

## CONVERGENCE ANALYSIS OF THE FEM COUPLED WITH FOURIER-MODE EXPANSION FOR THE ELECTROMAGNETIC SCATTERING BY BIPERIODIC STRUCTURES\*

GUANGHUI HU<sup>†</sup> AND ANDREAS RATHSFELD<sup>†</sup>

**Abstract.** Scattering of time-harmonic electromagnetic plane waves by a doubly periodic surface structure in  $\mathbb{R}^3$  can be simulated by a boundary value problem of the time-harmonic curl-curl equation. For a truncated FEM domain, non-local boundary conditions are required in order to satisfy the radiation conditions for the upper and lower half spaces. As an alternative to boundary integral formulations, to approximate radiation conditions and absorbing boundary methods, Huber et al. [SIAM J. Sci. Comput., 31 (2009), pp. 1500–1517] have proposed a coupling method based on an idea of Nitsche. In the case of profile gratings with perfectly conducting substrate, the authors have shown previously that a slightly modified variational equation can be proven to be equivalent to the boundary value problem and to be uniquely solvable. Now it is shown that this result can be used to prove convergence for the FEM coupled by truncated wave mode expansion. This result covers transmission gratings and gratings bounded by additional multi-layer systems.

**Key words.** electromagnetic scattering, diffraction gratings, convergence analysis, finite element methods, mortar technique

**AMS subject classifications.** 78A45, 78M10, 65N30, 35J20

**1. Introduction.** The diffraction of light by biperiodic gratings, e.g., by doubly periodic surface structures, can be simulated by the time-harmonic Maxwell equations. Eliminating the magnetic field, the electric field is the solution of a boundary value problem for the time-harmonic curl-curl equation. For finite element methods (FEM), this problem is reduced to a finite domain, where quasi-periodic lateral boundary conditions and non-local boundary conditions over the upper and lower boundary face are required. The first idea for the solution of the boundary value problem is to express the non-local boundary conditions by integral operators and to couple FEM with boundary elements (cf. [10, 18]). With this approach, for the solution of the boundary value problem, either the case of wave modes propagating parallel to the surface is to be excluded or standard methods for integral operators with non-trivial null space are to be applied. As an alternative to integral operators, a saddle point type formulation (cf., e.g., [1]) or absorbing boundary conditions (cf., e.g., [24]) can be used.

On the other hand, the radiation conditions mean that the solutions can be extended in the form of a Rayleigh series expansion of upward respectively downward radiating Fourier modes. So the idea to couple finite elements and Rayleigh expansions is natural. Huber et al. [15] propose such a method, where the finite elements and the Rayleigh series are coupled employing a mortar technique by Nitsche (cf. [20, 27]). In [14], the case of perfectly conducting profile gratings has been considered and the coupling terms of [15] have been slightly modified. It has been proved that the variational equation for the coupling of FEM and Rayleigh expansions is equivalent to the boundary value problem for scattering by gratings. If the last problem is uniquely solvable, then the operator of the variational equation is uniquely solvable, too. In the references [3, 4, 5, 6] similar solvability results for all frequencies except for a countable set of Rayleigh frequencies were obtained in periodic chiral structures, and the coupling of finite element and integral equation methods was proposed and analyzed. For a general coupling of finite elements and boundary elements we also refer to [13].

\*Received January 21, 2013. Accepted June 6, 2014. Published online on October 14, 2014. Recommended by U. Langer.

<sup>†</sup>Weierstrass Institute for Applied Analysis and Stochastics, Mohrenstr. 39, 10117 Berlin, Germany ({guanghui.hu, andreas.rathsfeld}@wias-berlin.de).

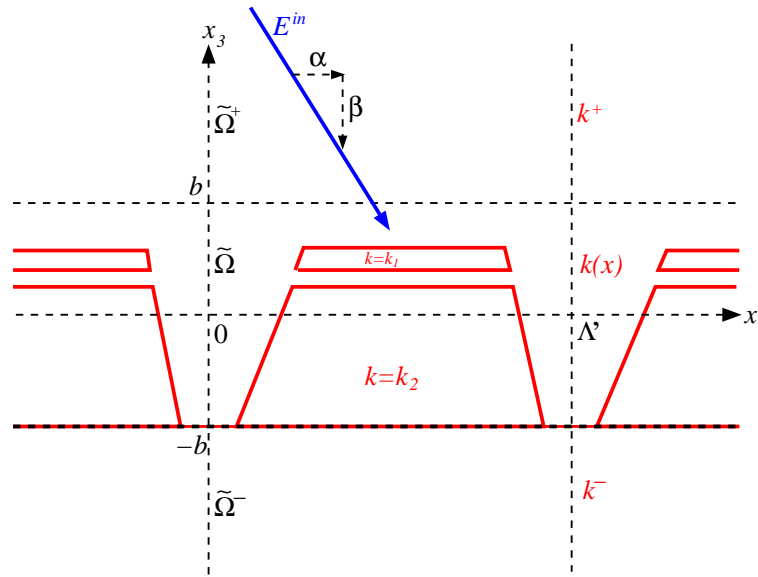


FIG. 3.1. Geometry of grating.

The subject of the present paper is to generalize the results on the variational formulation of [14] to the transmission problem, i.e., we analyze a formulation based on a mortar technique, which is a slight modification of the method proposed without proof in [15]. We show existence and uniqueness even in the case of resonance, where a Rayleigh frequency is allowed. Finally, we prove the convergence of the discretization scheme, i.e., for the coupling of the FEM (Nédélec’s edge elements) with truncated Rayleigh series expansions. Note that our mortar approach includes a natural coupling of Fourier modes with finite element functions and is easy to implement.

The plan of the paper is as follows. We formulate the boundary value problem and some solvability results in Section 3. In Section 4 we define the variational form and derive the Fredholm property for the operator corresponding to this form. The numerical discretization of the variational equation is introduced in Section 5. The stability and convergence of this method is proved. Of course, edge elements (cf., e.g., [17]) are employed for the FEM. In Section 6 we discuss the case of multi-layer systems beneath the grating structure. Instead of an extension of the FEM domain by the layers of the multi-layer system, we replace the down-going Fourier modes by special wave modes of the multi-layer system. Note that this idea goes back to the authors of [15]. The convergence analysis of Section 5 can be generalized to the multi-layer case too. Finally, we add a simple test showing that our method converges to the same solution as the 2D FEM for periodic 2D gratings and to the same solution as the method of [15].

**2. Preliminaries.** Throughout the paper, the symbols  $e_j$  ( $j = 1, 2, 3$ ) denote the unit coordinate vectors in the three-dimensional Cartesian coordinate system. The symbol  $(\cdot)^\top$  denotes the transpose of a vector in  $\mathbb{C}^2$  or  $\mathbb{C}^3$ , while the symbol  $\mathbf{a} \perp \mathbf{b}$  means the orthogonality of the vectors  $\mathbf{a} = (a_1, a_2, a_3)$ ,  $\mathbf{b} = (b_1, b_2, b_3) \in \mathbb{C}^3$  in the sense that  $\sum_{j=1}^3 a_j b_j = 0$ . Denote the unit sphere by  $\mathbb{S}^2 := \{x = (x_1, x_2, x_3)^\top \in \mathbb{R}^3 : \|x\| = 1\}$ , and define  $x' := (x_1, x_2)$  for  $x \in \mathbb{R}^3$ . The branch of the square root  $\sqrt{a}$  is chosen such that the imaginary part of  $\sqrt{a}$  is always positive, i.e.,  $\sqrt{a} = i\sqrt{-a}$  if  $a < 0$ .

**3. Diffraction problem.** Consider the scattering of a time-harmonic electromagnetic plane wave by a biperiodic structure (diffraction grating) which consists of at least two optical materials. By biperiodic or doubly periodic structure (cf. Fig. 3.1), we mean that the structure is periodic in two orthogonal directions  $x_1$  and  $x_2$  and bounded in  $x_3$ . The optical material inside the grating can be completely characterized by its dielectric coefficient and its magnetic permeability. For simplicity we assume that the medium is nonmagnetic with a constant magnetic permeability  $\mu(x) = \mu_0 > 0$  in  $\mathbb{R}^3$ . However, our arguments can be adapted to the case where  $\mu(x)$  is a periodic and piecewise constant function. The electric permittivity  $\epsilon(x)$  and the conductivity  $\sigma(x)$  are supposed to be  $\Lambda_j$ -periodic in  $x_j$  ( $j = 1, 2$ ) inside the grating and are homogeneous above and below the grating structure. More precisely, we assume that there exists a constant  $b > 0$  such that

$$\begin{aligned} \epsilon(x_1 + n_1\Lambda_1, x_2 + n_2\Lambda_2, x_3) &= \epsilon(x_1, x_2, x_3), \\ \sigma(x_1 + n_1\Lambda_1, x_2 + n_2\Lambda_2, x_3) &= \sigma(x_1, x_2, x_3), \end{aligned}$$

in  $\tilde{\Omega} := \{x : |x_3| < b\}$  for any  $n = (n_1, n_2) \in \mathbb{Z}^2$ , and

$$\begin{aligned} \epsilon(x) &= \epsilon_0^+ > 0, & \sigma(x) &= 0, & \text{in } x_3 > b, \\ \epsilon(x) &= \text{Re } \epsilon_0^-, & \sigma(x) &= \omega \text{Im } \epsilon_0^- > 0 & \text{in } x_3 < -b, \end{aligned}$$

with the circular frequency  $\omega > 0$ . Further, we restrict ourselves to the mostly used gratings, where  $\epsilon(x)$  and  $\sigma(x)$  are piecewise constant functions satisfying

$$(3.1) \quad 0 < \epsilon_0 < \epsilon(x) < \infty, \quad 0 \leq \sigma(x) < \infty \quad \text{in } \mathbb{R}^3.$$

