattice the periodicity vectors are equal to the unit vectors in $\mathbb{R}^2$ scaled with the length of the square. For an infinite PhC with hexagonal lattice, however, the periodicity vectors can either be chosen as two orthogonal vectors of different length or as two vectors of the same length that are not mutually orthogonal, *e. g.* $\mathbf{a}_1 = (1, 0)^{\mathrm{T}}$ and $\mathbf{a}_2 = (0.5, \sqrt{0.75})^{\mathrm{T}}$, [25].

Moreover, we consider a line defect of height $a_2^0$ characterized by the permittivity $\varepsilon_{\mathrm{defect}} : \mathbb{R} \times ] - a_2^0/2, a_2^0/2[ \to \mathbb{R}^+ \setminus \{0\}$ which is periodic in $\mathbf{a}_1$-direction, *i. e.* $\varepsilon_{\mathrm{defect}}(\mathbf{x} + \mathbf{a}_1) = \varepsilon_{\mathrm{defect}}(\mathbf{x})$. Usually, the line defect of a PhC wave-guide has constant permittivity, *i. e.* there are no holes in the guide. In fact, in the numerical examples in Section 7 the permittivity $\varepsilon_{\mathrm{defect}}$ in the guide is chosen to be constant.

Then we can define the permittivity $\varepsilon$ of the PhC wave-guide by

$$\varepsilon(\mathbf{x}) = \begin{cases} \varepsilon_{\mathrm{PhC}}^-(\mathbf{x}), & \text{if } \mathbf{x} \in \Omega^- := \mathbb{R} \times ] - \infty, -a_2^0/2[, \\ \varepsilon_{\mathrm{defect}}(\mathbf{x}), & \text{if } \mathbf{x} \in \Omega^0 := \mathbb{R} \times ] - a_2^0/2, a_2^0/2[, \\ \varepsilon_{\mathrm{PhC}}^+(\mathbf{x}), & \text{if } \mathbf{x} \in \Omega^+ := \mathbb{R} \times ] a_2^0/2, \infty[. \end{cases} \tag{1}$$

### 2.2. Model problem of finding guided modes

In PhC wave-guides there exist guided modes, sometimes also called trapped modes, which are eigensolutions of the time-harmonic Maxwell's equations and which propagate along the line-defect (*i. e.* along the $x_1$-axis) while decaying in the directions orthogonal to the line defect (*i. e.* along the $x_2$-axis).

It is a well known fact that in 2D the time-harmonic Maxwell's equations decouple into a transverse magnetic (TM) and a transverse electric (TE) mode that satisfy a 2D linear Helmholtz equation [1, 19]. For simplicity let us consider the TM-mode, for which

$$-\Delta E(\mathbf{x}) - \omega^2 \varepsilon(\mathbf{x}) E(\mathbf{x}) = 0, \qquad \mathbf{x} \in \mathbb{R}^2,$$

defines the electric field $E$ in $x_3$-direction.

Note that all results of this article can easily be transferred to the TE-mode for which the magnetic field $H$ in $x_3$-direction satisfies
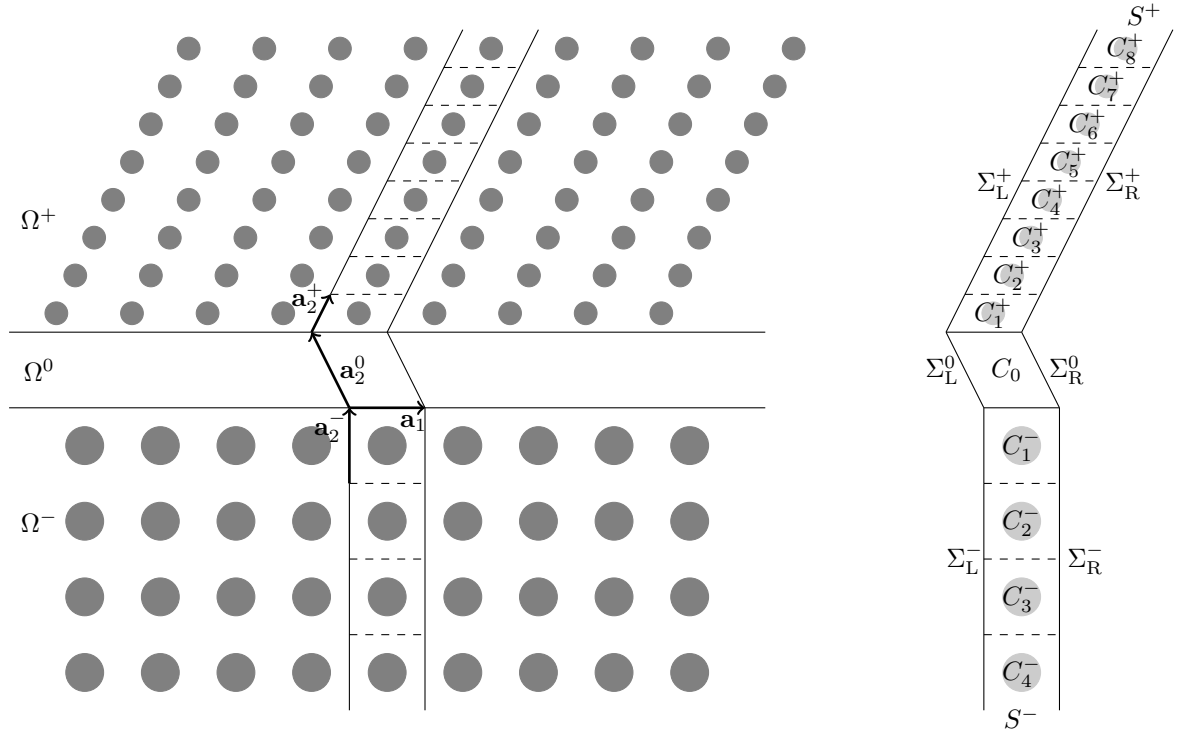
$$-\nabla \cdot \frac{1}{\varepsilon(\mathbf{x})} \nabla H(\mathbf{x}) - \omega^2 H(\mathbf{x}) = 0, \qquad \mathbf{x} \in \mathbb{R}^2.$$

Using the Floquet-Bloch theory [26] we deduce that guided modes are represented by Bloch modes $e_k$ that satisfy

$$-\Delta e_k(\mathbf{x}) - \omega^2 \varepsilon(\mathbf{x}) e_k(\mathbf{x}) = 0 \tag{2a}$$

in the infinite strip $S = S^+ \cup C_0 \cup S^- \subset \mathbb{R}^2$, *c. f.* Figure 1, with quasi-periodic boundary conditions

$$e_k |_{\Sigma_{\mathrm{R}}} = \mathrm{e}^{\mathrm{i}k|\mathbf{a}_1|} e_k |_{\Sigma_{\mathrm{L}}}, \qquad \partial_{\mathbf{n}} e_k |_{\Sigma_{\mathrm{R}}} = -\mathrm{e}^{\mathrm{i}k|\mathbf{a}_1|} \partial_{\mathbf{n}} e_k |_{\Sigma_{\mathrm{L}}}, \tag{2b}$$

(a) Sketch of the domain $\Omega = \Omega^+ \cup \Omega^0 \cup \Omega^-$ of the PhC wave-guide and its periodicity vectors.

(b) Sketch of the periodicity strip $S = S^+ \cup C_0 \cup S^-$

Figure 1: Sketch of the PhC wave-guide and the periodicity strip $S$.

where $\Sigma_\mathrm{L} = \Sigma_\mathrm{L}^+ \cup \Sigma_\mathrm{L}^0 \cup \Sigma_\mathrm{L}^- \subset \partial S$ denotes the boundary of $S$ at the left and $\Sigma_\mathrm{R} = \Sigma_\mathrm{R}^+ \cup \Sigma_\mathrm{R}^0 \cup \Sigma_\mathrm{R}^- \subset \partial S$ is the boundary of $S$ at the right, the parameter $k \in B$ is the so-called quasi-momentum in the one-dimensional Brillouin zone $B = [-\pi/|\mathbf{a}_1|, \pi/|\mathbf{a}_1|]$, and the operator $\partial_\mathbf{n}$ denotes the normal derivative, *i. e.* $\partial_\mathbf{n} = \mathbf{n} \cdot \nabla$ with the unit normal vector $\mathbf{n}$ outward to the domain $S$.

*2.3. Spectral properties*

In order to rewrite Eq. (2), that can be regarded as an eigenvalue problem in $\omega^2$ and $k$ respectively, in a rigorous way, we introduce the following function spaces. Let $H^1(S)$ be the usual space of square integrable functions in $S$ whose gradient is also square integrable. Then we define the periodic space

$$H_\mathrm{per}^1(S) := \left\{ u \in H^1(S) \text{ with } u|_{\Sigma_\mathrm{L}} = u|_{\Sigma_\mathrm{R}} \right\}.$$

Moreover, let $H^1(\Delta, S)$ be the subspace of $H^1(S)$ with functions whose Laplacian is square integrable. Then we define

$$H_\mathrm{per}^1(\Delta, S) := \left\{ u \in H^1(\Delta, S) \cap H_\mathrm{per}^1(S) \text{ with } \partial_\mathbf{n} u|_{\Sigma_\mathrm{L}} = -\partial_\mathbf{n} u|_{\Sigma_\mathrm{R}} \right\}.$$

Let $\Gamma_0^\pm = \partial S^\pm \cap \partial C_0$ denote the boundaries between the top and bottom half-strips $S^\pm$ and the channel cell $C_0$. Then we define $H_\mathrm{per}^{1/2}(\Gamma_0^\pm)$ as the trace of $H_\mathrm{per}^1(S^\pm)$ on the boundary $\Gamma_0^\pm$ and $H_\mathrm{per}^{-1/2}(\Gamma_0^\pm)$ as its dual space.

With these definitions and the substitution $e_k(\mathbf{x}) = \mathrm{e}^{\mathrm{i}kx_1} u(\mathbf{x})$, the eigenvalue problem (2) is equivalent to: find couples $(\omega^2, k) \in \mathbb{R}^+ \times B$ such that there exists a non-trivial $u \in H_\mathrm{per}^1(\Delta, S)$ that satisfies

$$-(\nabla + \mathrm{i}k\mathbf{a}_1) \cdot (\nabla + \mathrm{i}k\mathbf{a}_1)u(\mathbf{x}) - \omega^2 \varepsilon(\mathbf{x})u(\mathbf{x}) = 0, \qquad \mathbf{x} \in S. \qquad (3)$$

Then we call $(\omega^2, k)$ an *eigenvalue couple* of (3) with associated eigenfunction $u$. This eigenvalue problem is linear in $\omega^2$ when fixing $k \in B$, the so called $\omega$*-formulation*, and quadratic in $k$ when fixing $\omega \in \mathbb{R}^+$, the so called $k$*-formulation*. However, note that this problem is posed on the unbounded domain $S$.

In the remainder of this section we summarize the results of the spectral theory of this problem. To this end, we introduce the operators

$$\mathcal{A}_\mathrm{PhC}^+(k) := -\frac{1}{\varepsilon_\mathrm{PhC}^+}(\nabla + \mathrm{i}k\mathbf{a}_1) \cdot (\nabla + \mathrm{i}k\mathbf{a}_1),$$

3

$k \in B$, defined on $H^1_{\text{per}}(\Delta, S^+)$, and

$$\mathcal{A}^-_{\text{PhC}}(k) := -\frac{1}{\varepsilon^-_{\text{PhC}}}(\nabla + \mathrm{i}k\mathbf{a}_1) \cdot (\nabla + \mathrm{i}k\mathbf{a}_1),$$

$k \in B$, defined on $H^1_{\text{per}}(\Delta, S^-)$, that are related to eigenvalue problems of the form (3) posed on the infinite half strips $S^\pm$ with perfectly periodic permittivities $\varepsilon^\pm_{\text{PhC}}$ of the top and bottom PhC respectively. Furthermore, we introduce the operator

$$\mathcal{A}(k) := -\frac{1}{\varepsilon}(\nabla + \mathrm{i}k\mathbf{a}_1) \cdot (\nabla + \mathrm{i}k\mathbf{a}_1),$$

$k \in B$, defined on the function space $H^1_{\text{per}}(\Delta, S)$ of the whole strip $S$, that is related to the eigenvalue problem (3) with permittivity $\varepsilon$ as defined in Eq. (1).

Using the same arguments as in [26, 27] for operators with perfectly periodic coefficients and applying the Weyl theorem [28] as in [24], we can show

**Proposition 2.1.** *The essential spectrum* $\sigma_{\text{ess}}(k)$, $k \in B$, *of the operator* $\mathcal{A}(k)$ *satisfies*

$$\sigma_{\text{ess}}(k) = \sigma^+_{\text{ess}}(k) \cup \sigma^-_{\text{ess}}(k)$$

*with*

$$\sigma^\pm_{\text{ess}}(k) := \sigma(\mathcal{A}^\pm_{\text{PhC}}(k)) = \mathbb{R}^+ \setminus \bigcup_{n=1}^{N^\pm(k)} I^\pm_n(k),$$

*where* $I^\pm_n(k) \subset \mathbb{R}^+$ *are open intervals, the so called band gaps, and* $N^\pm(k)$ *is the number of band gaps.*

According to a recent result [29] on the Bethe-Sommerfeld conjecture [30] for periodic Maxwell operators in 2D, we note that the numbers $N^\pm(k)$ of band gaps are finite.

We conclude that there exists a finite number of band gaps $I_n(k) \subset \mathbb{R}^+$, $n = 1, \ldots, N(k)$, such that

$$\sigma_{\text{ess}}(k) = \mathbb{R}^+ \setminus \bigcup_{n=1}^{N(k)} I_n(k),$$

with the set of band gaps $\bigcup_{n=1}^{N(k)} I_n(k) = \bigcup_{n=1}^{N^+(k)} I^+_n(k) \cap \bigcup_{n=1}^{N^-(k)} I^-_n(k)$.

Using the theory of selfadjoint operators [28], we deduce

**Proposition 2.2.** *Inside the band gaps* $I_n(k)$, $n = 1, \ldots, N(k)$, $k \in B$, *of the operator* $\mathcal{A}(k)$ *related to the PhC wave-guide eigenvalue problem (3), there exist only isolated eigenvalues of finite multiplicity which can only accumulate at the boundaries of the band gaps* $I_n(k)$.

Let the isolated eigenvalues $\omega^2_m(k)$, $m = 1, \ldots, M(k)$, of $\mathcal{A}(k)$ inside the band gaps $I_n$, $n = 1, \ldots, N(k)$, be ordered such that

$$0 \le \omega^2_1(k) \le \ldots \le \omega^2_{M(k)}(k)$$

with $0 \le M(k) \le \infty$. Then we can define so called *dispersive curves* (of the first kind)

$$f^{(1)}_m(k) = \omega^2_m(k)$$

that are $\frac{2\pi}{|\mathbf{a}_1|}$-periodic, even and continuous [24]. On the other hand, using [31] we can also prove

**Proposition 2.3.** *There exists an alternative ordering* $m \mapsto \widetilde{m}(m)$ *of the eigenvalues* $\omega^2_m(k)$, $m = 1, \ldots, M(k)$ *such that the dispersive curves of the second kind*

$$f^{(2)}_m(k) = \omega^2_{\widetilde{m}(m)}(k),$$

$m = 1, \ldots, M(k)$, *are continuous and differentiable.*

Finally, we shall state a result related to the eigenfunctions of $\mathcal{A}(k)$ proven in [27].

**Proposition 2.4.** *Let* $k \in B$ *and let* $\omega^2 \in \sigma(k) \setminus \sigma_{\text{ess}}(k)$ *be an eigenvalue of (3). Then the associated eigenfunction* $u \in H^1_{\text{per}}(\Delta, S)$ *decays exponentially with* $|x_2|$, *where the decay rate is proportional to the distance of the eigenvalue* $\omega^2$ *to the essential spectrum* $\sigma_{\text{ess}}(k)$.

Note that this result is the mathematical justification for the super-cell method which we will explain in the following section.

*2.4. Model reduction using the the super-cell approach*

The frequently used super-cell method [18, 19] provides access to approximations of guided modes of the eigenvalue problem (3). Based on the observation in Proposition 2.4, that guided modes decay exponentially with $x_2 \to \pm\infty$, the eigenvalue problem (3) is posed on a bounded super-cell $S_{\mathrm{sc}} \subset S$ instead of the infinite strip $S$. The computational domain $S_{\mathrm{sc}}$ is obtained by simply cutting the infinite strip $S$ after $n \in \mathbb{N}$ periodicity cells of the PhCs on top and bottom, and prescribing periodic boundary conditions at the top and bottom boundaries $\Gamma_n^{\pm}$. Thus, the problem solved reads

$$-(\nabla + \mathrm{i}k\mathbf{a}_1) \cdot (\nabla + \mathrm{i}k\mathbf{a}_1)u(\mathbf{x}) - \omega^2 \varepsilon(\mathbf{x})u(\mathbf{x}) = 0, \qquad \mathbf{x} \in S_{\mathrm{sc}}, \tag{4a}$$

with

$$u|_{\Gamma_n^+} = u|_{\Gamma_n^-}, \qquad \partial_2 u|_{\Gamma_n^+} = \partial_2 u|_{\Gamma_n^-} . \tag{4b}$$

Since guided modes decay exponentially if $x_2 \to \pm\infty$, *c. f.* Proposition 2.4, it can be expected that the modelling error, that is introduced when prescribing periodic boundary conditions after a certain number of holes, is reasonably small. In fact, Soussi [18] showed that the solutions of the super-cell method converge exponentially towards the solutions of the exact problem (3) if the number $n$ of periodicity cells, that are included in the super-cell, tends to infinity.

The eigenvalue problem (4) of the super-cell method is again linear in $\omega^2$ ($\omega$-formulation) and quadratic in $k$ ($k$-formulation) allowing for standard numerical techniques of PhC band computation to be applied [19]. However, note that the eigenvalue problem (4) of the super-cell method has eigenvalues also inside the essential spectrum which have to be excluded. To do so, one needs to have access to the essential spectrum in advance, *i. e.* a full computation of the spectra $\sigma_{\mathrm{ess}}^{\pm}(k)$ of the operators $\mathcal{A}_{\mathrm{PhC}}^{\pm}(k)$ related to the PhCs on top and bottom of the guide is needed.

The super-cell method provides access to an approximation to the solutions of (3) while introducing a modelling error. However, as we shall see from the numerical results in Section 7, it is a challenging task to fix the number $n$ of periodicity cells to be included in the super-cell if no a priori knowledge about the confinement of the guided mode is available. Moreover, note that the number of eigenvalues inside the essential spectrum grows with the number of periodicity cells that are included in the super-cell. Using an iterative eigenvalue solver, this implies that one should restrict the eigenvalue computation — as far as possible — to the band gap by a shift and invert strategy, in order not to spoil the performance of the iterative solver.

In the following section we introduce a DtN approach [24] that provides access to an *exact* approximation in the sense that no modelling error is introduced. Furthermore, this new method automatically provides access to the essential spectra $\sigma_{\mathrm{ess}}^{\pm}(k)$ of the operators $\mathcal{A}_{\mathrm{PhC}}^{\pm}(k)$ related to the PhCs on top and bottom of the guide.

## 3. Non-linear eigenvalue problem using Dirichlet-to-Neumann operators

In this section we will define DtN operators for the two infinite half strips $S^{\pm}$ on top and bottom of the guide. These DtN operators then allow for a computation of the eigenvalues and eigenfunctions of (3) by solving a non-linear eigenvalue problem in the bounded domain $C_0$ in the wave-guide. Therefore, this method provides access to the eigenvalues of (3) without introducing a modelling error but for the price that the eigenvalue problem becomes non-linear.

We start by defining the DtN operators, continue with the computation of these operators, and finally, compute the derivatives of the DtN operators with respect to $\omega$ and $k$, which we need for the computation of the group velocity in Section 4.