Let  $\tilde{\Omega}^\pm := \{x : x_3 \gtrless \pm b\}$ . Suppose that a time-harmonic electromagnetic plane wave  $E^{in}(x)e^{-i\omega t}$  with  $E^{in}$  of the form

$$(3.2) \quad E^{in}(x) := q \exp(ik^+ x \cdot \hat{\theta}) = q \exp\left(i(x' \cdot \alpha - \beta x_3)\right), \quad i := \sqrt{-1}$$

is incident on the grating from  $\tilde{\Omega}^+$ . Here  $k^+ := \omega \sqrt{\epsilon_0^+ \mu_0}$  (respectively  $k^- := \omega \sqrt{\epsilon_0^- \mu_0}$ ) is defined as the wavenumber characterizing the homogenous medium in  $\tilde{\Omega}^+$  (respectively  $\tilde{\Omega}^-$ ). In (3.2), the symbol  $\hat{\theta}$  denotes the direction of incidence

$$\hat{\theta} := (\sin \theta_1 \cos \theta_2, \sin \theta_1 \sin \theta_2, -\cos \theta_1)^\top \in \mathbb{S}^2,$$

with the incident angles  $\theta_1 \in [0, \pi/2)$ ,  $\theta_2 \in [0, 2\pi)$ . Further, in (3.2), the three-dimensional vector  $q = (q_1, q_2, q_3)^\top \in \mathbb{S}^2$  stands for the direction of polarization satisfying  $q \perp \hat{\theta}$ , and

$$\alpha = (\alpha_1, \alpha_2)^\top := k(\sin \theta_1 \cos \theta_2, \sin \theta_1 \sin \theta_2)^\top \in \mathbb{R}^2, \quad \beta := k \cos \theta_1.$$

Eliminating the magnetic field from the reduced time-harmonic Maxwell's equations, we end up with the electric curl-curl equation

$$(3.3) \quad \text{curl curl } E(x) - k^2(x)E(x) = 0 \quad \text{for } x \in \mathbb{R}^3,$$

where  $k^2(x) := \omega^2 \mu_0 (\epsilon(x) + i\sigma(x)/\omega)$  and the electric field  $E$  in  $\tilde{\Omega}^+$  is the sum of the incident field  $E^{in}$  and the scattered field  $E^{sc}$ . The periodicity of the grating together with the form of  $E^{in}$  motivates us to look for  $\alpha$ -quasiperiodic solutions in the sense that  $E(x) \exp(-i\alpha \cdot x')$  is  $(\Lambda_1, \Lambda_2)$ -periodic in  $x'$ . In other words, it is required that

$$\begin{aligned} E(x_1 + \Lambda_1, x_2, x_3) &= \exp(i\Lambda_1 \alpha_1) E(x_1, x_2, x_3), \\ E(x_1, x_2 + \Lambda_2, x_3) &= \exp(i\Lambda_2 \alpha_2) E(x_1, x_2, x_3), \end{aligned}$$

for all  $x \in \mathbb{R}^3$ . Since the domain is unbounded in the  $x_3$ -direction, a radiation condition must be imposed. Noting that  $k(x) = k^\pm$  in  $\tilde{\Omega}^\pm$ , we suppose that the scattered field  $E^{sc}$  in  $\tilde{\Omega}^+$  and the electric field  $E$  in  $\tilde{\Omega}^-$  are composed of bounded outgoing plane waves in the form of

$$(3.4) \quad \begin{aligned} E^{sc}(x) &= \sum_{n \in \mathbb{Z}^2} E_n^+ \exp\left(i(\alpha_n \cdot x' + \beta_n^+ x_3)\right) \quad \text{for } x_3 > b, \quad E_n^+ \perp (\alpha_n, \beta_n^+)^T, \\ E(x) &= \sum_{n \in \mathbb{Z}^2} E_n^- \exp\left(i(\alpha_n \cdot x' - \beta_n^- x_3)\right) \quad \text{for } x_3 < -b, \quad E_n^- \perp (\alpha_n, -\beta_n^-)^T, \end{aligned}$$

where  $\alpha_n := (\alpha_n^{(1)}, \alpha_n^{(2)}) \in \mathbb{R}^2$ , with  $\alpha_n^{(j)} = \alpha_j + 2\pi n_j / \Lambda_j, j = 1, 2$ , for  $n = (n_1, n_2)^T \in \mathbb{Z}^2$ , and

$$\beta_n^\pm = \beta_n^\pm(k^\pm, \alpha) := \sqrt{(k^\pm)^2 - |\alpha_n|^2}.$$

We say that the scattered fields satisfy the radiation condition if expansions of the form (3.4) exist. These expansions are also referred to as the Rayleigh series expansions. The constant vectors  $E_n^\pm$  are called Rayleigh coefficients. Since  $\beta_n^\pm$  are real-valued only for finitely many indices  $n$ , we observe that only a finite number of wave modes in (3.4) propagate into the far field, while the remaining part consists of evanescent (or surface) waves decaying exponentially as  $x_3 \rightarrow \pm\infty$ . Thus, the above expansion for  $E^{sc}$  resp.  $E$  converges uniformly with all derivatives in the half space  $\{x_3 > a\}$  respectively  $\{x_3 < -a\}$  for any  $a > b$ .

Since the squared wave number  $k^2(x)$  is  $(\Lambda_1, \Lambda_2)$ -periodic in  $x'$  and both the incident and scattered fields are quasiperiodic, we can reduce the scattering problem to a single periodic cell. To this end, we introduce the following notation

$$\begin{aligned} \tilde{\Gamma}_b^\pm &:= \left\{ (x_1, x_2, x_3)^T \in \mathbb{R}^3: x_3 = \pm b \right\}, \\ \Gamma_b^\pm &:= \left\{ (x_1, x_2, x_3)^T \in \tilde{\Gamma}_b^\pm: 0 < x_j < \Lambda_j, j = 1, 2, \right\}, \\ \Omega^\pm &:= \left\{ (x_1, x_2, x_3)^T \in \tilde{\Omega}^\pm: 0 < x_j < \Lambda_j, j = 1, 2 \right\}, \\ \Omega &:= \left\{ x \in \tilde{\Omega}: 0 < x_j < \Lambda_j, j = 1, 2 \right\}. \end{aligned}$$

We next introduce some scalar and vector valued  $\alpha$ -quasiperiodic Sobolev spaces. Let  $H^s(\tilde{\Gamma}_b^\pm)$  be the complex-valued  $L^2$ -based Sobolev spaces of order  $s$  over  $\tilde{\Gamma}_b^\pm$ . Write

$$\begin{aligned} H_{loc}(\text{curl}, \tilde{\Omega}) &:= \left\{ G: \chi G, \text{curl}(\chi G) \in L^2(\tilde{\Omega})^3, \forall \chi \in C_0^\infty(\mathbb{R}^3) \right\}, \\ H_{loc}^s(\tilde{\Gamma}_b^\pm) &:= \left\{ G: \chi G \in H^s(\tilde{\Gamma}_b^\pm), \forall \chi \in C_0^\infty(\tilde{\Gamma}_b^\pm) \right\}, \\ H_{t,loc}^s(\tilde{\Gamma}_b^\pm) &:= \left\{ G \in H_{loc}^s(\tilde{\Gamma}_b^\pm): e_3 \cdot G = 0 \right\}, \\ H_{t,loc}^s(\text{Div}, \tilde{\Gamma}_b^\pm) &:= \left\{ G: G \in H_{t,loc}^s(\tilde{\Gamma}_b^\pm), \text{Div} G \in H_{t,loc}^s(\tilde{\Gamma}_b^\pm) \right\}, \\ H_{t,loc}^s(\text{Curl}, \tilde{\Gamma}_b^\pm) &:= \left\{ G: G \in H_{t,loc}^s(\tilde{\Gamma}_b^\pm), \text{Curl} G \in H_{t,loc}^s(\tilde{\Gamma}_b^\pm) \right\}, \end{aligned}$$

and

$$\begin{aligned}
 H(\text{curl}, \Omega) &:= \left\{ G|_{\Omega}: G \in H_{loc}(\text{curl}, \tilde{\Omega}), G \text{ is } \alpha\text{-quasiperiodic} \right\}, \\
 H_{qp}^s(\Omega) &:= \left\{ g|_{\Omega}: g \in H_{loc}^s(\tilde{\Omega}), g \text{ is } \alpha\text{-quasiperiodic} \right\}, \\
 H_t^s(\Gamma_b^{\pm}) &:= \left\{ G|_{\Gamma_b^{\pm}}: G \in H_{t,loc}^s(\tilde{\Gamma}_b^{\pm}), G \text{ is } \alpha\text{-quasiperiodic} \right\}, \\
 H_t^s(\text{Div}, \Gamma_b^{\pm}) &:= \left\{ G|_{\Gamma_b^{\pm}}: G \in H_{t,loc}^s(\text{Div}, \tilde{\Gamma}_b^{\pm}), G \text{ is } \alpha\text{-quasiperiodic} \right\}, \\
 H_t^s(\text{Curl}, \Gamma_b^{\pm}) &:= \left\{ G|_{\Gamma_b^{\pm}}: G \in H_{t,loc}^s(\text{Curl}, \tilde{\Gamma}_b^{\pm}), G \text{ is } \alpha\text{-quasiperiodic} \right\},
 \end{aligned}$$

where  $\text{Div}(\cdot)$  and  $\text{Curl}(\cdot)$  stand for the surface divergence and the surface scalar rotational operators, respectively. Note that, for  $x' \mapsto E(x', \pm b)$  in  $H_t^s(\Gamma_b^{\pm})$ ,  $s \in \mathbb{R}$ , we have the Fourier series expansion

$$\begin{aligned}
 E(x', \pm b) &= \sum_{n \in \mathbb{Z}^2} E_n^{\pm} \exp(i\alpha_n \cdot x'), \\
 E_n^{\pm} &:= (\Lambda_1 \Lambda_2)^{-1} \int_0^{\Lambda_1} \int_0^{\Lambda_2} E(x', \pm x_3) \exp(-i\alpha_n \cdot x') dx_1 dx_2 \in \mathbb{C}^3.
 \end{aligned}$$

Then, the spaces  $H_t^s(\Gamma_b^{\pm})$ ,  $H_t^s(\text{Div}, \Gamma_b^{\pm})$ , and  $H_t^s(\text{Curl}, \Gamma_b^{\pm})$  can be equipped with the following equivalent Sobolev norms

$$\begin{aligned}
 \|E\|_{H_t^s(\Gamma_b^{\pm})} &= \left( \sum_{n \in \mathbb{Z}^2} |E_n^{\pm}|^2 (1 + |\alpha_n|^2)^s \right)^{1/2}, \\
 \|E\|_{H_t^s(\text{Div}, \Gamma_b^{\pm})} &= \left( \sum_{n \in \mathbb{Z}^2} (1 + |\alpha_n|^2)^s (|E_n^{\pm}|^2 + |E_n^{\pm} \cdot (\alpha_n, 0)^{\top}|^2) \right)^{1/2}, \\
 \|E\|_{H_t^s(\text{Curl}, \Gamma_b^{\pm})} &= \left( \sum_{n \in \mathbb{Z}^2} (1 + |\alpha_n|^2)^s (|E_n^{\pm}|^2 + |E_n^{\pm} \times (\alpha_n, 0)^{\top}|^2) \right)^{1/2}.
 \end{aligned}$$

Recall that the space dual to  $H_t^s(\text{Div}, \Gamma_b^{\pm})$  with respect to the  $L^2$ -scalar product is  $H_t^s(\text{Div}, \Gamma_b^{\pm})' = H_t^{-s-1}(\text{Curl}, \Gamma_b^{\pm})$ , and that, for  $s = -1/2$ ,

$$\begin{aligned}
 H_t^{-1/2}(\text{Div}, \Gamma_b^{\pm}) &= \left\{ (e_3 \times E)|_{\Gamma_b^{\pm}}: E \in H(\text{curl}, \Omega) \right\}, \\
 H_t^{-1/2}(\text{Curl}, \Gamma_b^{\pm}) &= \left\{ (e_3 \times E)|_{\Gamma_b^{\pm}} \times e_3: E \in H(\text{curl}, \Omega) \right\}.
 \end{aligned}$$

Further, the trace mappings from  $H(\text{curl}, \Omega)$  to the tangential spaces  $H_t^{-1/2}(\text{Div}, \Gamma_b^{\pm})$  and  $H_t^{-1/2}(\text{Curl}, \Gamma_b^{\pm})$  are continuous and surjective (see [9, 17] and the references there). Finally, define our variational space

$$X = X_b := \left\{ E: \Omega \rightarrow \mathbb{C}^3: E \in H(\text{curl}, \Omega) \right\}$$

endowed with the norm

$$\|E\|_X := \|E\|_{H(\text{curl}, \Omega)} = \left( \|E\|_{L^2(\Omega)^3}^2 + \|\text{curl} E\|_{L^2(\Omega)^3}^2 \right)^{1/2}.$$

The boundary value problem for our scattering problem can be stated as follows.

(BVP): Given an incident electric field  $E^{in}$ , determine the quasiperiodic total electric field  $E \in H_{loc}(\text{curl}, \mathbb{R}^3)$  such that  $E(x)|_\Omega$  satisfies the curl-curl equation (3.3) in  $\Omega$  in the distributional sense and that the scattered field  $E^{sc} = E - E^{in}$  in  $x_3 > b$  as well as the transmitted field  $E$  in  $x_3 < -b$  admit a Rayleigh expansion of the form (3.4).

Introduce the set

$$(3.5) \quad \Upsilon_{\text{res}} := \Upsilon_{\text{res}}^+ \cup \Upsilon_{\text{res}}^-, \quad \Upsilon_{\text{res}}^\pm := \left\{ n \in \mathbb{Z}^2: \beta_n^\pm(k^\pm, \alpha) = 0 \right\}.$$

An incident angular frequency  $\omega$  with  $\Upsilon_{\text{res}} \neq \emptyset$  is called Rayleigh frequency. Note that the set  $\mathcal{F}$  of all Rayleigh frequencies depends on  $k^\pm$ ,  $\Lambda_1$ , and  $\Lambda_2$  but not on the shape of  $\Gamma$ .

Below we collect some uniqueness and existence results of (BVP) for a broad class of incident plane waves. Assume that the incident electric wave takes the form

$$(3.6) \quad E_{gen}^{in} := \sum_{n: \beta_n > 0} Q_n \exp(\alpha_n \cdot x' - \beta_n x_3),$$

where  $Q_n \in \mathbb{C}^3$  satisfies  $Q_n \perp (\alpha_n, -\beta_n)^\top$ . Note that  $E^{in}$  of (3.2) is of the form (3.6), where  $Q_n = q$ , for  $n = (0, 0)^\top$ , and  $Q_n = (0, 0, 0)^\top$  otherwise.

**THEOREM 3.1.** *Consider the scattering problem (BVP) with  $E^{in}$  replaced by  $E_{gen}^{in}$ .*

- (i) *There exists a unique solution to (BVP) for all  $\omega \in \mathbb{R}^+ \setminus \mathcal{D}$ , where  $\mathcal{D}$  is a discrete set with the only accumulating point at infinity.*
- (ii) *The problem (BVP) admits at least one solution for any  $\omega \in \mathbb{R}^+$ . Moreover, the far-field part of the solution scattered into the half space  $x_3 \geq \pm b$  is unique, i.e., the Rayleigh coefficients of the plane wave modes propagating into the half space  $x_3 \geq \pm b$  (namely, those  $E_n^\pm$  with  $\beta_n^\pm > 0$ ) are unique.*
- (iii) *There exists a small frequency  $\omega_0 > 0$  such that the problem (BVP) admits a unique solution for all  $\omega \in (0, \omega_0]$ .*