*3.1. Definition of the Dirichlet-to-Neumann operators*

As a first step towards the definition of the DtN operators, we introduce Dirichlet problems in the infinite half-strips $S^{\pm}$: for any $\varphi \in H_{\mathrm{per}}^{1/2}(\Gamma_0^{\pm})$ find $u^{\pm} \equiv u^{\pm}(\mathbf{x}; \omega, k, \varphi) \in H_{\mathrm{per}}^1(\Delta, S^{\pm})$ such that

$$-(\nabla + \mathrm{i}k\mathbf{a}_1) \cdot (\nabla + \mathrm{i}k\mathbf{a}_1)u^{\pm} - \omega^2 \varepsilon(\mathbf{x})u^{\pm} = 0, \qquad \mathbf{x} \in S^{\pm},$$
$$u^{\pm}|_{\Gamma_0^{\pm}} = \varphi. \tag{5}$$

Theorem 4.1 in [24] gives the following result.

**Theorem 3.1.** *If $\omega^2 \notin \sigma_{\mathrm{ess}}^{\pm}(k)$, the problem (5) is well-posed in $H_{\mathrm{per}}^1(\Delta, S^{\pm})$, except for a countable set of frequencies $\omega$.*

**Remark 3.2.** *There exist symmetry properties of the unit cells $C_n^\pm$ such that problem (5) is well-posed in $H_{\mathrm{per}}^1(\Delta, S^\pm)$ for all $\omega^2 \notin \sigma_{\mathrm{ess}}^\pm(k)$ if these symmetry properties are fulfilled [24]. However, if these symmetry properties are not fulfilled, the Dirichlet boundary condition can be replaced by the Robin boundary condition*

$$(\partial_{\mathbf{n}} + \mathrm{i}\alpha)\, u^\pm |_{\Gamma_0^\pm} = \varphi, \qquad \alpha \in \mathbb{R},$$

*at $\Gamma_0^\pm$ yielding a well-posed problem for any $\omega^2 \notin \sigma_{\mathrm{ess}}^\pm(k)$ [23, 32]. Taking the trace $(\partial_{\mathbf{n}} - \mathrm{i}\alpha)\, u^\pm$ on $\Gamma_0^\pm$ will then lead to Robin-to-Robin operators instead of DtN operators.*

In the sequel we shall assume that the problem (5) is well-posed. Then, for any $\varphi \in H_{\mathrm{per}}^{1/2}(\Gamma_0^\pm)$, the DtN operators $\Lambda^\pm(\omega, k) \in \mathcal{L}(H_{\mathrm{per}}^{1/2}(\Gamma_0^\pm), H_{\mathrm{per}}^{-1/2}(\Gamma_0^\pm))$ are defined as

$$\Lambda^\pm(\omega, k)\varphi = \mp \partial_2 u^\pm(\cdot\,; \omega, k, \varphi)|_{\Gamma_0^\pm}. \tag{6}$$

**Proposition 3.3** (Proposition 4.3 in [24])**.** *Let $k \in B$ and $\omega^2 \in \mathbb{R}^+ \setminus \sigma_{\mathrm{ess}}^\pm(k)$, then the DtN operators $\Lambda^\pm(\omega, k)$ are continuous from $H_{\mathrm{per}}^{1/2}(\Gamma_0^\pm)$ onto $H_{\mathrm{per}}^{-1/2}(\Gamma_0^\pm)$ and their norms are continuous with respect to $\omega \in \mathbb{R}^+ \setminus \sigma_{\mathrm{ess}}^\pm(k)$.*

Now we can state the main result of the DtN method.

**Theorem 3.4** (Theorem 4.4 in [24])**.** *Let $C_0 = S \cap \Omega^0$. Then the problem (3) is equivalent to: find couples $(\omega, k) \in \mathbb{R}^+ \times B$ such that there exists a non-trivial $u \in H_{\mathrm{per}}^1(\Delta, C_0)$ that satisfies*

$$-(\nabla + \mathrm{i}k\mathbf{a}_1) \cdot (\nabla + \mathrm{i}k\mathbf{a}_1)u(\mathbf{x}) - \omega^2 \varepsilon(\mathbf{x})u(\mathbf{x}) = 0, \qquad \mathbf{x} \in C_0, \tag{7a}$$

$$-\partial_2 u(\mathbf{x}) = \Lambda^+(\omega, k)u(\mathbf{x}), \qquad \mathbf{x} \in \Gamma_0^+, \tag{7b}$$

$$\partial_2 u(\mathbf{x}) = \Lambda^-(\omega, k)u(\mathbf{x}), \qquad \mathbf{x} \in \Gamma_0^-. \tag{7c}$$

Note that the problem (7) — in comparison to problem (3) — is posed on the bounded domain $C_0$ but it is non-linear with respect to $\omega$ and $k$.

The weak formulation of the non-linear eigenvalue problem (7) is: find couples $(\omega^2, k) \in \mathbb{R}^+ \times B$ such that there exists a non-trivial $u \in H_{\mathrm{per}}^1(C_0)$ that satisfies

$$\mathsf{a}_{C_0}(u, v; k) - \omega^2 \mathsf{m}_{C_0}(u, v) + \mathsf{d}_{C_0}(u, v; \omega, k) = 0 \tag{8}$$

for all $v \in H_{\mathrm{per}}^1(C_0)$, where the sesquilinear forms are defined as

$$\mathsf{a}_{C_0}(u, v; k) := \int_{C_0} (\nabla + \mathrm{i}k\mathbf{a}_1)u \cdot (\nabla - \mathrm{i}k\mathbf{a}_1)\overline{v}\, \mathrm{d}\mathbf{x}, \tag{9a}$$

$$\mathsf{m}_{C_0}(u, v) := \int_{C_0} \varepsilon u\overline{v}\, \mathrm{d}\mathbf{x}, \tag{9b}$$

$$\mathsf{d}_{C_0}(u, v; \omega, k) := \int_{\Gamma_0^+} \Lambda^+(\omega, k)u\, \overline{v}\, \mathrm{d}s(\mathbf{x}) + \int_{\Gamma_0^-} \Lambda^-(\omega, k)u\, \overline{v}\, \mathrm{d}s(\mathbf{x}). \tag{9c}$$

Furthermore — for simplicity of notation — we also define

$$\mathsf{b}_{C_0}(u, v; \omega, k) := \mathsf{a}_{C_0}(u, v; k) - \omega^2 \mathsf{m}_{C_0}(u, v). \tag{9d}$$

The DtN operators in weak form, that appear in Eq. (9c), satisfy

$$\int_{\Gamma_0^\pm} \Lambda^\pm(\omega, k)\varphi\, \overline{\psi}\, \mathrm{d}s(\mathbf{x}) = \int_{S^\pm} (\nabla + \mathrm{i}k\mathbf{a}_1)u^\pm(\cdot\,; \omega, k, \varphi) \cdot (\nabla - \mathrm{i}k\mathbf{a}_1)\overline{u^\pm}(\cdot\,; \omega, k, \psi)\, \mathrm{d}\mathbf{x}$$

$$- \omega^2 \int_{S^\pm} \varepsilon u^\pm(\cdot\,; \omega, k, \varphi)\overline{u^\pm}(\cdot\,; \omega, k, \psi)\, \mathrm{d}\mathbf{x} \quad (10)$$

for any $\varphi, \psi \in H_{\mathrm{per}}^{1/2}(\Gamma_0^\pm)$.

Again we note that (8) is a non-linear eigenvalue problem for both, the $\omega$-formulation and the $k$-formulation.

Before we elaborate on the computation of the DtN operators, let us introduce a "simplified" eigenvalue problem with fixed DtN operators. Let $(\omega_\Lambda^2, k_\Lambda) \in \mathbb{R}^+ \times B$ be arbitrary but fixed. Then the problem: find couples $(\omega^2, k) \in \mathbb{R}^+ \times B$ such that there exists a non-trivial $u \in H_{\mathrm{per}}^1(C_0)$ that satisfies

$$\mathsf{a}_{C_0}(u, v; k) - \omega^2 \mathsf{m}_{C_0}(u, v) + \mathsf{d}_{C_0}(u, v; \omega_\Lambda, k_\Lambda) = 0 \tag{11}$$

for all $v \in H_{\mathrm{per}}^1(C_0)$, is a linear eigenvalue problem in $\omega^2$, and a quadratic eigenvalue problem in $k$, whose solution coincides with the one of (8) if $(\omega^2, k) = (\omega_\Lambda^2, k_\Lambda)$. For the linear eigenvalue problem (11) in $\omega$-formulation we state the following important results.

**Proposition 3.5** (Corollary of Proposition 4.6 in [24]). *Let $(\omega_\Lambda^2, k) \in \mathbb{R}^+ \times B$ and $\omega_\Lambda^2 \notin \sigma_{\mathrm{ess}}(k)$. Then the eigenvalues $\omega_m^2(\omega_\Lambda, k)$, $m \in \mathbb{N}$, of the linear eigenvalue problem (11) in $\omega$-formulation are real.*

**Proposition 3.6** (Theorem 4.7 in [24]). *Let $\omega_m^2(\omega_\Lambda, k) \in \mathbb{R}$, $m \in \mathbb{N}$, denote the eigenvalues of the linear eigenvalue problem (11) in $\omega$-formulation with $\omega_m^2 \leq \omega_{m+1}^2$ for all $m \in \mathbb{N}$. Then the functions $g_m^{(1)}(\omega_\Lambda, k) = \omega_m^2(\omega_\Lambda, k)$ are continuous.*

Using the same argument as for Proposition 2.3 we can deduce

**Proposition 3.7.** *There exists an alternative ordering $m \mapsto \widetilde{m}(m)$ of the eigenvalues $\omega_m^2(\omega_\Lambda, k) \in \mathbb{R}$, $m \in \mathbb{N}$, of the linear eigenvalue problem (11) in $\omega$-formulation such that the functions $g_m^{(2)}(\omega_\Lambda, k) = \omega_{\widetilde{m}(m)}^2(\omega_\Lambda, k)$ are continuously differentiable.*

*3.2. Characterization of the Dirichlet-to-Neumann operators*

In Eq. (6) the DtN operators are defined via a Dirichlet problems (5) on an unbounded domain. In this subsection we summarize the results in [20, 24] how to compute the DtN operators via local cell problems, *i.e.* by solving Dirichlet problems on a single periodicity cell, and a stationary Ricatti equation.
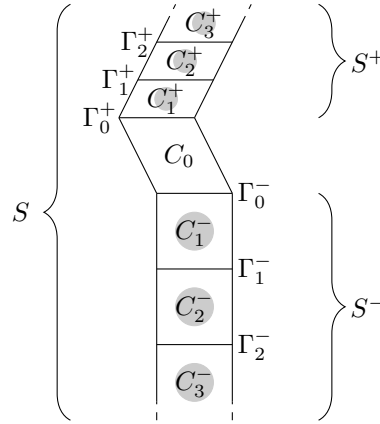


Figure 2: Sketch of the periodicity strip $S$ consisting of the two infinite half strips $S^\pm$ with their periodicity cells $C_n^\pm$, $n \geq 1$, and the channel cell $C_0$, as well as the interfaces $\Gamma_n^\pm$, $n \geq 0$ between the cells.

To this end, we note that the infinite strips $S^\pm$ on top and bottom of the guide can be expressed as union of an infinite number of periodicity cells $C_n^\pm$, $n \in \mathbb{N}$, *i.e.* $S^\pm = \bigcup_{n=1}^\infty C_n^\pm$, *c.f.* Figures 1(b) and 2. The top and bottom boundaries of the cell $C_n^\pm$ shall be denoted by $\Gamma_{n-1}^\pm$ and $\Gamma_n^\pm$, where the subscript "$n-1$" indicates the boundary that is closer to the guide, *i.e.* $\Gamma_0^\pm$ are the boundaries between the periodicity cells $C_1^\pm$ and the cell $C_0$ of the guide, see Figure 2.

We also note that — due to the periodicity and the infinity of the half strips — all cells $C_n^\pm$ can be identified by the first cell $C_1^\pm$ and all boundaries $\Gamma_n^\pm$ can be identified by the first boundary $\Gamma_0^\pm$.

Therefore, let us introduce shift operators $\mathcal{S}_n^\pm \in \mathcal{L}(C^\infty(\Gamma_0^\pm), C^\infty(\Gamma_n^\pm))$, $n \in \mathbb{N}$, defined by

$$\mathcal{S}_n^\pm \varphi(\mathbf{x}) = \varphi(\mathbf{x} \mp n\mathbf{a}_2^\pm) \tag{12}$$

By a density argument of $C^\infty(\Gamma_n^\pm)$ in $H_{\mathrm{per}}^{1/2}(\Gamma_n^\pm)$ and $H_{\mathrm{per}}^{-1/2}(\Gamma_n^\pm)$, respectively, we can extend the shift operators $\mathcal{S}_n^\pm$ to functions in $H_{\mathrm{per}}^{1/2}(\Gamma_n^\pm)$ and $H_{\mathrm{per}}^{-1/2}(\Gamma_n^\pm)$. For simplicity of notation we shall write $\mathcal{S}^\pm := \mathcal{S}_1^\pm$. Furthermore, we introduce the inverse $(\mathcal{S}^\pm)^{-1}$ of $\mathcal{S}^\pm$ which is simply given by

$$(\mathcal{S}^\pm)^{-1}\varphi(\mathbf{x}) = \varphi(\mathbf{x} \pm \mathbf{a}_2^\pm). \tag{13}$$

With the help of these operators we can express the trace of the unique solution $u^\pm(\cdot\,; \omega, k, \varphi)$ of the Dirichlet problem (5) at the edges $\Gamma_n^\pm$, $n \in \mathbb{N}$, as

$$u^\pm(\cdot\,; \omega, k, \varphi)|_{\Gamma_n^\pm} = \mathcal{S}_n^\pm(\mathcal{P}^\pm(\omega, k))^n \varphi,$$

with the *propagation operator* $\mathcal{P}^\pm(\omega, k) \in \mathcal{L}(H_{\mathrm{per}}^{1/2}(\Gamma_0^\pm))$ defined for any $\varphi \in H_{\mathrm{per}}^{1/2}(\Gamma_0^\pm)$ by

$$\mathcal{P}^\pm(\omega, k)\varphi = (\mathcal{S}^\pm)^{-1} u^\pm(\cdot\,; \omega, k, \varphi)|_{\Gamma_1^\pm}.$$

As shown in [20], the propagation operator $\mathcal{P}^\pm(\omega, k)$ is the unique solution of the so-called *Ricatti equation*

$$\mathcal{T}_{10}^\pm(\omega, k)(\mathcal{P}^\pm(\omega, k))^2 + (\mathcal{T}_{00}^\pm(\omega, k) + \mathcal{T}_{11}^\pm(\omega, k))\mathcal{P}^\pm(\omega, k) + \mathcal{T}_{01}^\pm = 0 \tag{14}$$

with spectral radius strictly less than 1. Here, the operators $\mathcal{T}_{ij}^{\pm}(\omega, k) \in \mathcal{L}(H_{\text{per}}^{1/2}(\Gamma_0^{\pm}), H_{\text{per}}^{-1/2}(\Gamma_0^{\pm}))$, $i, j = 0, 1$, are defined by

$$
\begin{aligned}
\mathcal{T}_{00}^{\pm}(\omega, k)\,\varphi &= \partial_{\mathbf{n}} u_0^{\pm}(\cdot\,; \omega, k, \quad \varphi))|_{\Gamma_0^{\pm}}, \\
\mathcal{T}_{01}^{\pm}(\omega, k)\,\varphi &= (\mathcal{S}^{\pm})^{-1} \partial_{\mathbf{n}} u_0^{\pm}(\cdot\,; \omega, k, \quad \varphi))|_{\Gamma_1^{\pm}}, \\
\mathcal{T}_{10}^{\pm}(\omega, k)\,\varphi &= \partial_{\mathbf{n}} u_1^{\pm}(\cdot\,; \omega, k, \mathcal{S}^{\pm}\varphi))|_{\Gamma_0^{\pm}}, \\
\mathcal{T}_{11}^{\pm}(\omega, k)\,\varphi &= (\mathcal{S}^{\pm})^{-1} \partial_{\mathbf{n}} u_1^{\pm}(\cdot\,; \omega, k, \mathcal{S}^{\pm}\varphi))|_{\Gamma_1^{\pm}},
\end{aligned}
$$

for any $\varphi \in H_{\text{per}}^{1/2}(\Gamma_0^{\pm})$, where $u_i^{\pm} \equiv u_i^{\pm}(\mathbf{x}; \omega, k, \varphi) \in H_{\text{per}}^1(\Delta, C_1^{\pm})$, $i = 0, 1$, solve

$$
-(\nabla + \mathrm{i}k\mathbf{a}_1) \cdot (\nabla + \mathrm{i}k\mathbf{a}_1)u_i^{\pm} - \omega^2 \varepsilon(\mathbf{x})u_i^{\pm} = 0, \qquad \mathbf{x} \in C_1^{\pm}, \tag{15a}
$$

with Dirichlet boundary data

$$
u_i^{\pm}|_{\Gamma_j^{\pm}} = \delta_{ij}\varphi. \tag{15b}
$$

Here and in the following, $\delta_{ij}$ denotes the usual Kronecker delta, i. e. $\delta_{ij} = 1$ if $i = j$, and $\delta_{ij} = 0$ if $i \neq j$. Then the DtN operators $\Lambda^{\pm}(\omega, k)$ are given by [20]

$$
\Lambda^{\pm}(\omega, k) = \mathcal{T}_{00}^{\pm}(\omega, k) + \mathcal{T}_{10}^{\pm}(\omega, k)\,\mathcal{P}^{\pm}(\omega, k). \tag{16}
$$

**Remark 3.8** (Well-posedness). *The Dirichlet cell problems (15) are well-posed except for a countable set of frequencies $\omega$, and hence, the operators $\mathcal{T}_{ij}$, $i, j = 0, 1$, are injective for almost any $\omega$ [20] . Moreover, we can show — using the Fredholm theory — that the operators $\mathcal{T}_{00}$, $\mathcal{T}_{11}$ and $\mathcal{T}_{00} + \mathcal{T}_{11}$ are isomorphisms from $H_{\text{per}}^{1/2}(\Gamma_0^{\pm})$ onto $H_{\text{per}}^{-1/2}(\Gamma_0^{\pm})$. On the other hand, the operators $\mathcal{T}_{01}$ and $\mathcal{T}_{10}$ are compact and hence, they are not bijective [32].*

**Remark 3.9** (Robin-type cell problems). *If we equip the local cell problems (15a) with Robin boundary data [23, 32]*

$$
(-\partial_{\mathbf{n}} + \mathrm{i}\alpha)\,u_i^{\pm}|_{\Gamma_j^{\pm}} = \delta_{ij}\varphi, \qquad \alpha \in \mathbb{R},
$$

*instead of Dirichlet boundary data (15b), the local cell problems are well-posed for all frequencies $\omega$. However, for the sake of simplicity we shall continue with Dirichlet cell problems and assume in the following that they are well-posed.*

Finally, let us introduce a weak formulation of the local cell problems (15) and the corresponding DtN like operators $\mathcal{T}_{ij}^{\pm}$, using a Dirichlet lifting ansatz.