The assertions (i) and (ii) follow from the existence and uniqueness of the magnetic field in the space  $H^1(\Omega)^3$ ; see [8, 7, 11, 25, 26]. Note that the constant magnetic permeability implies the piecewise  $H^1$ -regularity of the magnetic field, which is not true for the electric field. In the non-resonance case (i.e.,  $\Upsilon_{\text{res}} = \emptyset$ ), (i) and (ii) can also be proved by studying the following variational formulation for the electric field  $E$  in  $\Omega$ : find  $E \in X$  such that

$$(3.7) \quad \int_\Omega [\text{curl } E \cdot \text{curl } \bar{\varphi} - k^2(x) E \cdot \bar{\varphi}] dx - \int_{\Gamma_b^+} \mathcal{R}^+(e_3 \times E) \cdot (e_3 \times \bar{\varphi}) ds \\ + \int_{\Gamma_b^-} \mathcal{R}^-(e_3 \times E) \cdot (e_3 \times \bar{\varphi}) ds \\ = \int_{\Gamma_b^+} [(\text{curl } E^{in})_T - \mathcal{R}^+(e_3 \times E^{in})] \cdot (e_3 \times \bar{\varphi}) ds,$$

for all  $\varphi \in X$ , where  $(\cdot)_T := [e_3 \times (\cdot)]|_{\Gamma_b^+} \times e_3$  and the operators

$$\mathcal{R}^\pm: H_t^{-1/2}(\text{Div}, \Gamma_b^\pm) \rightarrow H_t^{-1/2}(\text{Curl}, \Gamma_b^\pm)$$

are the Dirichlet-to-Neumann maps defined by

$$(\mathcal{R}^\pm \tilde{E})(x') = \mp \sum_{n \in \mathbb{Z}^2} \frac{1}{i\beta_n^\pm} \left[ k^2 \tilde{E}_n^\pm - (\alpha_n \cdot \tilde{E}_n^\pm) \alpha_n \right] \exp(i\alpha_n \cdot x'),$$

for  $\tilde{E}(x') = \sum_{n \in \mathbb{Z}^2} \tilde{E}_n^\pm \exp(i\alpha_n \cdot x') \in H_t^{-1/2}(\text{Div}, \Gamma_b^\pm)$ ,  $\tilde{E}_n^\pm \in \mathbb{C}^2$ ; see [1, 2]. Note that the operator  $\mathcal{R}^+$  maps  $e_3 \times E^{sc}$  to  $(\text{curl } E^{sc})_T$  on  $\Gamma_b^+$  and that  $\mathcal{R}^-$  maps  $-e_3 \times E$  to the trace  $(e_3 \times \text{curl } E) \times e_3$  on  $\Gamma_b^-$ . If the incident frequency  $\omega$  is sufficiently small, then the set  $\Upsilon_{\text{res}}$  is always empty and one can prove that the sesquilinear form generated by the left-hand side of (3.7) is positive coercive over  $X \times X$  under the assumption (3.1). We refer to [14, Lemma 6.1] for the proof of the third assertion for perfectly conducting grating profiles using a variational formulation analogously to (3.7) but posed only in the upper half space. These results can be easily extended to transmission gratings.

There are two drawbacks in using (3.7) to compute the electric field. First, the transparent boundary operators  $\mathcal{R}^\pm$  do not make sense if  $\beta_n^\pm = 0$  (i.e., in the resonance case). Thus, Rayleigh frequencies must be excluded. Second, in practice,  $\mathcal{R}^\pm$  cannot be computed straightforwardly from (3.7). Instead, they must be approximated by taking sufficiently many terms in the expansions; see [7, Section 6] for the error estimates. Motivated by the variational formulations proposed in [15, 23] and based on the mortar technique of Nitsche (see Nitsche [20] and Sternberg [27]), we employ a consistent coupling of the electric field  $E$  on the interfaces  $\Gamma_b^\pm$  as a replacement of the Dirichlet-to-Neumann maps. This way we propose a more general variational formulation than (3.7) for the electric field, which allows us not only to handle (BVP) in the resonance case but also to approximate the non-local boundary operators on  $\Gamma_b^\pm$ . Numerical experiments and convergence rate for a similar variational formulation were already reported in [15]. The goals of this paper are to provide a theoretical justification of the modified Nitsche’s method and to prove the convergence of its numerical discretization using Nédélec’s finite elements.

**4. Variational formulation based on a coupling method.** In this section we propose a variational formulation equivalent to (BVP). We begin with the fact that any column vector  $E_n^+ \in \mathbb{C}^3$  satisfying  $(\alpha_n, \beta_n^+)^\top \perp E_n^+$  for some  $n = (n_1, n_2)^\top \in \mathbb{Z}^2$  can be represented as a linear combination of two vectors  $E_{n,0}^+, E_{n,1}^+ \in \mathbb{C}^3$ :

$$E_n^+ = C_{n,0}^+ E_{n,0}^+ + C_{n,1}^+ E_{n,1}^+, \quad C_{n,0}^+, C_{n,1}^+ \in \mathbb{C},$$

where

$$E_{n,0}^+ := \begin{cases} (-\alpha_n^{(2)}, \alpha_n^{(1)}, 0)^\top / |\alpha_n| \in \mathbb{S}^2 & \text{if } |\alpha_n| \neq 0, \\ (0, 1, 0)^\top & \text{otherwise,} \end{cases}$$

$$E_{n,1}^+ := \begin{cases} \frac{|\alpha_n|}{h_n^+} (\alpha_n, \beta_n^+)^\top \times E_{n,0}^+ = (-\alpha_n^{(1)} \beta_n^+, -\alpha_n^{(2)} \beta_n^+, |\alpha_n|^2)^\top / h_n^+ & \text{if } |\alpha_n| \neq 0, \\ (-1, 0, 0)^\top & \text{otherwise,} \end{cases}$$

with  $h_n^+ := |\alpha_n| \sqrt{|\alpha_n|^2 + |\beta_n^+|^2}$ . Obviously, it holds that  $(\alpha_n, \beta_n^+)^\top \perp E_{n,l}^+$ ,  $|E_{n,l}^+| = 1$ , for  $l = 0, 1, n \in \mathbb{Z}^2$ . One can observe further that  $E_{n,1}^+ \in \mathbb{S}^2$  if  $\beta_n^+ \in \mathbb{R}$ , and that  $E_{n,1}^+ = e_3$  if  $\beta_n^+ = 0$ . The above decomposition of  $E_n^+$  allows us to rewrite the Rayleigh expansion (3.4) for  $E^{sc}$  as (see also [23, Section 2.5])

$$E^{sc}(x) = \sum_{n \in \mathbb{Z}^2, l=1,2} C_{n,l}^+ U_{n,l}^+(x), \quad U_{n,l}^+ := E_{n,l}^+ \exp\left(i[\alpha_n \cdot x' + \beta_n^+ x_3]\right), \quad C_{n,l}^+ \in \mathbb{C},$$

for  $x_3 > b$ . Analogously, there holds

$$E(x) = \sum_{n \in \mathbb{Z}^2, l=1,2} C_{n,l}^- U_{n,l}^-(x), \quad U_{n,l}^- := E_{n,l}^- \exp\left(i[\alpha_n \cdot x' - \beta_n^- x_3]\right), \quad C_{n,l}^- \in \mathbb{C},$$

for  $x < -b$ , where  $E_{n,0}^- = E_{n,0}^+$  and

$$E_{n,1}^- := \begin{cases} \frac{|\alpha_n|}{h_n^-} (\alpha_n, -\beta_n^-)^\top \times E_{n,0}^- = (\alpha_n^{(1)} \beta_n^-, \alpha_n^{(2)} \beta_n^-, |\alpha_n|^2)^\top / h_n^- & \text{if } |\alpha_n| \neq 0, \\ (-1, 0, 0)^\top & \text{otherwise,} \end{cases}$$

with  $h_n^- := |\alpha_n| \sqrt{|\alpha_n|^2 + |\beta_n^-|^2}$ . Define the layers  $D^\pm$  of height one above  $\Gamma_b^+$  and below  $\Gamma_b^-$  by

$$\begin{aligned} D^+ &:= \{x \in \mathbb{R}^3: 0 < x_j < \Lambda_j, j = 1, 2, b < x_3 < b + 1\}, \\ D^- &:= \{x \in \mathbb{R}^3: 0 < x_j < \Lambda_j, j = 1, 2, -b - 1 < x_3 < -b\}. \end{aligned}$$

Now we introduce the Sobolev spaces  $Y_l^\pm$  as follows:

$$(4.1) \quad Y_l^\pm := \left\{ U \in H(\text{curl}, D^\pm): U(x) = \sum_{n \in \mathbb{Z}^2} C_{n,l}^\pm U_{n,l}^\pm(x), C_{n,l}^\pm \in \mathbb{C} \right\}, l = 0, 1.$$

Then we see that the function  $E^+(x) := E^{sc}|_{D^+}$  belongs to the space  $Y^+ := Y_0^+ \oplus Y_1^+$ , and that  $E^-(x) := E|_{D^-}$  belongs to the space  $Y^- := Y_0^- \oplus Y_1^-$ . Hence, the following problem is equivalent to (BVP):

(BVP'): Given an incident electric field  $E^{in}$ , find the  $\alpha$ -quasiperiodic fields  $(E, E^+, E^-)$  in  $\mathbb{H} := X \times Y^+ \times Y^-$  such that  $E$  satisfies the curl-curl equation (3.3) in  $\Omega$  in a distributional sense and the transmission conditions

$$\begin{aligned} e_3 \times (E - E^{in} - E^+) &= 0, & e_3 \times \text{curl}(E - E^{in} - E^+) &= 0 & \text{on } \Gamma_b^+, \\ e_3 \times (E - E^-) &= 0, & e_3 \times \text{curl}(E - E^-) &= 0 & \text{on } \Gamma_b^-. \end{aligned}$$