To this end, we start by introducing a Dirichlet lift $w_i \equiv w_i(\cdot\,; \varphi) \in H_{\text{per}}^1(C_1^{\pm})$ with $w_i|_{\Gamma_j^{\pm}} = \delta_{ij}\varphi$. Then the weak solutions $u_i^{\pm} \equiv u_i^{\pm}(\cdot\,; \omega, k, \varphi) \in H_{\text{per}}^1(C_1^{\pm})$ of the Dirichlet cell problems (15) can be decomposed into $u_i^{\pm}(\cdot\,; \omega, k, \varphi) = w_i(\cdot\,; \varphi) + u_{i,0}^{\pm}(\cdot\,; \omega, k, \varphi)$, where $u_{i,0}^{\pm} \equiv u_{i,0}^{\pm}(\cdot\,; \omega, k, \varphi) \in H_{\text{per},0}^1(C_1^{\pm}) := \{u \in H_{\text{per}}^1(C_1^{\pm}) \text{ with } u|_{\Gamma_i^{\pm}} = 0, \ i = 0, 1\}$ satisfies

$$
\mathsf{b}_{C_1^{\pm}}(u_{i,0}^{\pm}, v; \omega, k) = -\mathsf{b}_{C_1^{\pm}}(w_i, v; \omega, k) \tag{17}
$$

for all $v \in H_{\text{per},0}^1(C_1^{\pm})$, with the sesquilinear form $\mathsf{b}_{C_1^{\pm}}(\cdot, \cdot\,; \omega, k)$ defined by

$$
\mathsf{b}_{C_1^{\pm}}(u, v; \omega, k) := \int_{C_1^{\pm}} (\nabla + \mathrm{i}k\mathbf{a}_1)u \cdot (\nabla - \mathrm{i}k\mathbf{a}_1)\overline{v} - \omega^2 \varepsilon u\overline{v}\,\mathrm{d}\mathbf{x}. \tag{18}
$$

The DtN-like operators $\mathcal{T}_{ij}^{\pm}(\omega, k)$, $i, j = 0, 1$, then satisfy for any $\varphi, \psi \in H_{\text{per}}^{1/2}(\Gamma_0^{\pm})$

$$
\begin{aligned}
\int_{\Gamma_0^{\pm}} \mathcal{T}_{ij}^{\pm}(\omega, k)\varphi\,\overline{\psi}\,\mathrm{d}s(\mathbf{x}) &= \int_{\Gamma_0^{\pm}} \left((\mathcal{S}^{\pm})^{-1}\right)^j \left[\partial_{\mathbf{n}} u_i^{\pm}(\cdot\,; \omega, k, (\mathcal{S}^{\pm})^i\varphi)\right]\,\overline{\psi}\,\mathrm{d}s(\mathbf{x}) \\
&= \int_{\Gamma_j^{\pm}} \partial_{\mathbf{n}} u_i^{\pm}(\cdot\,; \omega, k, (\mathcal{S}^{\pm})^i\varphi)\,(\mathcal{S}^{\pm})^j\overline{\psi}\,\mathrm{d}s(\mathbf{x}) \\
&= \mathsf{b}_{C_1^{\pm}}(u_i^{\pm}(\cdot\,; \omega, k, (\mathcal{S}^{\pm})^i\varphi), w_j(\cdot\,; (\mathcal{S}^{\pm})^j\psi); \omega, k),
\end{aligned} \tag{19}
$$

where we used the relations (12) and (13), integration by parts and the fact that Eq. (15) implies

$$
\Delta u_i^{\pm}(\cdot\,; \omega, k, \varphi) = \left(k^2|\mathbf{a}_1|^2 - 2\mathrm{i}k\mathbf{a}_1 \cdot \nabla - \omega^2 \varepsilon\right) u_i^{\pm}(\cdot\,; \omega, k, \varphi).
$$

Note that from Eq. (17) it follows that $\mathsf{b}_{C_1^{\pm}}(u_i^{\pm}(\cdot\,; \omega, k, (\mathcal{S}^{\pm})^i\varphi), v; \omega, k) = 0$ if $v \in H_{\text{per},0}^1(C_1^{\pm})$. The term on the right hand side of Eq. (19) has exactly this form only that $w_j(\cdot\,; (\mathcal{S}^{\pm})^j\psi)$ has a non-vanishing trace at $\Gamma_j^{\pm}$.

### 3.3. Differentiability of the Dirichlet-to-Neumann operators

In the final part of this section, we will show that the DtN operators $\Lambda^\pm(\omega, k)$ are differentiable with respect to the frequency $\omega$ and the quasi-momentum $k$ inside the band gap. Furthermore, we shall explain how to compute the derivatives of the DtN operators via local cell problems. These derivatives will be needed in Section 4 where we derive a formula for the group velocity in which the derivatives of the DtN operators appear. Moreover, we will need the differentiability of the DtN operators with respect to the frequency and the quasi-momentum in Section 6 where we propose a direct procedure to solve the non-linear eigenvalue problem (8) that relies on this property.

Let us again assume that $\omega^2 \notin \sigma_{\text{ess}}^\pm(k)$ and let the problem (3) be well posed. Then the DtN operators are defined uniquely and can be computed using Eq. (16).

Let $u^\pm(\cdot\,; \omega, k, \varphi)$ be the unique solution in $H_{\text{per}}^1(\Delta, S^\pm)$ of the Dirichlet problem (5). Then we introduce $u_\omega^\pm(\cdot\,; \omega, k, \varphi)$ as the unique solution in $H_{\text{per}}^1(\Delta, S^\pm)$ of

$$-(\nabla + \mathrm{i}k\mathbf{a}_1) \cdot (\nabla + \mathrm{i}k\mathbf{a}_1)u_\omega^\pm - \omega^2 \varepsilon u_\omega^\pm = 2\omega\varepsilon u^\pm, \qquad \mathbf{x} \in S^\pm,$$
$$u_\omega^\pm|_{\Gamma_0^\pm} = 0, \tag{20}$$

and $u_k^\pm(\cdot\,; \omega, k, \varphi)$ as the unique solution in $H_{\text{per}}^1(\Delta, S^\pm)$ of

$$-(\nabla + \mathrm{i}k\mathbf{a}_1) \cdot (\nabla + \mathrm{i}k\mathbf{a}_1)u_k^\pm - \omega^2 \varepsilon u_k^\pm = 2\left(-k|\mathbf{a}_1|^2 + \mathrm{i}\mathbf{a}_1 \cdot \nabla\right)u^\pm, \qquad \mathbf{x} \in S^\pm,$$
$$u_k^\pm|_{\Gamma_0^\pm} = 0. \tag{21}$$

The functions $u_\omega^\pm(\cdot\,; \omega, k, \varphi)$ and $u_k^\pm(\cdot\,; \omega, k, \varphi)$ are well-defined for almost any $\omega^2 \notin \sigma_{\text{ess}}^\pm(k)$ thanks to the following proposition.

**Proposition 3.10.** *Let $\omega^2 \notin \sigma_{\text{ess}}^\pm(k)$ and let the problem (5) be well posed. Then the source problems (20) and (21) are well posed.*

The following result is then straightforward.

**Theorem 3.11.** *Suppose that $\omega^2 \notin \sigma_{\text{ess}}^\pm(k)$ and that the problem (5) is well posed in a neighbourhood of $\omega^2$. Then for any $\varphi \in H_{\text{per}}^{1/2}(\Gamma_0^\pm)$, $u(\cdot\,; \omega, k, \varphi)$ is Fréchet-differentiable with respect to $\omega$ and $k$, and*

$$\frac{\partial u^\pm(\mathbf{x}; \omega, k, \varphi)}{\partial \omega} := u_\omega^\pm(\mathbf{x}; \omega, k, \varphi) \qquad and \qquad \frac{\partial u^\pm(\mathbf{x}; \omega, k, \varphi)}{\partial k} := u_k^\pm(\mathbf{x}; \omega, k, \varphi).$$

*Proof.* Suppose that there exists a parameter $h_0$ such that the problem (5) is well posed for any $\omega^2 \in ((\omega - h_0)^2, (\omega + h_0)^2)$. It is easy to see that for all $h \in (0, h_0)$

$$e^\pm(h) = \frac{1}{h}\left(u^\pm(\mathbf{x}; \omega + h, k, \varphi) - u^\pm(\mathbf{x}; \omega, k, \varphi) - hu_\omega^\pm(\mathbf{x}; \omega, k, \varphi)\right)$$

is a solution in $H_{\text{per}}^1(\Delta, S^\pm)$ of

$$-(\nabla + \mathrm{i}k\mathbf{a}_1) \cdot (\nabla + \mathrm{i}k\mathbf{a}_1)e^\pm - \omega^2\varepsilon e^\pm = h^2\varepsilon u^\pm, \qquad \mathbf{x} \in S^\pm,$$
$$e^\pm|_{\Gamma_0^\pm} = 0.$$

Thanks to Proposition 3.10, this problem is well posed and then $\lim_{h\to 0} e^\pm(h) = 0$. The proof for the derivative with respect to $k$ uses exactly the same ideas. $\qquad\square$

Using the definition (6) of the DtN operators $\Lambda^\pm(\omega, k)$, we deduce their Fréchet-differentiability with respect to $\omega$ and $k$.

**Corollary 3.12.** *Suppose that $\omega^2 \notin \sigma_{\text{ess}}^\pm(k)$ and that the problem (5) is well posed in a neighbourhood of $\omega^2$. Then the DtN operators $\Lambda^\pm(\omega, k)$ are differentiable with respect to the frequency $\omega$ and the quasi-momentum $k$, and for all $\varphi \in H_{\text{per}}^{1/2}(\Gamma_0^\pm)$*

$$\frac{\partial \Lambda^\pm}{\partial \omega}(\omega, k)\varphi := \mp\partial_2 u_\omega^\pm(\cdot\,; \omega, k, \varphi)|_{\Gamma_0^\pm} \qquad and \qquad \frac{\partial \Lambda^\pm}{\partial k}(\omega, k)\varphi := \mp\partial_2 u_k^\pm(\cdot\,; \omega, k, \varphi)|_{\Gamma_0^\pm}.$$

**Remark 3.13.** *Iteratively repeating the same steps as above we can deduce that the DtN operators $\Lambda^\pm(\omega, k)$ are differentiable to any order with respect to the frequency $\omega$ and the quasi-momentum $k$ if $\omega^2 \notin \sigma_{\text{ess}}^\pm(k)$ and the problem (5) is well posed in a neighbourhood of $\omega^2$.*

For simplicity of notation, let us write

$$\Lambda_\omega^\pm(\omega,k) = \frac{\partial \Lambda^\pm}{\partial \omega}(\omega,k) \qquad \text{and} \qquad \Lambda_k^\pm(\omega,k) = \frac{\partial \Lambda^\pm}{\partial k}(\omega,k) \tag{22}$$

in the sequel.

Similarly to Eq. (10), we note that the DtN operators differentiated with respect to $\omega$ and $k$ satisfy

$$\int_{\Gamma_0^\pm} \Lambda_\omega^\pm(\omega,k)\varphi\,\overline{\psi}\,\mathrm{d}s(\mathbf{x}) = \int_{S^\pm} (\nabla + \mathrm{i}k\mathbf{a}_1)u_\omega^\pm(\varphi) \cdot (\nabla - \mathrm{i}k\mathbf{a}_1)\overline{u^\pm(\psi)}\,\mathrm{d}\mathbf{x}$$

$$- \omega^2 \int_{S^\pm} \varepsilon u_\omega^\pm(\varphi)\overline{u^\pm(\psi)}\,\mathrm{d}\mathbf{x} - 2\omega \int_{S^\pm} \varepsilon u^\pm(\varphi)\overline{u^\pm(\psi)}\,\mathrm{d}\mathbf{x} \tag{23a}$$

for any $\varphi, \psi \in H_{\mathrm{per}}^{1/2}(\Gamma_0^\pm)$, and — using integration by parts —

$$\int_{\Gamma_0^\pm} \Lambda_k^\pm(\omega,k)\varphi\,\overline{\psi}\,\mathrm{d}s(\mathbf{x}) = \int_{S^\pm} (\nabla + \mathrm{i}k\mathbf{a}_1)u_k^\pm(\varphi) \cdot (\nabla - \mathrm{i}k\mathbf{a}_1)\overline{u^\pm(\psi)}\,\mathrm{d}\mathbf{x}$$

$$- \omega^2 \int_{S^\pm} \varepsilon u_k^\pm(\varphi)\overline{u^\pm(\psi)}\,\mathrm{d}\mathbf{x} + 2k|\mathbf{a}_1|^2 \int_{S^\pm} u^\pm(\varphi)\overline{u^\pm(\psi)}\,\mathrm{d}\mathbf{x}$$

$$- \mathrm{i}|\mathbf{a}_1| \int_{S^\pm} \left(\partial_1 u^\pm(\varphi)\right)\overline{u^\pm(\psi)} - u^\pm(\varphi)\left(\partial_1 \overline{u^\pm(\psi)}\right)\,\mathrm{d}\mathbf{x} \tag{23b}$$

for any $\varphi, \psi \in H_{\mathrm{per}}^{1/2}(\Gamma_0^\pm)$.

Analogously to Corollary 3.12, we obtain the following result for the propagation operators.

**Corollary 3.14.** *Suppose that $\omega^2 \notin \sigma_{\mathrm{ess}}^\pm(k)$ and that the problem (5) is well posed in a neighbourhood of $\omega^2$. Then the propagation operators $\mathcal{P}^\pm(\omega,k)$ are differentiable with respect to the frequency $\omega$ and the quasi-momentum $k$, and for all $\varphi \in H_{\mathrm{per}}^{1/2}(\Gamma_0^\pm)$*

$$\frac{\partial \mathcal{P}^\pm}{\partial \omega}(\omega,k)\varphi := (\mathcal{S}^\pm)^{-1}u_\omega^\pm(\,\cdot\,;\omega,k,\varphi)|_{\Gamma_1^\pm} \qquad and \qquad \frac{\partial \mathcal{P}^\pm}{\partial k}(\omega,k)\varphi := (\mathcal{S}^\pm)^{-1}u_k^\pm(\,\cdot\,;\omega,k,\varphi)|_{\Gamma_1^\pm} \ .$$

Now we want to characterize first the derivatives of the propagation operators and then the derivatives of the DtN operators via solutions of local cell problems. In the following, we explain the characterization of the derivative with respect to $\omega$, the ideas for the derivatives with respect to $k$ are exactly the same.

To this end, we have to introduce the derivatives of the local DtN operators $\mathcal{T}_{ij}^\pm(\omega,k)$. Let us suppose that the Dirichlet cell problems (15) are well defined and let us introduce for all $\varphi \in H_{\mathrm{per}}^{1/2}(\Gamma_0^\pm)$ the unique solutions $u_{\omega,i}^\pm(\,\cdot\,;\omega,k,\varphi)$ in $H_{\mathrm{per}}^1(\Delta,C_1^\pm)$ of the new local cell problems

$$\begin{aligned} -(\nabla + \mathrm{i}k\mathbf{a}_1) \cdot (\nabla + \mathrm{i}k\mathbf{a}_1)u_{\omega,i}^\pm - \omega^2 \varepsilon u_{\omega,i}^\pm &= 2\omega\varepsilon u_i^\pm, \qquad \mathbf{x} \in C_1^\pm, \\ u_{\omega,i}^\pm|_{\Gamma_0^\pm \cup \Gamma_1^\pm} &= 0, \end{aligned} \tag{24}$$

where $u_i^\pm \equiv u_i^\pm(\,\cdot\,;\omega,k,\varphi)$ are the unique solutions of the local cell problems (15). Using exactly the same ideas as above, we can show that the operators $\mathcal{T}_{ij}^\pm(\omega,k)$ are Fréchet differentiable with respect to $\omega$ and for all $\varphi \in H_{\mathrm{per}}^{1/2}(\Gamma_0^\pm)$

$$\frac{\partial \mathcal{T}_{00}^\pm(\omega,k)}{\partial \omega}\varphi = \partial_\mathbf{n} u_{\omega,0}^\pm(\,\cdot\,;\omega,k,\quad \varphi))|_{\Gamma_0^\pm},$$

$$\frac{\partial \mathcal{T}_{01}^\pm(\omega,k)}{\partial \omega}\varphi = (\mathcal{S}^\pm)^{-1}\partial_\mathbf{n} u_{\omega,0}^\pm(\,\cdot\,;\omega,k,\quad \varphi))|_{\Gamma_1^\pm},$$

$$\frac{\partial \mathcal{T}_{10}^\pm(\omega,k)}{\partial \omega}\varphi = \partial_\mathbf{n} u_{\omega,1}^\pm(\,\cdot\,;\omega,k,\mathcal{S}^\pm\varphi))|_{\Gamma_0^\pm},$$

$$\frac{\partial \mathcal{T}_{11}^\pm(\omega,k)}{\partial \omega}\varphi = (\mathcal{S}^\pm)^{-1}\partial_\mathbf{n} u_{\omega,1}^\pm(\,\cdot\,;\omega,k,\mathcal{S}^\pm\varphi))|_{\Gamma_1^\pm} \ .$$

Finally, we can uniquely characterize the derivatives of the propagation operators $\mathcal{P}^\pm(\omega,k)$.

**Proposition 3.15.** *The derivatives of $\mathcal{P}^\pm(\omega,k)$ with respect to $\omega$ are the unique solutions in $\mathcal{L}(H_{\mathrm{per}}^{1/2}(\Gamma_0^\pm))$ of*

$$\left(\mathcal{T}_{10}^\pm(\omega,k)\mathcal{P}^\pm(\omega,k) + \mathcal{T}_{00}^\pm(\omega,k) + \mathcal{T}_{11}^\pm(\omega,k)\right)\frac{\partial \mathcal{P}^\pm(\omega,k)}{\partial \omega} + \mathcal{T}_{10}^\pm(\omega,k)\frac{\partial \mathcal{P}^\pm(\omega,k)}{\partial \omega}\mathcal{P}^\pm(\omega,k)$$

$$= -\frac{\partial \mathcal{T}_{10}^\pm(\omega,k)}{\partial \omega}\left(\mathcal{P}^\pm(\omega,k)\right)^2 - \left(\frac{\partial \mathcal{T}_{00}^\pm(\omega,k)}{\partial \omega} + \frac{\partial \mathcal{T}_{11}^\pm(\omega,k)}{\partial \omega}\right)\mathcal{P}^\pm(\omega,k) - \frac{\partial \mathcal{T}_{01}^\pm(\omega,k)}{\partial \omega}. \tag{25}$$

*Proof.* Differentiating Eq. (14) with respect to $\omega$, it is easy to see that the derivatives of $\mathcal{P}^{\pm}(\omega, k)$ with respect to $\omega$ are solutions of Eq. (25). To deduce uniqueness, it suffices to show that the operator

$$
\begin{aligned}
\mathcal{T}_{\omega,k} : \ \mathcal{L}(H_{\mathrm{per}}^{1/2}(\Gamma_0^{\pm})) \ &\to \ \mathcal{L}(H_{\mathrm{per}}^{1/2}(\Gamma_0^{\pm})) \\
X \ &\mapsto \ \left(\mathcal{T}_{10}^{\pm}(\omega, k)\mathcal{P}^{\pm}(\omega, k) + \mathcal{T}_{00}^{\pm}(\omega, k) + \mathcal{T}_{11}^{\pm}(\omega, k)\right) X + \mathcal{T}_{10}^{\pm}(\omega, k)X\mathcal{P}^{\pm}(\omega, k)
\end{aligned}
$$

is injective. However, injectivity of this operator was already proven in [33], where it occurs in the determination of the DtN operators for time domain problems. Finally, if there exist two solutions $\mathcal{P}_{\omega,1}$ and $\mathcal{P}_{\omega,2}$ of Eq. (25) then their difference satisfies $\mathcal{T}_{\omega,k}\left(\mathcal{P}_{\omega,1} - \mathcal{P}_{\omega,2}\right) = 0$ and by injectivity of $\mathcal{T}_{\omega,k}$, the two solutions are necessarily the same. $\qquad\square$

Differentiating Eq. (16), we can deduce that the derivatives of the DtN operators $\Lambda^{\pm}(\omega, k)$ with respect to $\omega$ read

$$
\Lambda_{\omega}^{\pm}(\omega, k) = \frac{\partial \mathcal{T}_{00}^{\pm}(\omega, k)}{\partial \omega} + \frac{\partial \mathcal{T}_{10}^{\pm}(\omega, k)}{\partial \omega}\mathcal{P}^{\pm}(\omega, k) + \mathcal{T}_{10}^{\pm}(\omega, k)\frac{\partial \mathcal{P}^{\pm}(\omega, k)}{\partial \omega}. \tag{26}
$$