Motivated by the arguments in [23, Section 3.2] and the variational formulations in [14, 15], we propose a new variational formulation that is equivalent to (BVP'). For the triples of functions  $(E, E^+, E^-)$ ,  $(V, V^+, V^-) \in \mathbb{H}$ , define the sesquilinear form  $a(\cdot, \cdot): \mathbb{H} \times \mathbb{H} \rightarrow \mathbb{C}$  by

$$\begin{aligned} &a\left((E, E^+, E^-), (V, V^+, V^-)\right) \\ &:= \int_{\Omega} \left\{ \text{curl } E \cdot \text{curl } \bar{V} - k^2(x) E \cdot \bar{V} \right\} dx \\ &\quad - \int_{\Gamma_b^+} \left\{ \text{curl } E^+ \cdot e_3 \times \bar{V} - e_3 \times (E - E^+) \cdot \text{curl } \bar{V}^+ \right\} ds \\ (4.2) \quad &+ \int_{\Gamma_b^-} \left\{ \text{curl } E^- \cdot e_3 \times \bar{V} - e_3 \times (E - E^-) \cdot \text{curl } \bar{V}^- \right\} ds \\ &\quad - \eta^+ \sum_{n \in \Upsilon^+} \left[ \int_{\Gamma_b^+} e_3 \times (E - E^+) \cdot (e_3 \times \bar{U}_{n,0}^+) ds \int_{\Gamma_b^+} e_3 \times V^+ \cdot (e_3 \times \bar{U}_{n,0}^+) ds \right] \\ &\quad - \eta^- \sum_{n \in \Upsilon^-} \left[ \int_{\Gamma_b^-} e_3 \times (E - E^-) \cdot (e_3 \times \bar{U}_{n,0}^-) ds \int_{\Gamma_b^-} e_3 \times V^- \cdot (e_3 \times \bar{U}_{n,0}^-) ds \right], \end{aligned}$$

where  $\eta^\pm > 0$  are constant factors. The set  $\Upsilon^\pm$  is a finite fixed subset of  $\mathbb{Z}^2$  with  $\Upsilon_{\text{res}}^\pm \subseteq \Upsilon^\pm$  (cf. (3.5)). Our variational formulation is to find  $(E, E^+, E^-) \in \mathbb{H}$  such that

$$(4.3) \quad a\left((E, E^+, E^-), (V, V^+, V^-)\right) = -a\left((0, E^{in}, 0), (V, V^+, V^-)\right)$$



for all  $(V, V^+, V^-) \in \mathbb{H}$ . Note that terms like  $\int_{\Gamma_b^\pm} \text{curl } E^\pm \cdot e_3 \times \bar{V} \, ds$  are bounded. Indeed, since  $E^\pm$  is the solution of the curl-curl equation in  $D^\pm$ , we get  $\text{curl } E^\pm \in H(\text{curl}, D^\pm)$  and  $(\text{curl } E^\pm)|_{\Gamma_b^\pm} \in H^{-1/2}(\text{Curl}, \Gamma_b^\pm)$ . Further, note that the second part of the second and third terms on the right-hand side of (4.2) both have opposite signs than the corresponding terms in [15]. Moreover, the integrals with the factor  $\eta^\pm$  in (4.2) are modifications of the following terms involved in the variational equation of [15]:

$$(4.4) \quad \eta^\pm \int_{\Gamma_b^\pm} e_3 \times (E - E^\pm) \cdot e_3 \times \overline{(V - V^\pm)} \, ds.$$

The expressions in (4.4) are not meaningful for general  $(E, E^+, E^-), (V, V^+, V^-) \in \mathbb{H}$  since both  $e_3 \times (E - E^\pm)$  and  $e_3 \times \overline{(V - V^\pm)}$  belong to the space  $H_t^{-1/2}(\text{Div}, \Gamma_b^\pm)$ . Integrals like  $\eta \int_{\Gamma_b^\pm} e_3 \times u \cdot e_3 \times \bar{v} \, ds$  in the mortar approach make sense for finite element methods, where  $u$  and  $v$  are finite element functions and  $\eta$  tends to zero with the mesh size. The idea employed in [23] is to replace the integral (4.4) by the Galerkin approximation

$$(4.5) \quad \sum_{\substack{n,l:|n|^2 < N \\ \beta_n^\pm \neq 0 \text{ or } l=0}} \left[ \eta^\pm \int_{\Gamma_b^\pm} e_3 \times (E - E^\pm) \cdot e_3 \times \bar{U}_{n,l}^\pm \, ds \int_{\Gamma_b^\pm} e_3 \times (V - V^\pm) \cdot e_3 \times \bar{U}_{n,l}^\pm \, ds \right]$$

$$(4.6) \quad + \eta^\pm \sum_{n:\beta_n^\pm=0} \left[ \int_{\Gamma_b^\pm} e_3 \times (E - E^\pm) \cdot \bar{U}_{n,0}^\pm \, ds \int_{\Gamma_b^\pm} e_3 \times (V - V^\pm) \cdot \bar{U}_{n,0}^\pm \, ds \right]$$

with a sufficiently large number  $N > 0$ . It is also mentioned in [23] that the summation in (4.5) and (4.6) can even be restricted to all  $n \in \mathbb{Z}^2$  with  $\beta_n^\pm = 0$ . In the present paper, we only use the terms of (4.5) with  $n \in \Upsilon^\pm$  and simplify them to get the last two terms in (4.2). Note that choosing  $\Upsilon^\pm$  larger than  $\Upsilon_{\text{res}}^\pm$  makes the numerical scheme more stable in the near-resonance case.

Arguing similarly to [14, Lemma 3.3], we can prove the equivalence of the variational formulation (4.3) and the problem (BVP'). Moreover, in the non-resonance case, i.e.,  $\Upsilon_{\text{res}} = \emptyset$ , and for  $\Upsilon = \Upsilon_{\text{res}}$ , the variational formulations (4.3) and (3.7) are equivalent; see [14, Remark 3.4]. Thus, the variational formulation (4.3) is indeed more general than (3.7). It is worth to mention that, using (4.3), we can also prove the solvability results in Theorem 3.1 since the arguments in [14] for perfectly conducting grating profiles can be easily adapted to transmission gratings. To prepare the convergence analysis of the finite element discretization, in this paper we only check the Fredholm property of the operator  $A : \mathbb{H} \rightarrow \mathbb{H}'$  generated by the bounded sesquilinear form  $a(\cdot, \cdot)$  defined in Section 4, i.e.,  $A$  is given by

$$(4.7) \quad a\left((E, E^+, E^-), (V, V^+, V^-)\right) = \left\langle A(E, E^+, E^-), (V, V^+, V^-) \right\rangle.$$

Here  $\mathbb{H}'$  denotes the space dual to  $\mathbb{H}$  with respect to the duality  $\langle \cdot, \cdot \rangle$  extending the scalar product in  $L^2(\Omega)^3 \times L^2(D^+)^3 \times L^2(D^-)^3$ . The rest of this section is devoted to verify the following theorem.

**THEOREM 4.1.** *The operator  $A$  defined by (4.7) is a Fredholm operator with index zero. First we recall the following definition.*

**DEFINITION 4.2.** *A bounded sesquilinear form  $l(\cdot, \cdot)$  given on some Hilbert space  $Y$  is called strongly elliptic if there exists a compact form  $\tilde{l}(\cdot, \cdot)$  and a constant  $c > 0$  such that*

$$\text{Re } l(u, u) \geq c \|u\|_Y^2 - \tilde{l}(u, u), \quad \forall u \in Y.$$

To prove Theorem 4.1, we need a periodic analogue of the Hodge decomposition of  $X$ .

LEMMA 4.3.

(i) We have  $X = X_0 \oplus X_1$ , where

$$X_1 := \left\{ \nabla p : p \in H_{qp}^1(\Omega) \right\} \subset X,$$

$$X_0 := \left\{ E_0 \in X : \int_{\Omega} k^2(x) \nabla p \cdot \bar{E}_0 \, dx = 0 \text{ for all } \nabla p \in X_1 \right\}$$

and the space  $X_0$  is compactly embedded into  $L^2(\Omega)^3$ .

(ii) We have  $\operatorname{div}(k^2(x)E_0) = 0$  in  $\Omega$  and  $e_3 \cdot E_0 = 0$  on  $\Gamma_b^{\pm}$  for any  $E_0 \in X_0$ .