The derivatives of the propagation operators $\mathcal{P}^{\pm}(\omega, k)$ and of the DtN operators $\Lambda^{\pm}(\omega, k)$ with respect to $k$ are characterized similarly by simply replacing all $\omega$-derivatives in Eqs. (25) and (26) by $k$-derivatives. On the other hand, the $k$-derivatives of the operators $\mathcal{T}_{ij}^{\pm}(\omega, k)$, $i, j = 0, 1$, are for all $\varphi \in H_{\mathrm{per}}^{1/2}(\Gamma_0^{\pm})$ given by

$$
\begin{aligned}
\frac{\partial \mathcal{T}_{00}^{\pm}(\omega, k)}{\partial k}\varphi &= \partial_{\mathbf{n}} u_{k,0}^{\pm}(\cdot\,; \omega, k, \quad \varphi))|_{\Gamma_0^{\pm}}, \\
\frac{\partial \mathcal{T}_{01}^{\pm}(\omega, k)}{\partial k}\varphi &= (\mathcal{S}^{\pm})^{-1}\partial_{\mathbf{n}} u_{k,0}^{\pm}(\cdot\,; \omega, k, \quad \varphi))|_{\Gamma_1^{\pm}}, \\
\frac{\partial \mathcal{T}_{10}^{\pm}(\omega, k)}{\partial k}\varphi &= \partial_{\mathbf{n}} u_{k,1}^{\pm}(\cdot\,; \omega, k, \mathcal{S}^{\pm}\varphi))|_{\Gamma_0^{\pm}}, \\
\frac{\partial \mathcal{T}_{11}^{\pm}(\omega, k)}{\partial k}\varphi &= (\mathcal{S}^{\pm})^{-1}\partial_{\mathbf{n}} u_{k,1}^{\pm}(\cdot\,; \omega, k, \mathcal{S}^{\pm}\varphi))|_{\Gamma_1^{\pm}},
\end{aligned}
$$

where $u_{k,i}^{\pm}(\mathbf{x}; \omega, k, \varphi)$ are the unique solutions in $H_{\mathrm{per}}^1(\Delta, C_1^{\pm})$, $i = 0, 1$, of

$$
\begin{aligned}
-(\nabla + \mathrm{i}k\mathbf{a}_1) \cdot (\nabla + \mathrm{i}k\mathbf{a}_1)u_{k,i}^{\pm} - \omega^2 \varepsilon u_{k,i}^{\pm} &= 2\left(k|\mathbf{a}_1| + \mathrm{i}\mathbf{a}_1 \cdot \nabla\right)u_i^{\pm}, \qquad \mathbf{x} \in C_1^{\pm}, \\
u_{k,i}^{\pm}|_{\Gamma_0^{\pm} \cup \Gamma_1^{\pm}} &= 0.
\end{aligned} \tag{27}
$$

**Remark 3.16.** *In contrast to the Dirichlet cell problems (15) to determine the DtN operators, the Dirichlet cell problems (24) and (27) to compute the $\omega$- and $k$-derivatives of the DtN operators have homogeneous Dirichlet boundary conditions but a source term that depends on the solutions $u_i^{\pm}$, $i = 0, 1$, of the original cell problems (15).*

## 4. Group velocity of guided modes

The group velocity $\frac{\partial \omega}{\partial k}$ of guided modes [34] is a quantity of particular interest when tailoring PhC wave-guides. In general, it is desired to construct so called *slow light wave-guides*, which are wave-guides whose band structure contains a frequency interval with a dispersive curve (of the second kind) with a small magnitude of its gradient, *i. e.* with a small magnitude of its group velocity, see for example [15, 35].

In this section we derive a formula for the group velocity $\frac{\partial \omega}{\partial k}$ of guided modes which can then be used in the process of optimizing the band structure of PhC wave-guides. Another application of this formula is described in Section 7, where we use the group velocity in order to compute the start values of the iterative Newton method to solve the non-linear eigenvalue problem (8).

We start with a symmetry property of the non-linear eigenvalue problem (8), and in the second part of this section we finally derive the formula for the group velocity.

*4.1. Symmetry of the non-linear eigenvalue problem in the Brillouin zone*

Let us consider the non-linear eigenvalue problem (8) with the sesquilinear forms $\mathsf{b}_{C_0}$ and $\mathsf{d}_{C_0}$. First we show

**Lemma 4.1.** *Let $k \in B$ and $\omega^2 \in \mathbb{R}^+ \setminus \sigma_{\mathrm{ess}}(k)$. Furthermore, let (5) be well-posed. Then*

$$
\mathsf{d}_{C_0}(u, v; \omega, k) = \mathsf{d}_{C_0}(\overline{v}, \overline{u}; \omega, -k). \tag{28}
$$

*Proof.* This is a direct consequence of the definition (9c) of the sesquilinear form $\mathsf{d}_{C_0}$ and the definition (10) of the DtN operators $\Lambda^{\pm}$. $\qquad\square$

**Lemma 4.2.** *Let $k \in B$ and $\omega^2 \in \mathbb{R}^+ \setminus \sigma_{\mathrm{ess}}(k)$. Furthermore, let (5) be well-posed. Then*

$$\overline{\mathsf{d}_{C_0}(u, v; \omega, k)} = \mathsf{d}_{C_0}(\overline{u}, \overline{v}; \omega, -k). \tag{29}$$

*Proof.* Using Lemma 4.1 and the fact that

$$\overline{\mathsf{d}_{C_0}(u, v; \omega, k)} = \mathsf{d}_{C_0}(v, u; \omega, \overline{k}) = \mathsf{d}_{C_0}(v, u; \omega, k), \tag{30}$$

which follows from the definition (9c) of $\mathsf{d}_{C_0}$, the definition (10) of the DtN operators $\Lambda^{\pm}$, and the fact that $\overline{k} = k$ if $k \in B \subset \mathbb{R}$, we can directly conclude Eq. (29). $\square$

**Lemma 4.3.** *For any $(\omega^2, k) \in \mathbb{R}^+ \times \mathbb{C}$*

$$\overline{\mathsf{b}_{C_0}(u, v; \omega, k)} = \mathsf{b}_{C_0}(\overline{u}, \overline{v}; \omega, -\overline{k}).$$

*Proof.* This follows directly from the definition (9d) of the sesquilinear form $\mathsf{b}_{C_0}$. $\square$

Now we are ready to show

**Proposition 4.4.** *Let $(\omega^2, k) \in \mathbb{R}^+ \times B$ with $\omega^2 \notin \sigma_{\mathrm{ess}}(k)$ be an eigenvalue couple of the non-linear eigenvalue problem (8) with associated eigenfunction $u \in H^1_{\mathrm{per}}(C_0)$. Then $(\omega^2, -k) \in \mathbb{R}^+ \times B$ is an eigenvalue couple of (8) with associated eigenfunction $\overline{u} \in H^1_{\mathrm{per}}(C_0)$.*

*Proof.* If $(\omega^2, k) \in \mathbb{R}^+ \times B$ with $\omega^2 \notin \sigma_{\mathrm{ess}}(k)$ is an eigenvalue couple of (8) with associated eigenfunction $u \in H^1_{\mathrm{per}}(C_0)$, then

$$\overline{\mathsf{b}_{C_0}(u, v; \omega, k)} + \overline{\mathsf{d}_{C_0}(u, v; \omega, k)} = \overline{\mathsf{b}_{C_0}(u, v; \omega, k) + \mathsf{d}_{C_0}(u, v; \omega, k)} = 0$$

for all $v \in H^1_{\mathrm{per}}(C_0)$. Using Lemmas 4.2 and 4.3 as well as the fact that $\overline{k} = k$ if $k \in B \subset \mathbb{R}$, we obtain

$$\mathsf{b}_{C_0}(\overline{u}, \overline{v}; \omega, -k) + \mathsf{d}_{C_0}(\overline{u}, \overline{v}; \omega, -k) = 0$$

for all $v \in H^1_{\mathrm{per}}(C_0)$, from which the result directly follows. $\square$

*4.2. Computation of the group velocity*

Now we come to the main theorem of this section and derive a formula for the group velocity of guided modes.

**Theorem 4.5.** *Let $k \in B$ and suppose $\omega^2(k) \in \mathbb{R}^+ \setminus \sigma_{\mathrm{ess}}(k)$ is an eigenvalue of the non-linear eigenvalue problem (8) with associated eigenmode $u \equiv u(\cdot\,; \omega(k), k) \in H^1_{\mathrm{per}}(C_0)$. Then the group velocity is real-valued and reads*

$$\frac{\partial \omega}{\partial k}(k) = \frac{2k|\mathbf{a}_1|^2 \int_{C_0} |u|^2 \, \mathrm{d}\mathbf{x} - 2|\mathbf{a}_1| \operatorname{Im}\left(\int_{C_0} u \, \partial_1 \overline{u} \, \mathrm{d}\mathbf{x}\right) + \int_{\Gamma_0^+} \Lambda_k^+ u \, \overline{u} \, \mathrm{d}s(\mathbf{x}) + \int_{\Gamma_0^-} \Lambda_k^- u \, \overline{u} \, \mathrm{d}s(\mathbf{x})}{2\omega \int_{C_0} \varepsilon |u|^2 \, \mathrm{d}\mathbf{x} - \int_{\Gamma_0^+} \Lambda_\omega^+ u \, \overline{u} \, \mathrm{d}s(\mathbf{x}) - \int_{\Gamma_0^-} \Lambda_\omega^- u \, \overline{u} \, \mathrm{d}s(\mathbf{x})}, \tag{31}$$

*with $\Lambda_\omega^{\pm}$ and $\Lambda_k^{\pm}$ as defined in (22). Alternatively, to Eq. (31) the group velocity can be written in the form*

$$\frac{\partial \omega}{\partial k}(k) = \frac{k|\mathbf{a}_1|^2 \int_S |u^S|^2 \, \mathrm{d}\mathbf{x} - |\mathbf{a}_1| \operatorname{Im}\left(\int_S u^S \, \partial_1 \overline{u^S} \, \mathrm{d}\mathbf{x}\right)}{\omega \int_S \varepsilon |u^S|^2 \, \mathrm{d}\mathbf{x}}, \tag{32}$$

*where $u^S \in H^1_{\mathrm{per}}(S)$ with $u^S|_{C_0} = u$ and $u^S|_{S^{\pm}} = u^{\pm}(\cdot\,; \omega(k), k, u|_{\Gamma_0^{\pm}})$.*

**Remark 4.6.** *For the computation of the group velocity only Eq. (31) is relevant, since the integrals in Eq. (32) are posed on the unbounded domain $S$. However, from Eq. (32) we can see that the denominator of the group velocity formula is strictly positive. Moreover, we can easily deduce the group velocity formula for the super-cell method, by simply replacing $S$ with the super-cell domain $S_{\mathrm{sc}} \subset S$.*

*Proof of Theorem 4.5.* If we assume that $(\omega^2(k), k) \in \mathbb{R}^+ \times B$ with $\omega^2 \notin \sigma_{\mathrm{ess}}(k)$ is an eigenvalue couple of the non-linear eigenvalue problem (8) with associated eigenmode $u(\cdot\,; k) \in H^1_{\mathrm{per}}(C_0)$, then we know from Proposition 2.3 that the group velocity $\omega_k = \frac{\partial \omega}{\partial k}$ exists at $(\omega^2(k), k)$ and is continuous. Moreover, we showed in Section 3.3 that the DtN operators $\Lambda^{\pm}(\omega, k)$ are differentiable. We can hence define

$u_k(\,\cdot\,;k) := \frac{\mathrm{d}u}{\mathrm{d}k}(\,\cdot\,;k) \in H^1_{\mathrm{per}}(C_0)$, and differentiate the non-linear eigenvalue problem (8) with respect to $k$, which yields the following equality

$$\int_{C_0} (\nabla + ik\mathbf{a}_1)u_k \cdot (\nabla - ik\mathbf{a}_1)\overline{v} - \omega^2 \varepsilon u_k \overline{v} \,\mathrm{d}\mathbf{x} + \int_{\Gamma_0^+} \Lambda^+(\omega,k)u_k\,\overline{v}\,\mathrm{d}s(\mathbf{x}) + \int_{\Gamma_0^-} \Lambda^-(\omega,k)u_k\,\overline{v}\,\mathrm{d}s(\mathbf{x})$$

$$= \int_{C_0} (\nabla + ik\mathbf{a}_1)u \cdot i\mathbf{a}_1\overline{v} - i\mathbf{a}_1 u \cdot (\nabla - ik\mathbf{a}_1)\overline{v}\,\mathrm{d}\mathbf{x} + 2\omega\omega_k \int_{C_0} \varepsilon u \overline{v}\,\mathrm{d}\mathbf{x}$$

$$- \int_{\Gamma_0^+} \Lambda_k^+(\omega,k)u\,\overline{v}\,\mathrm{d}s(\mathbf{x}) - \int_{\Gamma_0^-} \Lambda_k^-(\omega,k)u\,\overline{v}\,\mathrm{d}s(\mathbf{x})$$

$$- \omega_k \int_{\Gamma_0^+} \Lambda_\omega^+(\omega,k)u\,\overline{v}\,\mathrm{d}s(\mathbf{x}) - \omega_k \int_{\Gamma_0^-} \Lambda_\omega^-(\omega,k)u\,\overline{v}\,\mathrm{d}s(\mathbf{x}) \quad (33)$$

for all $v \in H^1_{\mathrm{per}}(C_0)$. It is easy to see that if we choose $v = u$ as test function and consider Lemma 4.1 and Proposition 4.4, the left hand side of Eq. (33) vanishes, and hence, all terms containing $u_k$ disappear and we get

$$2\omega\omega_k \int_{C_0} \varepsilon|u|^2 \,\mathrm{d}\mathbf{x} - \omega_k \int_{\Gamma_0^+} \Lambda_\omega^+(\omega,k)u\,\overline{u}\,\mathrm{d}s(\mathbf{x}) - \omega_k \int_{\Gamma_0^-} \Lambda_\omega^-(\omega,k)u\,\overline{u}\,\mathrm{d}s(\mathbf{x})$$

$$= -\int_{C_0} (\nabla + ik\mathbf{a}_1)u \cdot i\mathbf{a}_1\overline{u} - i\mathbf{a}_1 u \cdot (\nabla - ik\mathbf{a}_1)\overline{u}\,\mathrm{d}\mathbf{x}$$

$$+ \int_{\Gamma_0^+} \Lambda_k^+(\omega,k)u\,\overline{u}\,\mathrm{d}s(\mathbf{x}) + \int_{\Gamma_0^-} \Lambda_k^-(\omega,k)u\,\overline{u}\,\mathrm{d}s(\mathbf{x}). \quad (34)$$

Considering that $\mathbf{a}_1 = |\mathbf{a}_1|\,(1,0)^{\mathrm{T}}$, we can deduce that the integral over $C_0$ on the right hand side of Eq. (34) can be written in the form

$$-\int_{C_0} (\nabla + ik\mathbf{a}_1)u \cdot i\mathbf{a}_1\overline{u} - i\mathbf{a}_1 u \cdot (\nabla - ik\mathbf{a}_1)\overline{u}\,\mathrm{d}\mathbf{x} = i|\mathbf{a}_1| \int_{C_0} u\,(\partial_1\overline{u}) - (\partial_1 u)\,\overline{u}\,\mathrm{d}\mathbf{x} + 2k|\mathbf{a}_1|^2 \int_{C_0} |u|^2 \,\mathrm{d}\mathbf{x}. \quad (35)$$

The first integral of the right hand side of (35) is purely imaginary since $u\,(\partial_1\overline{u})$ is the complex conjugate of $(\partial_1 u)\,\overline{u}$. Using integration by parts and the fact that $u$ is periodic in $x_1$-direction, we can rewrite this purely imaginary integral in the form

$$\int_{C_0} u\,\partial_1\overline{u} - \partial_1 u\,\overline{u}\,\mathrm{d}\mathbf{x} = 2i\,\mathrm{Im}\left(\int_{C_0} u\,\partial_1\overline{u}\,\mathrm{d}\mathbf{x}\right).$$

Thus, we end up with

$$\omega_k \left(2\omega \int_{C_0} \varepsilon|u|^2 \,\mathrm{d}\mathbf{x} - \int_{\Gamma_0^+} \Lambda_\omega^+(\omega,k)u\,\overline{u}\,\mathrm{d}s(\mathbf{x}) - \int_{\Gamma_0^-} \Lambda_\omega^-(\omega,k)u\,\overline{u}\,\mathrm{d}s(\mathbf{x})\right)$$

$$= 2k|\mathbf{a}_1|^2 \int_{C_0} |u|^2 \,\mathrm{d}\mathbf{x} - 2|\mathbf{a}_1|\,\mathrm{Im}\left(\int_{C_0} u\,\partial_1\overline{u}\,\mathrm{d}\mathbf{x}\right) + \int_{\Gamma_0^+} \Lambda_k^+(\omega,k)u\,\overline{u}\,\mathrm{d}s(\mathbf{x}) + \int_{\Gamma_0^-} \Lambda_k^-(\omega,k)u\,\overline{u}\,\mathrm{d}s(\mathbf{x})$$

from which Eq. (31) directly follows.

Now we show that $\omega_k(k)$ is real-valued. From (23) it follows that

$$\int_{\Gamma_0^\pm} \Lambda_\omega^\pm u\,\overline{u}\,\mathrm{d}s(\mathbf{x}) = \int_{S^\pm} (\nabla + ik\mathbf{a}_1)u_\omega^\pm \cdot (\nabla - ik\mathbf{a}_1)\overline{u^\mp} - \omega^2 \varepsilon u_\omega^\pm \overline{u^\mp} \,\mathrm{d}\mathbf{x} - 2\omega \int_{S^\pm} \varepsilon|u^\pm|^2 \,\mathrm{d}\mathbf{x}.$$

Indeed, the first integral on the right hand side vanishes. Using integration by parts, considering that $u^\pm$ solves (5) and taking into account that the solution $u_\omega^\pm$ of the half-strip problem (20) vanishes on $\Gamma_0^\pm$ we deduce

$$\int_{S^\pm} (\nabla + ik\mathbf{a}_1)u_\omega^\pm \cdot (\nabla - ik\mathbf{a}_1)\overline{u^\mp} - \omega^2 \varepsilon u_\omega^\pm \overline{u^\mp} \,\mathrm{d}\mathbf{x}$$

$$= \overline{\int_{S^\pm} (-(\nabla + ik\mathbf{a}_1) \cdot (\nabla + ik\mathbf{a}_1)u^\pm - \omega^2 \varepsilon u^\pm)\,\overline{u_\omega^\mp}\,\mathrm{d}\mathbf{x} + \int_{\Gamma_0^\pm} (\mp\partial_2 u^\pm)\,\overline{u_\omega^\mp}\,\mathrm{d}s(\mathbf{x})} = 0$$

and hence,

$$\int_{\Gamma_0^\pm} \Lambda_\omega^\pm u\,\overline{u}\,\mathrm{d}s(\mathbf{x}) = -2\omega \int_{S^\pm} \varepsilon|u^\pm|^2 \,\mathrm{d}\mathbf{x} \in \mathbb{R}. \quad (36)$$

The same considerations apply to $\int_{\Gamma_0^\pm} \Lambda_k^\pm u\,\overline{u}\,\mathrm{d}s(\mathbf{x})$ and hence, the group velocity (31) is real valued.

From Eq. (36) and its analogon for $\int_{\Gamma_0^\pm} \Lambda_k^\pm u\,\overline{u}\,\mathrm{d}s(\mathbf{x})$ it is easy to see, that the formula (31) of group velocity is equivalent to the formula in Eq. (32). $\qquad\square$

## 5. Discretization

In Section 3 we introduced the non-linear eigenvalue problem for the exact computation of guided modes in weak form, *c.f.* Eq. (8). We also introduced a variational formulation (19) for the local cell problems to compute the DtN operators (16). In this section we now want to introduce a finite element discretization of the spaces involved, and describe the computation of the discrete DtN map.

### 5.1. High-order finite element method on curved cells

The finite element method (FEM) provides discrete subspaces of the Sobolev spaces $H^1_{\mathrm{per}}(C_0)$, $H^1_{\mathrm{per}}(C_1^{\pm})$ and $H^{1/2}_{\mathrm{per}}(\Gamma_0^{\pm})$ which we are going to describe in the following.

*The meshes.* For this, each of the rhomboid domains $C_1^{\pm}$ of the cell problems (17) and $C_0$ of the non-linear eigenvalue problem (8) are partitioned into non-overlapping possibly curved, triangular or quadrilateral subdomains — the geometrical cells (see for example the very coarse meshes in Figure 3). Each geometrical cell $K$ is defined as a smooth map $F_K$ of the reference cell $\widehat{K}(K)$, which is for triangular cells the convex hull of the points $(0,0)$, $(1,0)$ and $(0,1)$ and for quadrilaterals the square $[0,1]^2$. The sets of the geometrical cells — the meshes — are denoted by $\mathcal{M}(C_0)$ and $\mathcal{M}(C_1^{\pm})$, respectively. All these meshes are assumed to be periodic in direction $\mathbf{a}_1$, *i.e.* for each edge of a geometrical cell on $\Sigma_{\mathrm{L}}$ there is an edge on $\Sigma_{\mathrm{R}}$, which is only shifted by $\mathbf{a}_1$. In particular this means that the corresponding geometrical cells need to have the same parameterization on the boundary. We call two geometrical cells $K_1$ and $K_2$ of a mesh $\mathcal{M}$ neighbouring, if they share an edge or, in order to construct periodic basis functions, if an edge of $K_1$ coincides with an edge of $K_2$ shifted by $\mathbf{a}_1$ or $-\mathbf{a}_1$. The set of edges of $\mathcal{M}(C_0)$ lying on $\Gamma_0^{\pm}$ are the geometrical cells of the *interface mesh* $\mathcal{M}(\Gamma_0^{\pm})$. The geometrical cells $K$ in $\mathcal{M}(\Gamma_0^{\pm})$ can alternatively be defined by affine maps $F_K$ from the reference interval $\widehat{K} = [0,1]$. Two geometrical cells $K_1$ and $K_2$ in $\mathcal{M}(\Gamma_0^{\pm})$ are called neighbouring, if they share a vertex, or, to simply construct periodic basis function, if a vertex of $K_1$ coincides with a vertex of $K_2$ shifted by $\mathbf{a}_1$ or $-\mathbf{a}_1$. We assume matching meshes $\mathcal{M}(C_0)$ and $\mathcal{M}(C_1^{\pm})$, meaning that the edges on $\Gamma_0^{\pm}$ coincide. For a simple construction of the shift operators $\mathcal{S}^{\pm}$ we assume furthermore the meshes $\mathcal{M}(C_1^{\pm})$ to be periodic in direction $\mathbf{a}_2^{\pm}$.

*The FE spaces.* We define discrete subspaces of $H^1_{\mathrm{per}}(C_0)$, $H^1_{\mathrm{per}}(C_1^{\pm})$ and $H^{1/2}_{\mathrm{per}}(\Gamma_0^{\pm})$ as

$$S^{p,1}_{\mathrm{per}}(O, \mathcal{M}(O)) := \{v \in H^1_{\mathrm{per}}(O) \cap C^0(\overline{O}) : v|_K \circ F_K \in \mathcal{Q}_p(\widehat{K}(K)) \ \forall K \in \mathcal{M}(O)\},$$

where $p$ is a chosen polynomial degree, $O$ stands for one of the computational domains $C_0$, $C_1^{\pm}$ or $\Gamma_0^{\pm}$, and $\mathcal{Q}_p(\widehat{K})$ is the space of polynomials with maximal (or maximal total) degree $p$

$$\mathcal{Q}_p(\widehat{K}) = \begin{cases} \mathrm{span}\{\widehat{x}_1^k \widehat{x}_2^\ell, 0 \leq \max(k,\ell) \leq p\}, & \text{if } \widehat{K} \text{ is a quadrilateral}, \\ \mathrm{span}\{\widehat{x}_1^k \widehat{x}_2^\ell, 0 \leq k + \ell \leq p\}, & \text{if } \widehat{K} \text{ is a triangle}, \\ \mathrm{span}\{\widehat{x}^k, \quad 0 \leq k \leq p\}, & \text{if } \widehat{K} \text{ is an interval}. \end{cases}$$

Let us furthermore call $S^{p,1}_{\mathrm{per},0}(O, \mathcal{M}(O))$ the respective subspaces of $S^{p,1}_{\mathrm{per}}(O, \mathcal{M}(O))$ for $O = C_0, C_1^{\pm}$ with vanishing Dirichlet trace on $\Gamma_0^{\pm}$ and $\Gamma_1^{\pm}$, and $S^{p,1}_{\mathrm{per},0}(\Gamma_0^{\pm}, \mathcal{M}(\Gamma_0^{\pm})) = S^{p,1}_{\mathrm{per}}(\Gamma_0^{\pm}, \mathcal{M}(\Gamma_0^{\pm}))$.

If the maximal polynomial degree is $p = 1$ the basis functions are *hat functions* that take the value 1 at a single node of the mesh $\mathcal{M}$ and 0 at all other nodes, and we speak of *linear FEM*. For polynomial degrees larger than 1 the FEM is said to be of *high order* [36]. Besides the nodal hat functions, the basis of a high-order FEM consists of functions that can be identified to an edge and that vanish on (or more precisely in the closure of) all other edges in the mesh $\mathcal{M}$ — *the edge functions* —, and — for triangular and quadrilateral meshes — so-called *bubble functions* that can be identified to one cell and that are zero in (the closure of) all other cells of $\mathcal{M}$.

*Refinement strategies.* With the FE method the solution of the cell problems and the eigenmodes, solutions to the non-linear eigenvalue problems (8), are approximated. There are basically three strategies to improve the accuracy of an existing FE approximation:

- refine the mesh of the computational domain (*h*-FEM),

- increase the polynomial degree of the basis functions (*p*-FEM), or

- a combination of both (*hp*-FEM).

While $h$-FEM provides algebraic convergence, $p$-FEM converges exponentially if the solution is analytic in subdomains that are resolved exactly by the cells of the mesh [36]. This implies the need of curved cells in the mesh in order to exactly resolve computational domains that contain curved boundaries. For example the computational domain of PhC wave-guides as sketched in Figure 1 contains circular holes. Therefore, we use a mesh with circular cell boundaries, in order to perfectly resolve the computational domain (*i. e.* the cell $C_0$ of the guide and the single periodicity cells $C_1^{\pm}$ on top and bottom of $C_0$) with a coarse mesh, see Figure 3. Note, that in case of non-smooth material boundaries we can extend the previous and following definitions to $hp$-adaptive FE spaces.

### 5.2. Solution of the cell problems

From now on we assume that we have the same maximal polynomial degree $p$ for all previously defined spaces. We will denote the dimensions of the spaces $S_{\mathrm{per},0}^{p,1}(O, \mathcal{M}(O))$ by

$$N(O) := \dim S_{\mathrm{per},0}^{p,1}(O, \mathcal{M}(O)), \qquad O = C_0, C_1^{\pm}, \Gamma_0^{\pm}$$

and we conclude

$$\dim S_{\mathrm{per}}^{p,1}(C_0, \mathcal{M}(C_0)) = N(C_0) + N(\Gamma_0^+) + N(\Gamma_0^-),$$
$$\dim S_{\mathrm{per}}^{p,1}(C_1^{\pm}, \mathcal{M}(C_1^{\pm})) = N(C_1^{\pm}) + 2\, N(\Gamma_0^{\pm}).$$

To obtain discrete approximations of the operators defined in Section 3 and the algebraic systems for the variational formulations in $C_0$ and $C_1^{\pm}$ we introduce ordered lists of basis functions for the spaces described above. The Dirichlet traces $\varphi^{\pm} \in H_{\mathrm{per}}^{1/2}(\Gamma_0^{\pm})$ are approximated by

$$\varphi_h^{\pm}(\mathbf{x}) = \sum_{n=0}^{N(\Gamma_0^{\pm})-1} \varphi_n^{\pm} b_{\Gamma_0^{\pm},n}(\mathbf{x}),$$

where $b_{\Gamma_0^{\pm},n}$ is the $n$-th basis function of $S_{\mathrm{per}}^{p,1}(\Gamma_0^{\pm}, \mathcal{M}(\Gamma_0^{\pm}))$. The propagation operators on the discrete spaces are represented by the matrices $\mathbf{P}^{\pm}(\omega, k) \in \mathbb{C}^{N(\Gamma_0^{\pm}) \times N(\Gamma_0^{\pm})}$, *i. e.*

$$\mathcal{P}^{\pm}(\omega, k)\varphi_h^{\pm}(\mathbf{x}) = \sum_{n=0}^{N(\Gamma_0^{\pm})-1} \varphi_n^{\pm} \sum_{m=0}^{N(\Gamma_0^{\pm})-1} P_{mn}^{\pm}(\omega, k) b_{\Gamma_0^{\pm},m}(\mathbf{x}). \tag{37}$$

To obtain the propagation operator we need approximate solutions $u_{i,h}^{\pm}(\cdot\,; \omega, k, \varphi_h) \in S_{\mathrm{per}}^{p,1}(C_1^{\pm}, \mathcal{M}(C_1^{\pm}))$ to the cell problems (17). For this, we introduce basis functions $b_{C_1^{\pm},n}$ of $S_{\mathrm{per}}^{p,1}(C_1^{\pm}, \mathcal{M}(C_1^{\pm}))$, which are ordered such that the

- basis functions with index $n \in \mathcal{I}(C_1^{\pm}) := \{0, \dots, N(C_1^{\pm}) - 1\}$ vanish on $\Gamma_0^{\pm}$ and $\Gamma_1^{\pm}$,

- those with index $n \in \mathcal{I}(C_1^{\pm}, \Gamma_0^{\pm}) := \{N(C_1^{\pm}), \dots, N(C_1^{\pm}) + N(\Gamma_0^{\pm}) - 1\}$ vanish on $\Gamma_1^{\pm}$, but their traces on $\Gamma_0^{\pm}$ build a basis of $S_{\mathrm{per}}^{p,1}(\Gamma_0^{\pm}, \mathcal{M}(\Gamma_0^{\pm}))$, and

- those with index $n \in \mathcal{I}(C_1^{\pm}, \Gamma_1^{\pm}) := \{N(C_1^{\pm}) + N(\Gamma_0^{\pm}), \dots, N(C_1^{\pm}) + 2\, N(\Gamma_0^{\pm}) - 1\}$ vanish on $\Gamma_0^{\pm}$, but their traces on $\Gamma_1^{\pm}$ shifted to $\Gamma_0^{\pm}$ build a basis of $S_{\mathrm{per}}^{p,1}(\Gamma_0^{\pm}, \mathcal{M}(\Gamma_0^{\pm}))$ as well.

With this special ordering and the relation between the traces of the basis functions in $S_{\mathrm{per}}^{p,1}(C_1^{\pm}, \mathcal{M}(C_1^{\pm}))$ and $S_{\mathrm{per}}^{p,1}(\Gamma_0^{\pm}, \mathcal{M}(\Gamma_0^{\pm}))$

$$b_{\Gamma_0^{\pm},n} = \sum_{m=0}^{N(\Gamma_0^{\pm})} Q_{0,mn}^{\pm} b_{C_1^{\pm},m+N(C_1^{\pm})}|_{\Gamma_0^{\pm}} = \sum_{m=0}^{N(\Gamma_0^{\pm})} Q_{1,mn}^{\pm} b_{C_1^{\pm},m+N(C_1^{\pm})+N(\Gamma_0^{\pm})}|_{\Gamma_1^{\pm}} \tag{38}$$

with matrices $\mathbf{Q}_i^{\pm} \in \mathbb{R}^{N(\Gamma_0^{\pm}) \times N(\Gamma_0^{\pm})}$, $i = 0, 1$, we can define Dirichlet lifts for the basis functions $b_{\Gamma_0^{\pm},n} \in S_{\mathrm{per}}^{p,1}(\Gamma_0^{\pm}, \mathcal{M}(\Gamma_0^{\pm}))$ as

$$w_0^{\pm}(\cdot\,; b_{\Gamma_0^{\pm},n}) = \sum_{m=0}^{N(\Gamma_0^{\pm})-1} Q_{0,mn}^{\pm} b_{C_1^{\pm},m+N(C_1^{\pm})}, \qquad w_1^{\pm}(\cdot\,; b_{\Gamma_0^{\pm},n}) = \sum_{m=0}^{N(\Gamma_0^{\pm})-1} Q_{1,mn}^{\pm} b_{C_1^{\pm},m+N(C_1^{\pm})+N(\Gamma_0^{\pm})}.$$

We use an hierarchical family of shape functions proposed by Karniadakis and Sherwin [37] for all spaces $S_{\mathrm{per}}^{p,1}(O, \mathcal{M}(O))$ for which the matrices $\mathbf{Q}_i^{\pm}$, $i = 0, 1$, have the structure of permutation matrices with

entries $\pm 1$. Then, there are only entries $-1$ if the corresponding edge functions are of odd order and the global and local orientations of the edge, which are responsible for the direction, mismatch (see Figure 4).

Now, we can write the cell problems for $u_{i,h}^\pm$ as linear systems of equations

$$\mathbf{B}_{C_1^\pm}(C_1^\pm, C_1^\pm; \omega, k)\mathbf{u}_{i,0,h}^\pm(\omega, k, b_{\Gamma_0^\pm, n}) = -\mathbf{B}_{C_1^\pm}(C_1^\pm, \Gamma_i^\pm; \omega, k)\mathbf{Q}_i^\pm \mathbf{e}_n, \quad n = 0, \ldots, N(\Gamma_0^\pm) - 1,$$

where $\mathbf{u}_{i,0,h}^\pm(\omega, k, b_{\Gamma_0^\pm, n}) \in \mathbb{C}^{N(C_1^\pm)}$ is the coefficient vector of $u_{i,0,h}^\pm(\cdot\,; \omega, k, b_{\Gamma_0^\pm, n})$ with respect to the basis functions $b_{C_1^\pm}$, $\mathbf{e}_n$ is the $n$-th unit vector, and $\mathbf{B}_{C_1^\pm}(O_1, O_2, \omega, k)$ is the block with row indices $\mathcal{I}(O_1)$ and column indices $\mathcal{I}(O_2)$ of the system matrix $\mathbf{B}_{C_1^\pm}(\omega, k)$ related to the sesquilinear form (18) for the spaces $S_{\mathrm{per}}^{p,1}(C_1^\pm, \mathcal{M}(C_1^\pm))$.

We can collect the coefficient vectors $\mathbf{u}_{i,0,h}^\pm(\omega, k, b_{\Gamma_0^\pm, n})$ for $n = 0, \ldots, N(\Gamma_0^\pm) - 1$ in (rectangular) matrices $\mathbf{U}_{i,0,h}^\pm \in \mathbb{C}^{N(C_1^\pm) \times N(\Gamma_0^\pm)}$ which can be defined by

$$\mathbf{B}_{C_1^\pm}(C_1^\pm, C_1^\pm; \omega, k)\mathbf{U}_{i,0,h}^\pm(\cdot\,; \omega, k) = -\mathbf{B}_{C_1^\pm}(C_1^\pm, \Gamma_i^\pm; \omega, k)\mathbf{Q}_i^\pm,$$

Inserting the basis functions $b_{\Gamma_0^\pm, n}$ of $S_{\mathrm{per}}^{p,1}(\Gamma_0^\pm, \mathcal{M}(\Gamma_0^\pm))$ into (19) defines the matrices $\mathbf{T}_{ij}^\pm$, $i, j = 0, 1$ as

$$\begin{aligned}
\mathbf{T}_{ij}^\pm(\omega, k) &= (\mathbf{Q}_j^\pm)^{\mathrm{T}}\mathbf{B}_{C_1^\pm}(\Gamma_j^\pm, C_1^\pm)\mathbf{U}_{i,0,h}^\pm + (\mathbf{Q}_j^\pm)^{\mathrm{T}}\mathbf{B}_{C_1^\pm}(\Gamma_j^\pm, \Gamma_i^\pm)\mathbf{Q}_i^\pm \\
&= (\mathbf{Q}_j^\pm)^{\mathrm{T}}\left(\mathbf{B}_{C_1^\pm}(\Gamma_j^\pm, \Gamma_i^\pm) - \mathbf{B}_{C_1^\pm}(\Gamma_j^\pm, C_1^\pm)\mathbf{B}_{C_1^\pm}(C_1^\pm, C_1^\pm)^{-1}\mathbf{B}_{C_1^\pm}(C_1^\pm, \Gamma_i^\pm)\right)\mathbf{Q}_i^\pm,
\end{aligned} \tag{39}$$

which are equal to Schur complement matrices applied to "permutation" matrices $\mathbf{Q}_i^\pm$. Here we omitted the $(\omega, k)$-dependence in the matrices $\mathbf{B}_{C_1^\pm}(\cdot, \cdot\,; \omega, k)$.
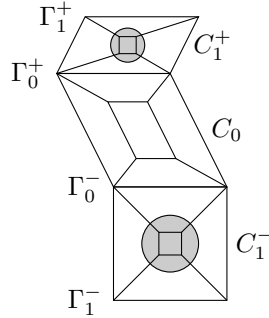


Figure 3: Sketch of the mesh of the computational domain with the cell $C_0$ of the guide, the single periodicity cells $C_1^\pm$ on top and bottom of $C_0$, the interfaces $\Gamma_0^\pm$ between them, and the boundaries $\Gamma_1^\pm$ at top and bottom.

An important issue of the discretization of eigenvalue problems is its stability, *i. e.* the existence of a minimal dimension of the FE space, such that the standard asymptotic convergence estimates hold for any dimension larger than this threshold. To the best of our knowledge, this issue has not yet been solved for the specific non-linear eigenvalue problem (8). However, numerical evidence shows that the standard asymptotic convergence estimates hold true.

Thus, we can use $p$-FEM on a coarse grid as our method of choice for the computation of guided modes in PhC wave-guides with smooth material boundaries and can expect exponential convergence.

*5.3. Solution of the discrete Ricatti equation*

The Ricatti equation (14) is fulfilled for any $\varphi \in H_{\mathrm{per}}^{1/2}(\Gamma_0^\pm)$ the operators are applied to. A discrete Ricatti equation results if we apply the operators to a basis of the discrete space $S_{\mathrm{per}}^{p,1}(\Gamma_0^\pm, \mathcal{M}(\Gamma_0^\pm))$ and take the duality product with this basis. Using the definition of the matrices $\mathbf{P}^\pm(\omega, k)$ in (37) and $\mathbf{T}_{ij}^\pm(\omega, k)$ in (39) we can write this as a linear system of equations

$$\mathbf{T}_{10}^\pm(\omega, k)(\mathbf{P}^\pm(\omega, k))^2 + (\mathbf{T}_{00}^\pm(\omega, k) + \mathbf{T}_{11}^\pm(\omega, k))\mathbf{P}^\pm(\omega, k) + \mathbf{T}_{01}^\pm(\omega, k) = \mathbf{0}. \tag{40}$$

Considering that the discretization preserves the periodicity properties of $C_1^\pm$ in $\mathbf{a}_2$-direction we deduce that the propagation matrix $\mathbf{P}^\pm(\omega, k)$ is the unique matrix satisfying Eq. (40) with eigenvalues whose magnitude is strictly less than 1.

In [20] Joly et. al. proposed a modified Newton method to solve the matrix valued problem (40) where the spectral constraint $\rho(\mathbf{P}^\pm(\omega, k)) < 1$ is integrated implicitly into the algorithm. This modified Newton
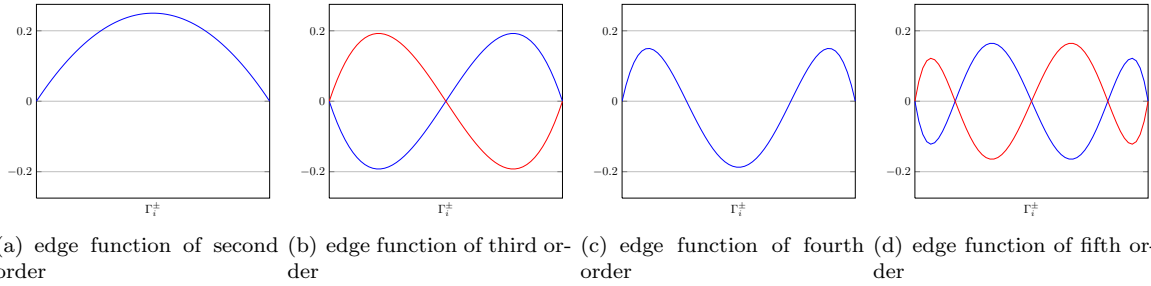
(a) edge function of second order  (b) edge function of third order  (c) edge function of fourth order  (d) edge function of fifth order

Figure 4: Basis functions according to Karniadakis and Sherwin [37] of the trace spaces $H_{\mathrm{per}}^{1/2}(\Gamma_i^\pm)$, $i = 0, 1$, with no $h$-refinement corresponding to the mesh in Figure 3. Note that due to periodicity, the first order basis function, which is not shown in this figure, is constant with value 1. While the basis functions of even order are uniquely defined, the basis functions of odd order are not unique (blue and red curve) and depend on the local and global orientation of the edge.

method only requires the matrix $\mathbf{T}_{00}^\pm(\omega, k) + \mathbf{T}_{11}^\pm(\omega, k)$ to be invertible, which is guaranteed by the fact that the corresponding linear operator $\mathcal{T}_{00}^\pm(\omega, k) + \mathcal{T}_{11}^\pm(\omega, k)$ is an isomorphism, see Remark 3.8.