*Proof.* See, e.g., [3, Section 3.1] for the proof of the first assertion in more general periodic chiral structures and [17, Section 4.4] in the case of non-periodic structures where  $k^2(x)$  is allowed to be a complex-valued function. Using integration by parts, it follows from the definition of  $X_0$  that  $\operatorname{div}(k^2(x)E_0) = 0$  in  $\Omega$  and  $e_3 \cdot k^2(x)E_0 = 0$  on  $\Gamma_b^{\pm}$ . Since  $k^2(x)$  is a non-vanishing piecewise constant function in  $\Omega$ , we obtain  $e_3 \cdot E_0 = 0$  on  $\Gamma_b^{\pm}$ .  $\square$

By Lemma 4.3 and the definitions of  $Y_l^{\pm}$ , we can decompose our space  $\mathbb{H}$  into six subspaces

$$\mathbb{H} = (X_0 \oplus X_1) \times (Y_0^+ \oplus Y_1^+) \times (Y_0^- \oplus Y_1^-).$$

For  $(E, E^+, E^-), (V, V^+, V^-) \in \mathbb{H}$ , we may assume that

$$E = \nabla p + E_0, \quad E^{\pm} = E_0^{\pm} + E_1^{\pm}, \quad \text{where } \nabla p \in X_1, E_0 \in X_0, E_l^{\pm} \in Y_l^{\pm}, l=0,1,$$

$$V = \nabla \xi + V_0, \quad V^{\pm} = V_0^{\pm} + V_1^{\pm}, \quad \text{where } \nabla \xi \in X_1, V_0 \in X_0, V_l^{\pm} \in Y_l^{\pm}, l=0,1.$$

For the analysis of the form  $a$ , we define several sesquilinear forms as follows. Let

$$a_1(\nabla p, \nabla \xi) := \int_{\Omega} k^2(x) \nabla p \cdot \nabla \bar{\xi} \, dx, \quad \forall \nabla p, \nabla \xi \in X_1,$$

$$a_2(E_0, V_0) := \int_{\Omega} \{ \operatorname{curl} E_0 \cdot \operatorname{curl} \bar{V}_0 - k^2(x) E_0 \cdot \bar{V}_0 \} \, dx, \quad \forall E_0, V_0 \in X_0,$$

$$a_3^{\pm}(E_0^{\pm}, V_0^{\pm}) := \pm \int_{\Gamma_b^{\pm}} e_3 \times E_0^{\pm} \cdot \operatorname{curl} \bar{V}_0^{\pm} \, ds, \quad \forall E_0^{\pm}, V_0^{\pm} \in Y_0^{\pm},$$

$$a_4^{\pm}(E_1^{\pm}, V_1^{\pm}) := \pm \int_{\Gamma_b^{\pm}} e_3 \times E_1^{\pm} \cdot \operatorname{curl} \bar{V}_1^{\pm} \, ds, \quad \forall E_1^{\pm}, V_1^{\pm} \in Y_1^{\pm},$$

and let

$$a_5^{\pm} \left( (E, E^+, E^-), (V, V^+, V^-) \right) := \pm \int_{\Gamma_b^{\pm}} e_3 \times E \cdot \operatorname{curl} \bar{V}^{\pm} \, ds,$$

$$a_6^{\pm} \left( (E, E^+, E^-), (V, V^+, V^-) \right) := -\eta^{\pm} \sum_{n \in \Upsilon^{\pm}} \left\{ \int_{\Gamma_b^{\pm}} e_3 \times (E - E^{\pm}) \cdot (e_3 \times \bar{U}_{n,0}^{\pm}) \, ds \overline{\int_{\Gamma_b^{\pm}} (e_3 \times V^{\pm}) \cdot (e_3 \times \bar{U}_{n,0}^{\pm}) \, ds} \right\},$$

for any  $(E, E^+, E^-), (V, V^+, V^-) \in \mathbb{H}$ . For brevity we write

$$a_5^{\pm} \left( (E, E^+, E^-), (V, V^+, V^-) \right) = a_5^{\pm}(E, V^{\pm}), \quad \forall E \in X, V^{\pm} \in Y^{\pm}.$$

LEMMA 4.4. For any  $\nabla\xi \in X_1$  and  $V_0^\pm \in Y_0^\pm$ , we have  $a_5^\pm(\nabla\xi, V_0^\pm) = 0$ .

*Proof.* Assume that  $\nabla\xi \in X_1$  and  $V_0^+ \in Y_0^+$ . Without loss of generality,  $\xi$  can be assumed to be smooth. We can expand the function  $\xi(x)$  into the series

$$\xi(x) = \sum_{n \in \mathbb{Z}^2} f_n(x_3) \exp(i\alpha_n \cdot x'), \quad f_n \in C^2(\mathbb{R}^+),$$

in a sufficiently small neighborhood of  $\Gamma_b^+$ . This implies that

$$(e_3 \times \nabla\xi)|_{\Gamma_b^+} = \sum_{n \in \mathbb{Z}^2} i f_n(b) (-\alpha_n^{(2)}, \alpha_n^{(1)}, 0)^\top \exp(i\alpha_n \cdot x').$$

Making use of  $\operatorname{curl} U_{n,0}^+ = i U_{n,1}^+ \sqrt{|\alpha_n|^2 + |\beta_n^+|^2}$  (see [14, Lemma 3.1]), and recalling the definition of  $U_{n,1}^+$  and the sesquilinear form  $a_5^+$ , we end up with the identity

$$a_5^+(\nabla\xi, V_0^+) = \int_{\Gamma_b^+} (e_3 \times \nabla\xi) \cdot \operatorname{curl} \overline{V_0^+} ds = 0.$$

The proof for  $a_5^-$  can be carried out analogously.  $\square$

Note that the last proof is a new and simpler proof of [14, Lemma 4.3]. Using Lemmas 4.3 and 4.4 and the definition of  $a$ , a simple calculation implies (see Table 4.1)

$$\begin{aligned} & a\left((E, E^+, E^-), (V, V^+, V^-)\right) \\ &= a\left((\nabla p + E_0, E_0^+ + E_1^+, E_0^- + E_1^-), (\nabla\xi + V_0, V_0^+ + V_1^+, V_0^- + V_1^-)\right) \\ &= -a_1(\nabla p, \nabla\xi) + a_2(E_0, V_0) - a_3^+(E_0^+, V_0^+) - a_4^+(E_1^+, V_1^+) + a_5^+(E_0, V_0^+) \\ &\quad - \overline{a_5^+(V_0, E_0^+)} + a_5^+(E_0, V_1^+) - \overline{a_5^+(V_0, E_1^+)} + a_5^+(\nabla p, V_1^+) - \overline{a_5^+(\nabla\xi, E_1^+)} \\ &\quad + a_6^+\left((E, E^+, E^-), (V, V^+, V^-)\right) + a_3^-(E_0^-, V_0^-) + a_4^-(E_1^-, V_1^-) - a_5^-(E_0, V_0^-) \\ &\quad + \overline{a_5^-(V_0, E_0^-)} - a_5^-(E_0, V_1^-) + \overline{a_5^-(V_0, E_1^-)} - a_5^-(\nabla p, V_1^-) + \overline{a_5^-(\nabla\xi, E_1^-)} \\ &\quad + a_6^-\left((E, E^+, E^-), (V, V^+, V^-)\right). \end{aligned}$$

*Proof of Theorem 4.1.* Obviously, we have

- $a_1$  is coercive on  $X_1$ , i.e., there exists some constant  $C > 0$  such that

$$\operatorname{Re} [a_1(\nabla p, \nabla p)] \geq C \|\nabla p\|_X, \quad \forall \nabla p \in X_1.$$

- $a_2$  is strongly elliptic over  $X_0$  due to the estimate

$$\operatorname{Re} [a_2(E_0, E_0)] \geq \|E_0\|_X - [1 + \|k^2\|_{L^\infty(\Omega)}] \|E_0\|_{L^2(\Omega)^3}^2,$$

for any  $E_0 \in X_0$  and the compact imbedding of  $X_0$  into  $L^2(\Omega)^3$  (see Lemma 4.3).

- $a_6^\pm$  are compact forms over  $\mathbb{H}$  since each of them corresponds to a finite rank operator over  $\mathbb{H}$ .

To demonstrate the Fredholm property of the sesquilinear form  $a$ , we now need to study the other forms  $a_3^\pm$ ,  $a_4^\pm$ , and  $a_5^\pm$ . Concerning  $a_3^+$  and  $a_4^+$ , it is shown in [14, Lemma 4.5] that there exist compact forms  $\tilde{a}_3^+ : Y_0^+ \times Y_0^+ \rightarrow \mathbb{C}$  and  $\tilde{a}_4^+ : Y_1^+ \times Y_1^+ \rightarrow \mathbb{C}$  such that

$$(4.8) \quad \begin{aligned} -\operatorname{Re} a_3^+(E_0^+, E_0^+) &\geq C_3^+ \|E_0^+\|_{H(\operatorname{curl}, D^+)}^2 - \tilde{a}_3^+(E_0^+, E_0^+), \quad \forall E_0^+ \in Y_0^+, \\ \operatorname{Re} a_4^+(E_1^+, E_1^+) &\geq C_4^+ \|E_1^+\|_{H(\operatorname{curl}, D^+)}^2 - \tilde{a}_4^+(E_1^+, E_1^+), \quad \forall E_1^+ \in Y_1^+, \end{aligned}$$

for some constants  $C_3^+, C_4^+ > 0$ , i.e., the sesquilinear forms  $-a_3^+$  and  $a_4^+$  are strongly elliptic over  $Y_0^+$  and  $Y_1^+$ , respectively. The proof of the estimates in (4.8) can be easily extended to the sesquilinear forms  $a_3^-$  and  $a_4^-$ . That is, we can find compact forms  $\tilde{a}_3^- : Y_0^- \times Y_0^- \rightarrow \mathbb{C}$  and  $\tilde{a}_4^- : Y_1^- \times Y_1^- \rightarrow \mathbb{C}$  such that

$$(4.9) \quad \begin{aligned} \operatorname{Re} a_3^-(E_0^-, E_0^-) &\geq C_3^- \|E_0^-\|_{H(\operatorname{curl}, D^-)}^2 - \tilde{a}_3^-(E_0^-, E_0^-), \quad \forall E_0^- \in Y_0^-, \\ -\operatorname{Re} a_4^-(E_1^-, E_1^-) &\geq C_4^- \|E_1^-\|_{H(\operatorname{curl}, D^-)}^2 - \tilde{a}_4^-(E_1^-, E_1^-), \quad \forall E_1^- \in Y_1^-, \end{aligned}$$

for some constants  $C_3^-, C_4^- > 0$ . Hence the strong ellipticity of  $a_3^-$  and  $-a_4^-$  follows. Finally, in view of [14, Lemma 4.7] we have

- $a_5^+$  is compact over  $X_0 \times Y_1^+$ ,

and analogously

- $a_5^-$  is compact over  $X_0 \times Y_1^-$ .