Another method that was sketched in [20] is based on a spectral decomposition of the propagation matrix $\mathbf{P}^\pm(\omega, k)$. This spectral decomposition has two main advantages compared to the modified Newton method: first, its computational performance is better, and second, its results have a physical meaning as we will see later in Definition 5.3 and Remark 5.4. Even though it has not been proven that the propagation matrix $\mathbf{P}^\pm(\omega, k)$ is diagonalizable — in fact Hohage and Soussi [38] showed that the propagation operator $\mathcal{P}^\pm(\omega, k)$ of the TM mode is of Jordan type — we will use this spectral method because in practise it seems that the matrix is always diagonalizable. But also if this should not be the case, and the propagation matrix is of Jordan type, we can still use this spectral method in a generalized form by identifying the Jordan blocks and computing the Jordan chains. See [32] for more details. Thus, we seek eigenvalues $\mu^\pm(\omega, k) \in \mathbb{C}$ with magnitude strictly less than 1 and their corresponding eigenvectors $\boldsymbol{\psi}^\pm(\omega, k) \in \mathbb{C}^{N(\Gamma_0^\pm)}$ of the quadratic eigenvalue problem

$$\left[ \mathbf{T}_{10}^\pm(\omega, k) \left( \mu^\pm(\omega, k) \right)^2 + \left( \mathbf{T}_{00}^\pm(\omega, k) + \mathbf{T}_{11}^\pm(\omega, k) \right) \mu^\pm(\omega, k) + \mathbf{T}_{01}^\pm(\omega, k) \right] \boldsymbol{\psi}^\pm(\omega, k) = 0, \qquad (41)$$

which can be transformed into the generalized linear eigenvalue problem

$$\begin{pmatrix} -(\mathbf{T}_{00}^\pm(\omega, k) + \mathbf{T}_{11}^\pm(\omega, k)) & -\mathbf{T}_{01}^\pm(\omega, k) \\ \mathbf{I} & \mathbf{0} \end{pmatrix} \boldsymbol{\Psi}^\pm(\omega, k) = \mu^\pm(\omega, k) \begin{pmatrix} \mathbf{T}_{10}^\pm(\omega, k) & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{pmatrix} \boldsymbol{\Psi}^\pm(\omega, k), \qquad (42)$$

c.f. [39], with

$$\boldsymbol{\Psi}^\pm(\omega, k) = \begin{pmatrix} \mu^\pm(\omega, k) \boldsymbol{\psi}^\pm(\omega, k) \\ \boldsymbol{\psi}^\pm(\omega, k) \end{pmatrix}. \qquad (43)$$

Now let us come to an important symmetry property of the eigenvalues of the propagation matrix $\mathbf{P}^\pm(\omega, k)$. To this end, we first show

**Lemma 5.1.** *The matrices $\mathbf{T}_{ij}^\pm(\omega, k)$, $i, j = 1, 2$, corresponding to the linear operators $\mathcal{T}_{i,j}^\pm$, $i, j = 0, 1$, are Hermitian, i.e. they satisfy*

$$\mathbf{T}_{ij}^\pm(\omega, k)^{\mathrm{T}} = \overline{\mathbf{T}_{ji}^\pm(\omega, k)}, \qquad i, j = 0, 1. \qquad (44)$$

*Proof.* First we note that the system matrix $\mathbf{B}_{C_1^\pm}(\omega, k)$ is Hermitian which can be seen directly from the definition of the underlying sesquilinear form (18) considering that $\varepsilon, \omega^2, k \in \mathbb{R}$. Moreover, we recall that the matrices $\mathbf{Q}_i^\pm$, $i = 0, 1$, are real valued. Then the result follows directly from the definition of the matrices $\mathbf{T}_{ij}^\pm(\omega, k)$ in (39). $\square$

Using Lemma 5.1 it is easy to see that the quadratic eigenvalue problem (41) satisfies

**Proposition 5.2.** *If $\mu^\pm(\omega, k) \in \mathbb{C} \setminus \{0\}$ is an eigenvalue of (41), then $\left( \overline{\mu^\pm(\omega, k)} \right)^{-1}$ is also an eigenvalue.*

*Proof.* Taking the complex conjugate of (41) and inserting (44) yields

$$\left[ \mathbf{T}_{01}^\pm(\omega, k) \left( \overline{\mu^\pm(\omega, k)} \right)^2 + \left( \mathbf{T}_{00}^\pm(\omega, k) + \mathbf{T}_{11}^\pm(\omega, k) \right) \overline{\mu^\pm(\omega, k)} + \mathbf{T}_{10}^\pm(\omega, k) \right]^{\mathrm{T}} \overline{\boldsymbol{\psi}^\pm(\omega, k)} = 0.$$

Multiplying with $\left(\overline{\mu^\pm(\omega,k)}\right)^{-2}$ and taking the transpose gives

$$\overline{\boldsymbol{\psi}^\pm(\omega,k)}^{\mathrm{T}} \left[\mathbf{T}_{10}^\pm(\omega,k)\left(\overline{\mu^\pm(\omega,k)}\right)^{-2} + \left(\mathbf{T}_{00}^\pm(\omega,k) + \mathbf{T}_{11}^\pm(\omega,k)\right)\left(\overline{\mu^\pm(\omega,k)}\right)^{-1} + \mathbf{T}_{01}^\pm(\omega,k)\right] = 0.$$

This implies that there exists a vector $\widetilde{\boldsymbol{\psi}}^\pm(\omega,k) \in \mathbb{C}^{N(\Gamma_0^\pm)}$ such that

$$\left[\mathbf{T}_{10}^\pm(\omega,k)\left(\overline{\mu^\pm(\omega,k)}\right)^{-2} + \left(\mathbf{T}_{00}^\pm(\omega,k) + \mathbf{T}_{11}^\pm(\omega,k)\right)\left(\overline{\mu^\pm(\omega,k)}\right)^{-1} + \mathbf{T}_{01}^\pm(\omega,k)\right]\widetilde{\boldsymbol{\psi}}^\pm(\omega,k) = 0,$$

and hence, $\left(\overline{\mu^\pm(\omega,k)}\right)^{-1}$ is an eigenvalue of (41) with (right) eigenvector $\widetilde{\boldsymbol{\psi}}^\pm(\omega,k)$ and left eigenvector $\overline{\boldsymbol{\psi}^\pm(\omega,k)}^{\mathrm{T}}$. $\qquad\square$

An advantage of the spectral decomposition, that also contributes to its better computational performance compared to the modified Newton method, is that we can directly determine whether $\omega^2$ is inside the discrete approximation of the essential spectrum $\sigma_{\mathrm{ess}}(k)$.

**Definition 5.3.** *We call the approximative essential spectrum $\sigma_{\mathrm{ess}}^{h,\pm}(k)$ (of the operator $\mathcal{A}^\pm(k)$) related to the eigenvalue problem (3) in discretized form the set of numbers $\omega^2$ for which the quadratic eigenvalue problem (41) has eigenvalues with magnitude 1. Furthermore, let $\sigma_{\mathrm{ess}}^h(k) := \sigma_{\mathrm{ess}}^{h,+}(k) \cup \sigma_{\mathrm{ess}}^{h,-}(k)$.*

With the help of Propositions 5.2 and Definition 5.3 it is now clear how to compute the spectral decomposition of the propagation matrix $\mathbf{P}^\pm(\omega,k)$. We solve the general eigenvalue problem (42) for its $2N(\Gamma_0^\pm)$ eigenvalues $\mu^\pm(\omega,k)$. If there exist eigenvalues with magnitude equal to 1 we stop our computation as we know from Definition 5.3 that this means that $\omega^2$ is in the approximative essential spectrum $\sigma_{\mathrm{ess}}^h(k)$. Otherwise, and in accordance to Proposition 5.2, the $2N(\Gamma_0^\pm)$ eigenvalues $\mu^\pm(\omega,k)$ split into $N(\Gamma_0^\pm)$ eigenvalues with magnitude strictly less than 1 and $N(\Gamma_0^\pm)$ eigenvalues with magnitude strictly larger than 1. While discarding the $N(\Gamma_0^\pm)$ eigenvalues with magnitude strictly larger than 1, the $N(\Gamma_0^\pm)$ eigenvalues $\mu^\pm(\omega,k)$ with magnitude strictly less than 1 and their corresponding eigenvectors $\boldsymbol{\psi}^\pm(\omega,k)$ form the spectral decomposition of the propagation matrix $\mathbf{P}^\pm(\omega,k)$.

Note at this point that we do not introduce a modelling error when we compute the propagation matrix since all $2N(\Gamma_0^\pm)$ eigenvalues of the general eigenvalue problem (42) are computed. Thus, the only error that we expect is due to the choice of the discretization.

**Remark 5.4.** *Finally, we remark that — assuming $\mathbf{P}^\pm(\omega,k)$ is diagonalizable — the eigenvectors of $\mathbf{P}^\pm(\omega,k)$ form a basis of the traces of the discretized evanescent PhC modes.*

*5.4. Definition of the discrete DtN operator*

In the variational formulation (8) the sesquilinear form $\mathsf{d}_{C_0}$ is related to the solution in the two semi-infinite strips which is represented by the Dirichlet-to-Neumann maps $\Lambda^\pm(\omega,k)$. When inserting the basis functions $b_{\Gamma_0^\pm,n}$ of $S_{\mathrm{per}}^{p,1}(\Gamma_0^\pm, \mathcal{M}(\Gamma_0^\pm))$ into each of the two integrals in $\mathsf{d}_{C_0}$ and using the definition of the DtN operators (16) we obtain the matrices

$$\mathbf{D}^\pm(\omega,k) = \mathbf{T}_{00}^\pm(\omega,k) + \mathbf{T}_{10}^\pm(\omega,k)\mathbf{P}^\pm(\omega,k) \in \mathbb{C}^{N(\Gamma_0^\pm) \times N(\Gamma_0^\pm)}.$$

However, stating the variational formulation (8) in $S_{\mathrm{per}}^{p,1}(C_0, \mathcal{M}(C_0)) \subset H_{\mathrm{per}}^1(C_0)$ we have to insert in $\mathsf{d}_{C_0}$ rather the traces of the basis functions $b_{C_0,n}$, $n = 0, \ldots, N(C_0) - 1$, on $\Gamma_0^\pm$. Ordering these basis functions in the same way as for $S_{\mathrm{per}}^{p,1}(C_1^\pm, \mathcal{M}(C_1^\pm))$, *i.e.* the basis functions $b_{C_0,n}$ of $S_{\mathrm{per}}^{p,1}(C_0, \mathcal{M}(C_0))$

- with index $0, \ldots, N(C_0) - 1$ vanish on $\Gamma_0^\pm$,

- with index $N(C_0), \ldots, N(C_0) + N(\Gamma_0^+) - 1$ vanish on $\Gamma_0^-$, but their traces on $\Gamma_0^+$ build a basis of $S_{\mathrm{per}}^{p,1}(\Gamma_0^+, \mathcal{M}(\Gamma_0^+))$, and

- with index $N(C_0) + N(\Gamma_0^+), \ldots, N(C_0) + N(\Gamma_0^+) + N(\Gamma_0^-) - 1$ vanish on $\Gamma_0^+$, but their traces on $\Gamma_0^-$ build a basis of $S_{\mathrm{per}}^{p,1}(\Gamma_0^-, \mathcal{M}(\Gamma_0^-))$,

we obtain analogously to Eq. (38)

$$b_{\Gamma_0^\pm,n} = \sum_{m=0}^{N(\Gamma_0^\pm)-1} Q_{mn}^\pm b_{C_0,m+N(C_0)}\big|_{\Gamma_0^\pm}$$

with matrices $\mathbf{Q}^{\pm} \in \mathbb{R}^{N(\Gamma_0^{\pm}) \times N(\Gamma_0^{\pm})}$. The discrete form of the non-linear eigenvalue problem (8) then reads

$$\left(\mathbf{A}_{C_0}(k) - \omega^2 \mathbf{M}_{C_0} + \mathbf{D}_{C_0}(\omega, k)\right) \mathbf{u}(\omega, k) = \mathbf{0} \tag{45}$$

where $\mathbf{u}(\omega, k) \in \mathbb{C}^N$, with $N := N(C_0) + N(\Gamma_0^+) + N(\Gamma_0^-)$, is the coefficient vector of the eigenmode $u(\omega, k)$, $\mathbf{A}_{C_0}(k)$ and $\mathbf{M}_{C_0}$ are the matrices related to the sesquilinear forms (9a) and (9b), respectively, and $\mathbf{D}_{C_0}(\omega, k)$ is a block matrix of the form

$$\mathbf{D}_{C_0}(\omega, k) = \begin{pmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & (\mathbf{Q}^+)^{\mathrm{T}} \mathbf{D}^+(\omega, k) \mathbf{Q}^+ & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & (\mathbf{Q}^-)^{\mathrm{T}} \mathbf{D}^-(\omega, k) \mathbf{Q}^- \end{pmatrix} \in \mathbb{C}^{N \times N}.$$

## 6. Numerical solution of the non-linear eigenvalue problem

In this section we introduce techniques to solve the non-linear eigenvalue problem (8). It is split into two parts: in the first part we show the application of a direct procedure to solve (8) and in the second part we introduce iterative solution techniques.

### 6.1. Direct Procedure

The direct procedure to solve the non-linear eigenvalue problem (8) is based on a recent work by Effenberger and Kressner [40], who proposed a linearization of non-linear eigenvalue problems using the Chebyshev interpolation. Let us directly consider the discrete non-linear eigenvalue problem (45) in $\omega$-formulation

$$\mathbf{N}_k(\omega)\mathbf{u} = 0, \tag{46}$$

with $\mathbf{N}_k : \omega \mapsto \mathbf{A}_{C_0}(k) - \omega^2 \mathbf{M}_{C_0} - \mathbf{D}_{C_0}(\omega, k)$, and in $k$-formulation

$$\mathbf{N}_\omega(k)\mathbf{u} = 0, \tag{47}$$

with $\mathbf{N}_\omega : k \mapsto \mathbf{A}_{C_0}(k) - \omega^2 \mathbf{M}_{C_0} - \mathbf{D}_{C_0}(\omega, k)$. According to Remark 3.13, and assuming that $\omega^2 \notin \sigma_{\mathrm{ess}}^{h,\pm}(k)$ and that the problem (5) is well posed in a neighbourhood of $\omega^2$, we can deduce that both complex matrix-valued functions $\mathbf{N}_k$ and $\mathbf{N}_\omega$ are differentiable to any order and hence, the Chebyshev interpolation of non-linear eigenvalue problems, as described in [40], can be applied to (46) and (47).

For simplicity let us concentrate on the $\omega$-formulation and describe the procedure to linearize $\mathbf{N}_k$. We try to find a polynomial approximation $\mathbf{N}_k^d$ of order $d$ for the non-linear function $\mathbf{N}_k$. To this end, we first fix an interval $I_\omega = [\omega_a, \omega_b] \subset \mathbb{R}^+ \setminus \sigma_{\mathrm{ess}}^h(k)$ on which we project the $d + 1$ Chebyshev nodes $\cos(\frac{i+0.5}{d+1}\pi)$, $i = 0, \ldots, d$, of the first kind to obtain the $d + 1$ projected Chebyshev nodes

$$\omega_i = \frac{\omega_b - \omega_a}{2} \cos\left(\frac{i + 0.5}{d + 1} \pi\right) + \frac{\omega_a + \omega_b}{2} \in I_\omega, \qquad i = 0, \ldots, d.$$

Let $c_j : I_\omega \to \mathbb{R}$, $\omega \mapsto c_j(\omega)$, $j = 0, \ldots, d$, denote the first $d + 1$ Chebyshev polynomials defined on the interval $I_\omega$, i. e.

$$\begin{aligned} c_0(\omega) &= 1, \\ c_1(\omega) &= \omega, \\ c_{j+2}(\omega) &= 2\omega c_{j+1}(\omega) - c_j(\omega), \qquad j = 0, \ldots, d - 2. \end{aligned} \tag{48}$$

Then we approximate

$$\mathbf{N}_k(\omega) \approx \mathbf{N}_k^d(\omega) = \sum_{j=0}^d \mathbf{C}_j c_j(\omega) \tag{49}$$

where the $d + 1$ matrices $\mathbf{C}_j \in \mathbb{C}^{N \times N}$ are given by the interpolation condition

$$\mathbf{N}_k(\omega_i) = \sum_{j=0}^d \mathbf{C}_j c_j(\omega_i) = \sum_{j=0}^d \mathbf{C}_j \cos \frac{j(i + 0.5)\pi}{d + 1}$$

for all $i = 0, \ldots, d$, which can be solved efficiently for $\mathbf{C}_j$ using the discrete cosine transformation [41] of the second type, i. e.

$$\mathbf{C}_0 = \frac{1}{d + 1} \sum_{i=0}^d \mathbf{N}_k(\omega_i),$$

$$\mathbf{C}_j = \frac{2}{d + 1} \sum_{i=0}^d \mathbf{N}_k(\omega_i) \cos \frac{j(i + 0.5)\pi}{d + 1}, \qquad j = 1, \ldots, d.$$

Substituting $\mathbf{u}_j(\omega) := c_j(\omega)\mathbf{u}$ and searching the kernel of $\mathbf{N}_k^d(\omega)$ defined in (49), we obtain the polynomial eigenvalue problem

$$\sum_{j=0}^{d} \mathbf{C}_j \mathbf{u}_j(\omega) = 0,$$

that can be linearized into the general linear eigenvalue problem

$$\begin{pmatrix} \mathbf{0} & \mathbf{I} & & & \\ \mathbf{I} & \mathbf{0} & \mathbf{I} & & \\ & \ddots & \ddots & & \ddots \\ & & \mathbf{I} & \mathbf{0} & \mathbf{I} \\ -\mathbf{C}_0 & \cdots & -\mathbf{C}_{d-3} & \mathbf{C}_d - \mathbf{C}_{d-2} & -\mathbf{C}_{d-1} \end{pmatrix} \begin{pmatrix} \mathbf{u}_0 \\ \vdots \\ \mathbf{u}_{d-1} \end{pmatrix} = \omega \begin{pmatrix} \mathbf{I} & & & \\ & 2\mathbf{I} & & \\ & & \ddots & \\ & & & 2\mathbf{I} \\ & & & & 2\mathbf{C}_d \end{pmatrix} \begin{pmatrix} \mathbf{u}_0 \\ \vdots \\ \mathbf{u}_{d-1} \end{pmatrix} \quad (50)$$

of dimension $d\,N$, where we used the three term recurrence relation (48) of the Chebyshev polynomials.

Applying a shift and invert strategy to compute the eigenvalues of (50) the matrix of the left hand side needs to be inverted. However, due to the structure of this matrix its inverse can be determined by simply inverting a matrix of size $N \times N$, [40]. To explain this, let us assume we want to solve

$$\begin{pmatrix} \mathbf{0} & \mathbf{I} & & & \\ \mathbf{I} & \mathbf{0} & \mathbf{I} & & \\ & \ddots & \ddots & & \ddots \\ & & \mathbf{I} & \mathbf{0} & \mathbf{I} \\ -\mathbf{C}_0 & \cdots & -\mathbf{C}_{d-3} & \mathbf{C}_d - \mathbf{C}_{d-2} & -\mathbf{C}_{d-1} \end{pmatrix} \begin{pmatrix} \mathbf{x}_0 \\ \vdots \\ \mathbf{x}_{d-1} \end{pmatrix} = \begin{pmatrix} \mathbf{y}_0 \\ \vdots \\ \mathbf{y}_{d-1} \end{pmatrix}.$$

We can then deduce that

$$\mathbf{x}_1 = \mathbf{y}_0, \qquad \mathbf{x}_{2j+1} = \mathbf{y}_{2j} - \mathbf{x}_{2j-1}, \quad j = 1, \ldots, \lfloor d/2 \rfloor - 1,$$

and

$$\mathbf{x}_{2j} = \widetilde{\mathbf{y}}_{2j-1} + (-1)^j \mathbf{x}_0, \quad j = 1, \ldots, \lfloor (d-1)/2 \rfloor,$$

with

$$\widetilde{\mathbf{y}}_1 = \mathbf{y}_1, \qquad \widetilde{\mathbf{y}}_{2j+1} = \mathbf{y}_{2j+1} - \widetilde{\mathbf{y}}_{2j-1}, \quad j = 1, \ldots, \lfloor (d-1)/2 \rfloor - 1.$$

Finally, $\mathbf{x}_0$ is determined as solution of the $N$-dimensional linear system

$$\left( \sum_{j=0}^{\lfloor d/2 \rfloor} (-1)^{j+1} \mathbf{C}_{2j} \right) \mathbf{x}_0 = \mathbf{y}_{d-1} + \left( \sum_{j=0}^{\lfloor d/2 \rfloor - 1} \mathbf{C}_{2j+1} \mathbf{x}_{2j+1} \right) + \left( \sum_{j=1}^{\lfloor (d-1)/2 \rfloor} \mathbf{C}_{2j} \widetilde{\mathbf{y}}_{2j-1} \right) - \mathbf{C}_d \mathbf{b},$$

with

$$\mathbf{b} = \begin{cases} \widetilde{\mathbf{y}}_{d-3}, & \text{if } d \text{ is even}, \\ \mathbf{x}_{d-2}, & \text{if } d \text{ is odd}. \end{cases}$$

Thus, it simply remains to invert the matrix $\left( \sum_{j=0}^{\lfloor d/2 \rfloor} (-1)^{j+1} \mathbf{C}_{2j} \right) \in \mathbb{C}^{N \times N}$.

*6.2. Iterative Procedure*

In Proposition 3.6 we showed that the functions $g_m^{(1)}(\omega_\Lambda, k) = \omega_m^2(\omega_\Lambda, k)$, $m \in \mathbb{N}$, *i.e.* the eigenvalues of the linear eigenvalue problem (11) in $\omega$-formulation with respect to $(\omega_\Lambda^2, k_\Lambda) \in \mathbb{R}^+ \times B$ are continuous and hence, the distance functions

$$d_{\omega,m}^{(1)}(\omega_\Lambda, k) = \omega_\Lambda^2 - g_m^{(1)}(\omega_\Lambda, k) = \omega_\Lambda^2 - \omega_m^2(\omega_\Lambda, k) \in \mathbb{R}$$

are continuous for all $m \in \mathbb{N}$. On the other hand, Proposition 3.7 guarantees that there exists an alternative ordering $m \mapsto \widetilde{m}(m)$ of the eigenvalues such that the functions $g_m^{(2)}(\omega_\Lambda, k) = \omega_{\widetilde{m}(m)}^2(\omega_\Lambda, k)$, $m \in \mathbb{N}$, are differentiable and hence, the Newton method can be applied to find the roots of

$$d_{\omega,m}^{(2)}(\omega_\Lambda, k) = \omega_\Lambda^2 - g_m^{(2)}(\omega_\Lambda, k) = \omega_\Lambda^2 - \omega_{\widetilde{m}(m)}^2(\omega_\Lambda, k), \quad (51)$$

which are then eigenvalues of the non-linear eigenvalue problem (8). Numerical results for this procedure can be found in [24] where, however, the Newton method is directly applied to the continuous distance function $d_{\omega,m}^{(1)}$.

Let us now introduce the global distance function

$$d_\omega(\omega_\Lambda, k) = d^{(2)}_{\omega,m^*}(\omega_\Lambda, k) \tag{52}$$

where

$$m^* = \arg\min_{m \in \mathbb{N}} |d^{(2)}_{\omega,m}(\omega_\Lambda, k)|.$$

As shown in the numerical results in Section 7 this function is not continuous and hence not differentiable, however, we shall also see in Section 7 that the numerical results are reasonable when using the derivative of the continuously differentiable distance function $d^{(2)}_{\omega,m^*}(\cdot, k)$ with respect to $\omega_\Lambda$ in the Newton method applied to the global distance function $d_\omega(\omega_\Lambda, k)$. This implies that we save computational effort by simply applying the Newton method to $d_\omega(\cdot, k)$ instead of to the first, say, $M \in \mathbb{N}$ distance functions $d^{(2)}_{\omega,m}(\cdot, k)$, $m = 1, \ldots, M$, corresponding to the $M$ eigenvalues $\omega_m^2$ of smallest magnitude as done exemplarily for $M = 1$ in [24].

---

**Algorithm 1** Newton's method applied to global distance function

---

Fix $k \in B$ and choose start value $\omega^{(0)} \in \mathbb{R}^+$.
**for** $n = 0, \ldots$ **do**
 Set $\omega_\Lambda = \omega^{(n)}$.
 **if** $\omega_\Lambda^2 \in \sigma^h_{\mathrm{ess}}(k)$, *i. e.* propagation operators $\mathcal{P}^\pm(\omega_\Lambda, k)$ are not defined and there exists a solution of
 Eq. (14) with spectral radius equal to one, **then**
  **exit**
 **end if**
 Compute DtN maps $\Lambda^\pm(\omega_\Lambda, k)$.
 Solve linear eigenvalue problem (11) for $\omega^2$ keeping $\Lambda^\pm(\omega, k) = \Lambda^\pm(\omega_\Lambda, k)$ fixed.
 Evaluate global distance function $d_\omega(\omega_\Lambda, k)$.
 **if** $d_\omega(\omega_\Lambda, k) \approx 0$ **then**
  **exit**
 **end if**
 Compute new value $\omega^{(n+1)} = \omega^{(n)} - \left( \frac{\partial}{\partial \omega_\Lambda} d^{(2)}_{\omega,m^*}(\omega_\Lambda, k) \right)^{-1} d_\omega(\omega_\Lambda, k)$.
**end for**

---

The iterative scheme then works as shown in Algorithm 1, where the computation of the derivative of $d^{(2)}_{\omega,m^*}(\cdot, k)$ with respect to $\omega_\Lambda$ follows similar ideas as the computation of the group velocity in Section 4. Again we note that Proposition 3.7 guarantees that the eigenvalues $\omega^2_{\widetilde{m}(m)}(\omega_\Lambda, k)$ of the linear eigenvalue problem (11) are continuously differentiable with respect to $\omega_\Lambda$. Thus, we can differentiate (11) which yields

$$\int_{C_0} 2\omega_{\widetilde{m}}(\omega_\Lambda, k) \left( \frac{\partial}{\partial \omega_\Lambda} \omega_{\widetilde{m}}(\omega_\Lambda, k) \right) \varepsilon u_{\widetilde{m}} \overline{v} \, \mathrm{d}\mathbf{x} - \int_{\Gamma_0^+} \Lambda^+_\omega(\omega_\Lambda, k) u_{\widetilde{m}} \overline{v} \, \mathrm{d}s(\mathbf{x}) - \int_{\Gamma_0^-} \Lambda^-_\omega(\omega_\Lambda, k) u_{\widetilde{m}} \overline{v} \, \mathrm{d}s(\mathbf{x})$$

$$= \int_{C_0} (\nabla + \mathrm{i}k\mathbf{a}_1) \frac{\partial u_{\widetilde{m}}}{\partial \omega_\Lambda} \cdot (\nabla - \mathrm{i}k\mathbf{a}_1)\overline{v} - \varepsilon \omega_{\widetilde{m}}^2(\omega_\Lambda, k) \frac{\partial u_{\widetilde{m}}}{\partial \omega_\Lambda} \overline{v} \, \mathrm{d}\mathbf{x}$$

$$+ \int_{\Gamma_0^+} \Lambda^+(\omega, k) \frac{\partial u_{\widetilde{m}}}{\partial \omega_\Lambda} \overline{v} \, \mathrm{d}s(\mathbf{x}) + \int_{\Gamma_0^-} \Lambda^-(\omega, k) \frac{\partial u_{\widetilde{m}}}{\partial \omega_\Lambda} \overline{v} \, \mathrm{d}s(\mathbf{x}), \quad (53)$$

for all $v \in H^1_{\mathrm{per}}(C_0)$, where $u_{\widetilde{m}}$ denotes the eigenfunction of (11) that corresponds to the eigenvalue $\omega^2_{\widetilde{m}(m)}(\omega_\Lambda, k)$. The right hand side of (53) vanishes if we test with $v = u_{\widetilde{m}}$ since $(\Lambda^\pm(\omega, k))^* = \Lambda^\pm(\omega, k)$, *c.f.* Eq. (30). Hence, we obtain

$$\frac{\partial}{\partial \omega_\Lambda} \omega_{\widetilde{m}}(\omega_\Lambda, k) = \frac{\int_{\Gamma_0^+} \Lambda^+_\omega(\omega_\Lambda, k) u_{\widetilde{m}} \overline{u_{\widetilde{m}}} \, \mathrm{d}s(\mathbf{x}) + \int_{\Gamma_0^-} \Lambda^-_\omega(\omega_\Lambda, k) u_{\widetilde{m}} \overline{u_{\widetilde{m}}} \, \mathrm{d}s(\mathbf{x})}{2\omega_{\widetilde{m}}(\omega_\Lambda, k) \int_{C_0} \varepsilon |u_{\widetilde{m}}|^2 \, \mathrm{d}\mathbf{x}}$$

and consequently, using Eq. (51),

$$\frac{\partial}{\partial \omega_\Lambda} d^{(2)}_{\omega,m}(\omega_\Lambda, k) = 2\omega_\Lambda - \frac{\int_{\Gamma_0^+} \Lambda^+_\omega(\omega_\Lambda, k) u_{\widetilde{m}} \overline{u_{\widetilde{m}}} \, \mathrm{d}s(\mathbf{x}) + \int_{\Gamma_0^-} \Lambda^-_\omega(\omega_\Lambda, k) u_{\widetilde{m}} \overline{u_{\widetilde{m}}} \, \mathrm{d}s(\mathbf{x})}{\int_{C_0} \varepsilon |u_{\widetilde{m}}|^2 \, \mathrm{d}\mathbf{x}}.$$

The selection of the start value $\omega^{(0)}$ is of particular importance for the convergence of the Newton method. Recall that in Section 2 we introduced two kinds of dispersive curves, the continuous but not necessarily differentiable dispersive curves $f_m^{(1)}(k)$ of the first kind and the continuously differentiable

dispersive curves $f_m^{(2)}(k)$ of the second kind, see Proposition 2.3. The group velocity introduced in Section 4 is the derivative of the dispersive curves of the second kind and can hence be used to estimate the behaviour of these curves around a vicinity of a known eigenvalue $\omega^2(k)$ at $k \in B$. That means in order to compute the eigenvalue $\omega^2(k + k_h)$ at $k + k_h \in B$, with $|k_h|$ reasonably small, we set the start value $\omega^{(0)}(k + k_h)$ in Algorithm 1 to $\omega^{(0)}(k + k_h) = \omega(k) + k_h \frac{\partial \omega}{\partial k}(k)$. However, note that this procedure cannot be used as a path following algorithm without taking into account that two dispersive curves can change their behaviour when they come close. It is for example possible that two curves appear to intersect but if one zooms into that region it becomes evident that they come very close but do not intersect and form a so called *mini-stop band* [42].

An iterative procedure for the $k$-formulation can be set up analogously. However, we have to note again that smoothness of the global distance function

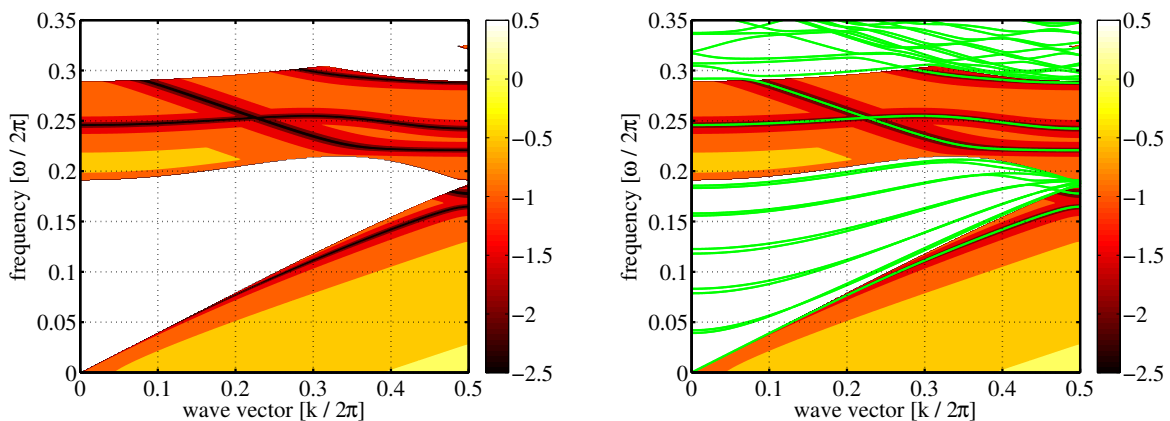$$d_k(\omega, k_\Lambda) = d_{k,m^*}(\omega, k_\Lambda) = k_\Lambda - k_{m^*}(\omega, k_\Lambda)$$

with

$$m^* = \arg \min_{m \in \mathbb{N}} |k_\Lambda - k_m(\omega, k_\Lambda)|,$$

is not guaranteed.

## 7. Numerical results

In this section we present numerical results of the proposed methods. The numerical example is based on a PhC wave-guide as described in [19], *i. e.* we study the TE mode of a PhC wave-guide of hexagonal lattice structure with $\mathbf{a}_1^0 = \mathbf{a}_1^+ = \mathbf{a}_1^- = (1,0)^{\mathrm{T}}$ and $\mathbf{a}_2^0 = \mathbf{a}_2^+ = \mathbf{a}_2^- = (0.5, \sqrt{0.75})^{\mathrm{T}}$, and with air holes ($\varepsilon = 1$) of radius 0.31 in a homogeneous and isotropic dielectric material of relative permittivity $\varepsilon = 11.4$.

The application of the super-cell method to this test setting was shown in [19]. We will briefly explain the major results of that study. Using finite elements with curved cells it is possible to create coarse meshes of the computational domain that provide an exact resolution of the domain with its circular holes, see Section 5. Hence, $p$-FEM of the super-cell method with no further $h$-refinement yields exponential convergence in both formulations, the $\omega$- and the $k$-formulation, when neglecting the modelling error introduced by the super-cell approach, [19]. However, note that this coarse mesh is relatively large compared to the mesh of the DtN method as shown in Figure 3 since it includes a certain number of periodicity cells (typically more than three) on both sides of the guide. Moreover, one has to have access to the essential spectrum in order to interpret the results of the super-cell method, *i. e.* all eigenvalues of the super-cell method that lie in the essential spectrum have to be excluded. Hence, a full computation of the essential spectrum is needed in advance.



(a) Magnitude of global distance function $d_\omega$ in logarithmic scale.

(b) Magnitude of global distance function $d_\omega$ in logarithmic scale compared to the results of super-cell method (green).

Figure 5: Magnitude of global distance function $d_\omega$ of $\omega$-formulation in logarithmic scale evaluated on a grid of $500 \times 500$ $(\omega, k)$ points, left and right. Results of the super-cell method with five rows of periodicity cells on top and bottom are given on the right (green lines).

### 7.1. Numerical results of the iterative procedure

Now let us come to our first numerical results. To give an orientation, in Figure 5(a) the magnitude of the global distance function $d_\omega$ of the $\omega$-formulation, *c. f.* Eq. (52), is plotted in logarithmic scale.

The essential spectrum is left blank, while dark lines indicate a small magnitude of the global distance function. In Figure 5(b) these results are compared with the results of the super-cell method using five rows of holes on top and bottom of the guide. For both computations, the DtN based computation of the global distance function and the computation of the super-cell method, the polynomial degree is set to $p = 7$. One can see that the dark lines, indicating small magnitudes of the global distance function, match well with the green lines of the super-cell method inside the band gap.



(a) $d_\omega$ with respect to $\omega$ for fixed $k = 0.3 \cdot 2\pi$.

(b) $d_k$ with respect to $\omega$ for fixed $k = 0.3 \cdot 2\pi$.

(c) $d_\omega$ with respect to $k$ for fixed $\omega = 0.25 \cdot 2\pi$.

(d) $d_k$ with respect to $k$ for fixed $\omega = 0.25 \cdot 2\pi$.

Figure 6: Global distance functions $d_\omega$ (left) and $d_k$ (right) evaluated on an equidistant mesh of $k$ in the interval $[0, \pi]$ for a fixed value of $\omega = 0.25 \cdot 2\pi$ (top), and evaluated on an equidistant mesh of $\omega$ in the interval $[0.22 \cdot 2\pi, 0.30 \cdot 2\pi]$ for a fixed value of $k = 0.3 \cdot 2\pi$ (bottom). The roots of the global distance functions are marked with red dots.

Resolving the results of the global distance function on a fine scale, as presented in Figure 6, one can see that the global distance functions $d_\omega$ and $d_k$ are not continuous (the sign swaps), but in a neighbourhood of their roots they are smooth. Note that when two dispersive curves cross, there will not exist a smooth neighbourhood of the global distance function. However, it will be continuous at this point, since it will tend to zero from both sides. Hence, the application of the Newton method to find the roots of the global distance functions, as proposed in Section 6, is reasonable as long as the start value is sufficiently close to the root. Note, that in certain cases the application of the iterative scheme to the distance function $d_\omega$ compared to $d_k$ or vice versa might be preferable as one can see in Figures 6(c) and 6(d), that shows that the distance function $d_\omega$ is clearly superior compared to $d_k$ when computing the eigenvalue near $k = 0.3835 \cdot 2\pi$ since it is less sensitive towards the start value of the iteration. In the following numerical results we will use the global distance function $d_\omega$.

Now we show convergence results of the Newton method applied to the global distance function $d_\omega$ of the $\omega$-formulation. Unless otherwise stated, the reference solution is computed by setting the polynomial degree to $p = 20$ and applying the same iterative scheme. In Figure 7(b) the convergence of the absolute errors of a well confined mode and a mode close to the band edge (see Figure 7(a)) are shown for an increasing polynomial degree. The real part of the magnetic field component of these two modes can be seen in Figures 9(a) and 9(b) respectively. Note that there is no value for the error of the mode close to the band edge for the lowest polynomial degree $p = 3$ since for this degree the mode is inside the essential spectrum and can therefore not be captured. As expected for $p$-FEM, we can observe exponential

(a) Location of well confined mode (blue) and mode close to band edge (green) in band structure.