To prove the Fredholm property of the variational formulation (4.3), it suffices to verify that the operator corresponding to the sesquilinear form  $a - a_6^+ - a_6^-$  is Fredholm over  $\mathbb{H}$  with index zero. For this purpose, we define the spaces  $\mathbb{H}_j = X_j \times Y_j^+ \times Y_j^-$  for  $j = 0, 1$ , so that we can rewrite  $\mathbb{H} = X \times Y^+ \times Y^- = \mathbb{H}_0 \oplus \mathbb{H}_1$ . Define the sesquilinear forms

$$\begin{aligned} b_0 &\left( (E_0, E_0^+, E_0^-), (V_0, V_0^+, V_0^-) \right) \\ &:= a_2(E_0, V_0) - a_3^+(E_0^+, V_0^+) + a_3^-(E_0^-, V_0^-) \\ &\quad + \overline{a_5^+(E_0, V_0^+) - a_5^+(V_0, E_0^+) - a_5^-(E_0, V_0^-) + a_5^-(V_0, E_0^-)}, \end{aligned}$$

for all  $(E_0, E_0^+, E_0^-), (V_0, V_0^+, V_0^-) \in \mathbb{H}_0$ , and

$$\begin{aligned} b_1 &\left( (\nabla p, E_1^+, E_1^-), (\nabla \xi, V_1^+, V_1^-) \right) \\ &:= -a_1(\nabla p, \nabla \xi) - a_4^+(E_1^+, V_1^+) + a_4^-(E_1^-, V_1^-) \\ &\quad + \overline{a_5^+(\nabla p, V_1^+) - a_5^+(\nabla \xi, E_1^+) - a_5^-(\nabla p, V_1^-) + a_5^-(\nabla \xi, E_1^-)}, \end{aligned}$$

for all  $(\nabla p, E_1^+, E_1^-), (\nabla \xi, V_1^+, V_1^-) \in \mathbb{H}_1$ . Now split the form in Table 4.1 in blocks corresponding to the splitting  $\mathbb{H} = \mathbb{H}_1 \times \mathbb{H}_2$ . Then the restriction to  $\mathbb{H}_1$  is the form  $b_0$  with the strongly elliptic quadratic form

$$\begin{aligned} \operatorname{Re} b_0 &\left( (E_0, E_0^+, E_0^-), (E_0, E_0^+, E_0^-) \right) \\ &= \operatorname{Re} a_2(E_0, E_0) - \operatorname{Re} a_3^+(E_0^+, E_0^+) + \operatorname{Re} a_3^-(E_0^-, E_0^-). \end{aligned}$$

The restriction to  $\mathbb{H}_1$  is the form  $b_1$ , and  $-b_1$  has the strongly elliptic quadratic form

$$\begin{aligned} -\operatorname{Re} b_1 &\left( (\nabla p, E_1^+, E_1^-), (\nabla p, E_1^+, E_1^-) \right) \\ &= \operatorname{Re} a_1(\nabla p, \nabla p) + \operatorname{Re} a_4^+(E_1^+, E_1^+) - \operatorname{Re} a_4^-(E_1^-, E_1^-). \end{aligned}$$

Consequently, the diagonal blocks of the splitting into  $2 \times 2$  blocks of size  $3 \times 3$  correspond to Fredholm operators with index zero. On the other hand, the full form in Table 4.1 differs from the diagonal block matrix only by compact terms. Hence the form  $a$  generates a Fredholm operator with index zero.  $\square$

TABLE 4.1  
 The diagram for the sesquilinear form  $a - a_6^+ - a_6^-$  over  $\mathbb{H} \times \mathbb{H}$ , where  $\mathbb{H} = X \times Y^+ \times Y^-$ .

		$\mathbb{H}_0 := X_0 \times Y_0^+ \times Y_0^-$		
		$X_0(E_0)$	$Y_0^+(E_0^+)$	$Y_0^-(E_0^-)$
$\mathbb{H}_0$	$X_0(V_0)$	$a_2(E_0, V_0)$	$-\overline{a_5^+(V_0, E_0^+)}$	$\overline{a_5^-(V_0, E_0^-)}$
	$Y_0^+(V_0^+)$	$a_5^+(E_0, V_0^+)$	$-a_3^+(E_0^+, V_0^+)$	0
	$Y_0^-(V_0^-)$	$-a_5^-(E_0, V_0^-)$	0	$a_3^-(E_0^-, V_0^-)$
$\mathbb{H}_1$	$X_1(\nabla\xi)$	0	0	0
	$Y_1^+(V_1^+)$	$a_5^+(E_0, V_1^+)$	0	0
	$Y_1^-(V_1^-)$	$-a_5^-(E_0, V_1^-)$	0	0

		$\mathbb{H}_1 := X_1 \times Y_1^+ \times Y_1^-$		
		$X_1(\nabla p)$	$Y_1^+(E_1^+)$	$Y_1^-(E_1^-)$
$\mathbb{H}_0$	$X_0(V_0)$	0	$-\overline{a_5^+(V_0, E_1^+)}$	$\overline{a_5^-(V_0, E_1^-)}$
	$Y_0^+(V_0^+)$	0	0	0
	$Y_0^-(V_0^-)$	0	0	0
$\mathbb{H}_1$	$X_1(\nabla\xi)$	$-a_1(\nabla p, \nabla\xi)$	$-\overline{a_5^+(\nabla\xi, E_1^+)}$	$\overline{a_5^-(\nabla\xi, E_1^-)}$
	$Y_1^+(V_1^+)$	$a_5^+(\nabla p, V_1^+)$	$-a_4^+(E_1^+, V_1^+)$	0
	$Y_1^-(V_1^-)$	$-a_5^-(\nabla p, V_1^-)$	0	$a_4^-(E_1^-, V_1^-)$

**5. Numerical analysis of the Finite Element Method.**

**5.1. Finite element spaces and the FEM.** As mentioned in the introduction, we assume that the optical medium in  $\mathbb{R}^3$  is piecewise smooth. For the convergence analysis, we suppose that the interface between any two different materials is a polyhedral surface. Let  $\tau_h = \tau_h(\Omega)$  be a partition of  $\overline{\Omega}$  by tetrahedrons  $K$  of diameter  $h_K$ , i.e.,  $\overline{\Omega} = \cup_{K \in \tau_h} \overline{K}$ , where

$h$  denotes the maximum diameter of the elements in  $\tau_h$ . Of course, we suppose that  $\epsilon$  and  $k$  are constant over each  $K \in \tau_h$ . We will use standard Nédélec's edge elements (cf. [17]) and analyze convergence for  $h \rightarrow 0$ . For each element  $K \in \tau_h$  and  $k > 1$ , denote by  $P_k$  the polynomials of maximal total degree  $k$  and by  $\tilde{P}_k$  the homogeneous polynomials of total degree  $k$ . Define the subspace  $S_k$  of homogeneous vector polynomials of degree  $k$  by  $S_k := \{\mathbf{p} \in (\tilde{P}_k)^3 \mid x \cdot \mathbf{p}(x) = 0\}$ . The curl conforming edge elements of Nédélec rely on the use of the vector polynomial space  $R_K := (P_{k-1})^3 \oplus S_k$ . More precisely, the Nédélec finite element space of edge elements of degree  $k$  are defined as follows.

DEFINITION 5.1. Let  $X_h \subset X$  be the set of functions  $E_h : \Omega \rightarrow \mathbb{C}^3$  such that:

- (i) For any  $K \in \tau_h$ , we have  $E_h|_K \in R_K$ .
- (ii) For any edge  $e$  of the FE partition and for any  $K, K' \in \tau_h$  s.t.  $e \subseteq \overline{K} \cap \overline{K}'$ , we have  $\int_e (E_h|_K) \cdot \tau \, q \, de = \int_e (E_h|_{K'}) \cdot \tau \, q \, de$  for any  $q \in P_{k-1}$ . Here,  $\tau$  is the unit vector pointing into the direction of  $e$ .
- (iii) For any face  $f$  of the FE partition and for any  $K, K' \in \tau_h$  such that  $f \subseteq K \cap K'$ , there holds  $\int_f (E_h|_K) \cdot \mathbf{q} \, ds = \int_f (E_h|_{K'}) \cdot \mathbf{q} \, ds$  for any  $\mathbf{q} \in (P_{k-2})^3$  with  $\mathbf{q} \cdot \nu_f = 0$ . Here,  $\nu_f$  denotes the normal to the face  $f$ .