(b) Convergence of well confined mode (blue) and mode close to band edge (green) with respect to polynomial degree $p$.
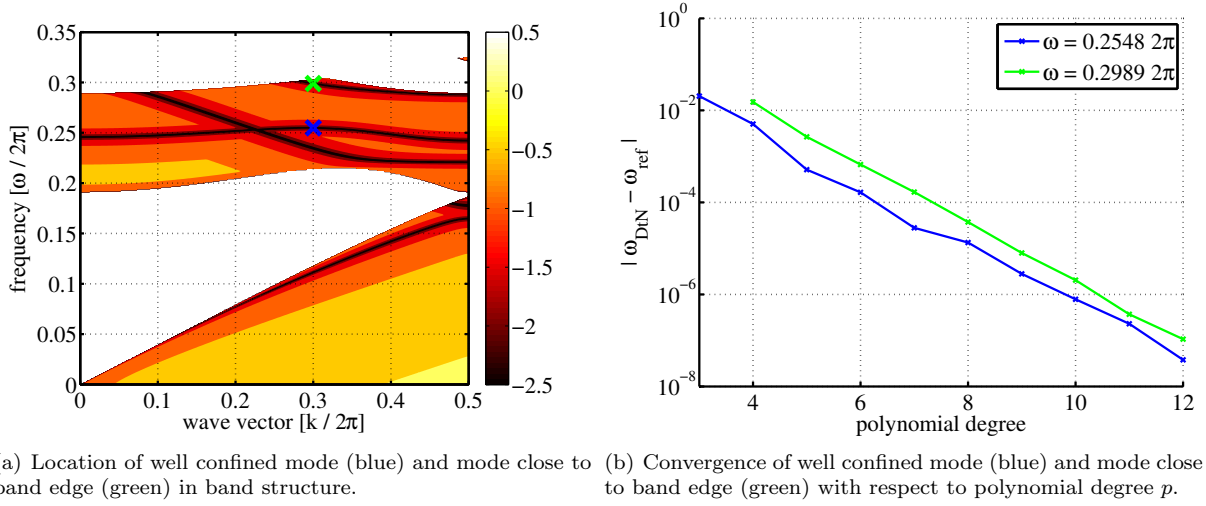
Figure 7: Location and convergence analysis of a well confined mode (blue) and a mode close to the band edge (green) computed with the Newton method applied to the global distance function $d_\omega$ of the $\omega$-formulation. The reference solution $\omega_{\text{ref}}$ is computed taking a polynomial degree of $p = 20$.

convergence with the same convergence rate for both modes.



(a) Fixed polynomial degree $p = 7$.



(b) Fixed number of periodicity cells $n = 3$ (dashed lines) and $n = 7$ (solid lines).

Figure 8: Convergence of the well confined mode (blue) and the mode close to the band edge (green) with respect to the number of the periodicity cells $n$ while keeping the polynomial degree $p$ fixed (left), and with respect to the polynomial degree $p$ while keeping the number of periodicity cells $n$ fixed (right). The reference solution $\omega_{\text{ref}}$ is computed using the iterative DtN method with a polynomial degree $p = 7$ (left) and $p = 20$ (right).

## 7.2. Comparison to the numerical results of the super-cell method

The convergence results of the super-cell method are presented in Figure 8. On the left we observe an exponential convergence of the results of the super-cell method with polynomial degree $p = 7$ towards the solution of the iterative DtN method of the same polynomial degree when increasing the number of periodicity cells on top and bottom. However, the rates of convergence differ significantly. The rate of the mode close to the band edge (green) is much smaller than the rate of the well confined mode (blue), see Figure 8(a). On the other hand, the figure on the right, where the number of periodicity cells is kept fixed to $n = 3$ and $n = 7$ while the polynomial degree is increased from $p = 3$ to $p = 12$, shows that the error of the super-cell method only converges exponentially towards the solution of the iterative DtN method with polynomial degree $p = 20$ until a certain error plateau is reached. This error plateau, which is due to the modelling error introduced by the super-cell approach, is significantly larger for the mode close the the band edge compared to the well-confined mode, see Figure 8(b). These results clearly demonstrate that the super-cell method is a good approximation of the exact DtN method for well-confined modes but for modes close to the band edge it produces errors of significantly larger orders.

24

(a) Well confined, odd mode at $\omega = 0.2548 \cdot 2\pi$.  (b) Odd mode close to band edge at $\omega = 0.2989 \cdot 2\pi$.  (c) Even mode at $\omega = 0.2347 \cdot 2\pi$.

Figure 9: Real part of the magnetic field component for three guided modes at $k = 0.3 \cdot 2\pi$. Computed with a polynomial degree of $p = 7$.

### 7.3. Numerical results of the direct procedure

Let us now come to the numerical results of the direct procedure to solve the non-linear eigenvalue problem. The Chebyshev interpolation requires a priori knowledge of the essential spectrum, since the differentiability of the non-linear eigenvalue problem can only be guaranteed outside the essential spectrum, see Section 6. Thus, this method is particularly interesting, if we apply it to the $k$-formulation in a band gap of the hole Brillouin zone, i.e. an interval of frequencies $\omega$ that are outside the essential spectrum for all values of $k$ in the Brillouin zone $B$. The convergence of the Chebyshev interpolation is shown in Figure 10 where the results of the direct procedure to compute the eigenvalues in the band gap $[0.22 \cdot 2\pi, 0.28 \cdot 2\pi]$ using the Chebyshev interpolation is compared to a reference solution computed with the iterative method. We observe an exponential convergence of the mean error of the eigenvalues computed at a sample of 200 frequencies in the band gap.
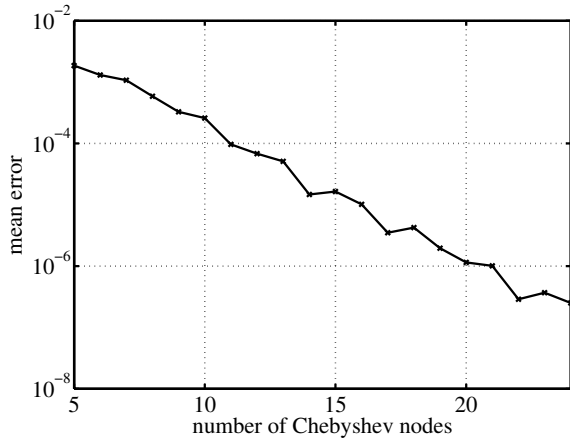


Figure 10: Convergence of the mean error $\frac{1}{J} \sum_{j=1}^{J} |k_{\mathrm{Cheb},j} - k_{\mathrm{ref},j}|$ of the eigenvalues $k_{\mathrm{Cheb},j}$ using the Chebyshev interpolation in the hole Brillouin zone of a sample of 200 frequencies in the band gap $[0.22 \cdot 2\pi, 0.28 \cdot 2\pi]$ with respect to the number of Chebyshev nodes $d$. The polynomial degree is $p = 7$ and the reference solutions $k_{\mathrm{ref},j}$ are computed using the iterative DtN method with the same polynomial degree.
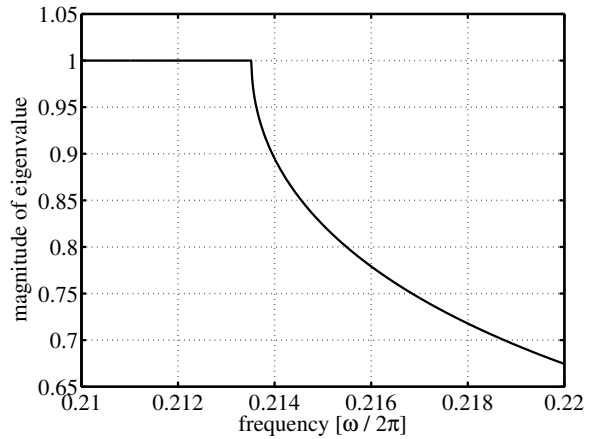
Figure 11: Magnitude of the eigenvalue of the propagation operator for $k = 0.3 \cdot 2\pi$ near the band edge at $\omega \approx 0.2135 \cdot 2\pi$. The eigenvalue has magnitude strictly less than 1 in the band gap and equal to 1 in the essential spectrum.
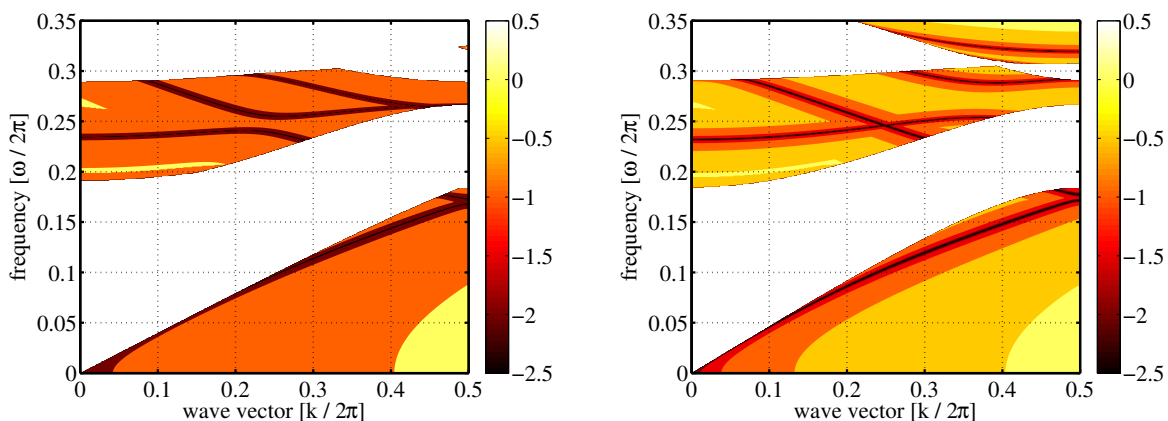
Analogously, we can observe exponential convergence towards the iterative solution if the direct procedure is applied to the $\omega$-formulation. However, we have to take care that the interval is sufficiently far away from the band edge, since the magnitude of the eigenvalue $\mu$ of the propagation operator, that

changes from $|\mu| = 1$ to $|\mu| < 1$ at the band edge, has a root-like singularity at the band edge as Figure 11 shows, and hence, its derivative becomes arbitrarily large near the band edge. This drawback of the Chebyshev interpolation is of smaller significance in the $k$-formulation as long as the chosen frequency interval is not arbitrarily close to the band edge at any $k \in B$, as demonstrated in the convergence analysis in Figure 10.

### 7.4. Numerical results for a PhC wave-guide with different periodicity on top and bottom of the guide

Finally, we demonstrate that the DtN method can be used to compute the band structure of a PhC wave-guide with different periodicity on top and bottom of the guide. We study a PhC wave-guide with the same isotropic dielectric material of relative permittivity $\varepsilon = 11.4$ as above. By convention, we have $\mathbf{a}_1^0 = \mathbf{a}_1^+ = \mathbf{a}_1^- = (1,0)^{\mathrm{T}}$. As before, we choose a hexagonal lattice structure on top of the guide with $\mathbf{a}_2^+ = (0.5, \sqrt{0.75})^{\mathrm{T}}$, and circular holes of radius 0.31. Below the guide, we consider a PhC with square lattice, $i.e.$ $\mathbf{a}_2^- = (0,1)^{\mathrm{T}}$, and circular holes of radius 0.4. For the guide we choose homogeneous medium and $\mathbf{a}_2^0 = (0.25, \sqrt{0.875})^{\mathrm{T}}$, $c.f.$ Figure 1(a).

The numerical results of the global distance function $d_\omega$ of the TE-mode can be found in Figure 12(a). For comparison, Figure 12(b) shows the global distance function $d_\omega$ of the TE-mode of the PhC wave-guide with square lattice on top and bottom of the guide. On the other hand, the global distance function of the PhC wave-guide with hexagonal lattice on top and bottom of the guide can be seen in Figure 5(a). Obviously, the essential spectrum (blank area) of the joint PhC wave-guide is the union of the essential spectra of the two wave-guides with conforming periodicities on top and bottom, $c.f.$ Proposition 2.1.



(a) Joint PhC wave-guide with hexagonal lattice (of radius 0.31) on top and square lattice (of radius 0.4) on bottom of the guide.

(b) PhC wave-guide with square lattice (of radius 0.4).

Figure 12: Magnitude of global distance function $d_\omega$ of $\omega$-formulation in logarithmic scale evaluated on a grid of $500 \times 500$ $(\omega, k)$ points of the joint PhC wave-guide (left) and the PhC wave-guide with square lattice (right).

## 8. Conclusion

We showed the numerical discretization of the DtN approach in [24] for the exact computation of guided modes in PhC wave-guides using high-order FEM. DtN maps for periodic media are computed by solving local Dirichlet problems and a quadratic eigenvalue problem. Using these DtN maps we transformed the eigenvalue problem for the computation of guided modes, that is posed on an unbounded domain, to a non-linear eigenvalue problem in a unit cell. We proposed an iterative and a direct procedure to solve the resulting non-linear eigenvalue problem. Numerical examples showed an exponential convergence for $p$-FEM independent of the confinement of the guided mode which stands in contrast to the super-cell method for which we showed a significant dependence on the confinement of the guided mode. We introduced a formula for the group velocity of guided modes that — on the one hand — is used to select the start values in the iterative procedure to solve the non-linear eigenvalue problem, and — on the other hand — the formula for the group velocity can deal as an objective function in the optimization process of PhC wave-guides to obtain slow light modes.

An interesting topic of future research is the extension of the proposed method to complex valued wave vectors $k$ and the computation of leaky modes, where the characterization of the DtN operators remains an open question.

**References**

[1] J. Joannopoulos, Photonic crystals: Molding the flow of light, Princeton University Press, Princeton, NJ, USA, 2008.

[2] P. Kuchment, The mathematics of photonic crystals, SIAM, Philadelphia, PA, USA, 2001, Ch. 7, pp. 207–272.

[3] K. Busch, Photonic band structure theory: assessment and perspectives, C. R. Physique 3 (1) (2002) 53–66.

[4] E. Istrate, A. Green, E. Sargent, Behavior of light at photonic crystal interfaces, Phys. Rev. B 71 (19) (2005) 195122.

[5] W. Jiang, R. T. Chen, X. Lu, Theory of light refraction at the surface of a photonic crystal, Phys. Rev. B 71 (2005) 245115.

[6] D. Boffi, M. Conforti, L. Gastaldi, Modified edge finite elements for photonic crystals, Numer. Math. 105 (2) (2006) 249–266.

[7] Z. Hu, Y. Y. Lu, Improved Dirichlet-to-Neumann map method for modeling extended photonic crystal devices, Opt. Quantum Electron. 40 (11-12) (2008) 921–932.

[8] C. Engström, M. Richter, On the spectrum of an operator pencil with applications to wave propagation in periodic and frequency dependent materials, SIAM J. Appl. Math. 70 (1) (2009) 231–247.

[9] V. Hoang, M. Plum, C. Wieners, A computer-assisted proof for photonic band gaps, Z. Angew. Math. Phys. 60 (6) (2009) 1035–1052.

[10] R. Norton, R. Scheichl, Convergence analysis of planewave expansion methods for 2d Schrödinger operators with discontinuous periodic potentials, SIAM J. Numer. Anal. 47 (6) (2010) 4356–4380.

[11] H. Brandsmeier, K. Schmidt, C. Schwab, A multiscale hp-FEM for 2D photonic crystal bands, J. Comput. Phys. 230 (2).

[12] S. Giani, I. G. Graham, Adaptive finite element methods for computing band gaps in photonic crystals, Numer. Math. 121 (1) (2012) 31–64.

[13] P. Kuchment, B. Ong, On guided waves in photonic crystal waveguides, in: P. Kuchment (Ed.), Waves in periodic and random media, Vol. 339 of Contemp. Math., American Math. Society, Providence, RI, USA, 2004, pp. 105–115.

[14] T. F. Krauss, Why do we need slow light?, Nat. Photonics 2 (2008) 448–450.

[15] J. Li, T. P. White, L. O'Faolain, A. Gomez-Iglesias, T. F. Krauss, Systematic design of flat band slow light in photonic crystal waveguides, Opt. Express 16 (9) (2008) 6227–6232.

[16] M. Soljacic, J. D. Joannopoulos, Enhancement of nonlinear effects using photonic crystals, Nat. Mater. 3 (2004) 211–219.

[17] D. Givoli, Recent advances in the DtN FE Method, Arch. Comput. Methods Eng. 6 (2) (1999) 71–116.

[18] S. Soussi, Convergence of the supercell method for defect modes calculations in photonic crystals, SIAM J. Numer. Anal. 43 (3) (2005) 1175–1201.

[19] K. Schmidt, R. Kappeler, Efficient computation of photonic crystal waveguide modes with dispersive material, Opt. Express 18 (7) (2010) 7307–7322.

[20] P. Joly, J.-R. Li, S. Fliss, Exact boundary conditions for periodic waveguides containing a local perturbation, Commun. Comput. Phys. 1 (6) (2006) 945–973.

[21] S. Fliss, P. Joly, Exact boundary conditions for time-harmonic wave propagation in locally perturbed periodic media, Appl. Numer. Math. 59 (9) (2009) 2155–2178.

[22] S. Fliss, E. Cassan, D. Bernier, New approach to describe light refraction at the surface of a photonic crystal, JOSA B 27 (2010) 1492–1503.

[23] S. Fliss, P. Joly, J.-R. Li, Exact boundary conditions for wave propagation in periodic media containing a local perturbation, no. 1, Bentham Science, 2010, Ch. 5, pp. 108–134.

[24] S. Fliss, A Dirichlet-to-Neumann approach for the exact computation of guided modes in photonic crystal waveguides, accepted for publication in SIAM J. Sci. Comput. (2012).

[25] C. Kittel, Introduction to solid state physics, 8th Edition, Wiley, New York, NY, USA, 2004.

[26] P. Kuchment, Floquet theory for partial differential equations, Birkhäuser, Basel, Switzerland, 1993.

[27] A. Figotin, A. Klein, Localized classical waves created by defects, J. Stat. Phys. 86 (1997) 165–177.

[28] M. Reed, B. Simon, Methods of modern mathematical physics, Vol. 1–4, Academic Press, New York, NY, USA, 1972–1978.

[29] M. Vorobets, On the Bethe–Sommerfeld conjecture for certain periodic Maxwell operators, J. Math. Anal. Appl. 377 (1) (2011) 370–383.

[30] H. Bethe, A. Sommerfeld, Elektronentheorie der Metalle, Heidelberger Taschenbücher, Springer, Berlin & Heidelberg, Germany, 1967.

[31] T. Katō, Perturbation theory for linear operators, Grundlehren der mathematischen Wissenschaften, Springer, Berlin & Heidelberg, Germany, 1995.

[32] S. Fliss, Etude mathématique et numérique de la propagation des ondes dans des milieux périodiques localement perturbés, Ph.D. thesis, École Doctorale de l'École Polytechnique, Palaiseau Cedex, France (May 2009).

[33] J. Coatléven, Helmholtz equation in periodic media with a line defect, J. Comput. Phys. 231 (4) (2012) 1675–1704.

[34] L. Brillouin, Wave propagation and group velocity, Pure and Applied Physics, Vol. 8, Academic Press, New York, NY, USA, 1960.

[35] R. Kappeler, P. Kaspar, P. Friedli, H. Jäckel, Design proposal for a low loss in-plane active photonic crystal waveguide with vertical electrical carrier injection, Opt. Express 20 (8) (2012) 9264–9275.

[36] C. Schwab, $p$- and $hp$-finite element methods: Theory and applications in solid and fluid mechanics, Oxford University Press, Oxford, UK, 1998.

[37] G. Karniadakis, S. Sherwin, Spectral/hp element methods for computational fluid dynamics, Numerical Mathematics and Scientific Computation, Oxford University Press, Oxford, UK, 2005.

[38] T. Hohage, S. Soussi, Riesz bases and Jordan form of the translation operator in semi-infinite periodic waveguides, J. Math. Pures Appl., in press (2012).

[39] F. Tisseur, K. Meerbergen, The quadratic eigenvalue problem, SIAM Rev. 43 (2) (2001) 235–286.

[40] C. Effenberger, D. Kressner, Chebyshev interpolation for nonlinear eigenvalue problems, BIT 52 (4) (2012) 933–951.

[41] N. Ahmed, T. Natarajan, K. Rao, Discrete cosine transform, IEEE Trans. Comput. 100 (1) (1974) 90–93.

[42] S. Olivier, M. Rattier, H. Benisty, C. Weisbuch, C. Smith, R. De la Rue, T. Krauss, U. Oesterle, R. Houdre, Mini-stopbands of a one-dimensional system: The channel waveguide in a two-dimensional photonic crystal, Phys. Rev. B 63 (11) (2001) 1133111–1133114.