To define the discretized spaces for  $Y_l^\pm$ , for some constant  $C > 0$ , we introduce the finite set  $\Upsilon_h := \{n \in \mathbb{Z}^2 : |n| \leq C/h\}$ . Then set

$$Y_h^\pm := Y_{h,0}^\pm \oplus Y_{h,1}^\pm, \quad Y_{h,l}^\pm := \text{span} \left\{ U_{n,l}^\pm : n \in \Upsilon_h \right\}, \quad l = 0, 1.$$

The discretized full space is defined as  $\mathbb{H}_h := X_h \times Y_h^+ \times Y_h^-$ . Now the finite element approximation associated to (4.3) can be formulated as follows: find  $(E_h, E_h^+, E_h^-) \in \mathbb{H}_h$  such that

$$(5.1) \quad a \left( (E_h, E_h^+, E_h^-), (V_h, V_h^+, V_h^-) \right) = -a \left( (0, E^{in}, 0), (V_h, V_h^+, V_h^-) \right),$$

for all  $(V_h, V_h^+, V_h^-) \in \mathbb{H}_h$ .

**5.2. Auxiliary notation and facts.** Let  $(F, F^+, F^-) \in \mathbb{H}'_h$  be defined as the right-hand side of equation (5.1), and let  $P_h := (P^{X_h}, P^{Y_h^+}, P^{Y_h^-})$  be the orthogonal projection of  $\mathbb{H}$  onto  $\mathbb{H}_h$ . Then we obtain the operator equation of the FEM

$$(5.2) \quad A_h(E_h, E_h^+, E_h^-) = (P_h)^*(F, F^+, F^-), \quad A_h := (P_h)^* A|_{\mathbb{H}_h},$$

where  $A : \mathbb{H} \rightarrow \mathbb{H}'$  is given in (4.7). Note that the operators  $A_h : \mathbb{H}_h \rightarrow \mathbb{H}'_h$  are uniformly bounded in  $h > 0$ . It follows from [17, Lemma 10.10] that  $P^{X_h} \rightarrow I$  in  $X$ , and by the definitions of  $Y_h^\pm$ , we see that  $P^{Y_h^\pm} \rightarrow I$  in  $Y_h^\pm$ . This implies strong convergence of  $P_h$  to  $I$  in  $\mathbb{H}$ . Consequently, the convergence  $(P_h)^* \rightarrow I$  holds in  $\mathbb{H}'$  and  $A_h P_h \rightarrow A$ .

DEFINITION 5.2. The operators  $A_h : \mathbb{H}_h \rightarrow \mathbb{H}'_h$  are called stable if there exists an  $h_0 > 0$  such that  $A_h$  is invertible for all  $h \leq h_0$  and if  $\|A_h^{-1}\| \leq c$  for some constant  $c > 0$  independent of  $h \in (0, h_0)$ .

Note that the operator norm of the inverse operator  $A_h^{-1}$  can be computed as

$$\|A_h^{-1}\|^{-1} = \inf_{\substack{(0,0,0) \neq (E_h, E_h^+, E_h^-) \\ (E_h, E_h^+, E_h^-) \in \mathbb{H}_h}} \sup_{\substack{(0,0,0) \neq (V_h, V_h^+, V_h^-) \\ (V_h, V_h^+, V_h^-) \in \mathbb{H}_h}} \frac{\left| a \left( (E_h, E_h^+, E_h^-), (V_h, V_h^+, V_h^-) \right) \right|}{\| (E_h, E_h^+, E_h^-) \|_{\mathbb{H}} \| (V_h, V_h^+, V_h^-) \|_{\mathbb{H}}}.$$

DEFINITION 5.3. We say that the FEM for (4.3) is convergent if, for any  $(F, F^+, F^-)$  in  $\mathbb{H}'$  and for all  $h < h_0$ , the approximate solution  $(E_h, E_h^+, E_h^-)$  to

$$(5.3) \quad a \left( (E_h, E_h^+, E_h^-), (V_h, V_h^+, V_h^-) \right) = \left\langle (F, F^+, F^-), (V_h, V_h^+, V_h^-) \right\rangle,$$

for all  $(V_h, V_h^+, V_h^-) \in \mathbb{H}_h$  exists and is unique, and if  $(E_h, E_h^+, E_h^-)$  converges strongly in  $\mathbb{H}$  to the exact solution  $(E, E^+, E^-)$  of the continuous variational problem

$$A(E, E^+, E^-) = (F, F^+, F^-).$$

Now we recall two well-known results on the convergence and perturbations (cf., e.g., [21, Chapter 1], [22]), which are our main tools for analyzing the discrete variational problem (5.1). Lemma 5.4 is a simple consequence of the Banach-Steinhaus theorem and, for the reader's convenience, we provide a short proof of Lemma 5.5.

LEMMA 5.4. *Suppose the strong convergence  $P_h \rightarrow I$ . Then the finite element scheme (5.3) is convergent if and only if the operators  $A_h$  defined in (5.2) are stable.*

LEMMA 5.5. *Suppose that  $P_h \rightarrow I$ . Furthermore, suppose that the operators  $B_h : \mathbb{H}_h \rightarrow \mathbb{H}'_h$  are stable, and that the convergence  $B_h P_h \rightarrow B$  holds as  $h \rightarrow 0$  for some operator  $B : \mathbb{H} \rightarrow \mathbb{H}'$ . Moreover, let  $T : \mathbb{H} \rightarrow \mathbb{H}'$  be a compact operator such that  $C := B + T$  is invertible. Let the operators  $C_h : \mathbb{H}_h \rightarrow \mathbb{H}'_h$  be small perturbations of  $B_h + (P_h)^* T|_{\mathbb{H}_h}$ , i.e.,*

$$C_h = B_h + (P_h)^* T|_{\mathbb{H}_h} + D_h, \quad \|D_h\| \rightarrow 0 \quad \text{as } h \rightarrow 0.$$

Then the operators  $C_h$  are stable.

*Proof.* The small perturbations  $D_h$  can be treated by the usual Neumann series argument. Hence, it suffices to prove that the operators  $B_h + (P_h)^* T|_{\mathbb{H}_h} : \mathbb{H}_h \rightarrow \mathbb{H}'_h$  are stable, i.e., that the inverse operators of  $B_h + (P_h)^* T|_{\mathbb{H}_h}$  exist and are uniformly bounded.

We first show that  $B^{-1}$  exists. Since  $C$  is invertible and  $T$  is compact,  $B$  is a Fredholm operator with index zero. Hence, we only need to show that  $\text{Ker} B = \{0\}$ . Noting that the operators  $B_h$  are stable, we get, for any  $u \in \mathbb{H}$ , that  $\|P_h u\|_{\mathbb{H}} = \|B_h^{-1} B_h P_h u\|_{\mathbb{H}} \leq c \|B_h P_h u\|_{\mathbb{H}'}$  with a constant  $c > 0$  independent of  $h$ . Letting  $h \rightarrow 0$ , we obtain  $\|u\|_{\mathbb{H}} \leq c \|B u\|_{\mathbb{H}'}$ , which implies  $\text{Ker} B = \{0\}$ .

Now the pointwise convergence  $B_h^{-1} (P_h)^* \rightarrow B^{-1}$  is easy to see, and thus the norm convergence  $\|B_h^{-1} (P_h)^* - B^{-1}\| \rightarrow 0$  as  $h \rightarrow 0$  follows. A simple calculation shows

$$\begin{aligned} B_h + (P_h)^* T|_{\mathbb{H}_h} &= B_h [I|_{\mathbb{H}_h} + B_h^{-1} (P_h)^* T|_{\mathbb{H}_h}] \\ &= B_h \{P_h (I + B^{-1} T)|_{\mathbb{H}_h} + P_h [B_h^{-1} (P_h)^* - B^{-1}] T|_{\mathbb{H}_h}\}. \end{aligned}$$

To prove the stability of  $B_h + (P_h)^* T|_{\mathbb{H}_h}$ , we only need to prove that of  $P_h (I + B^{-1} T)|_{\mathbb{H}_h}$  because the second term in the curly brackets of the previous identity tends to zero as  $h \rightarrow 0$ . From the invertibility of  $C$ , the existence of  $(I + B^{-1} T)^{-1}$  follows. Then, we can verify that

$$\begin{aligned} &[P_h (I + B^{-1} T)^{-1}|_{\mathbb{H}_h}] [P_h (I + B^{-1} T)|_{\mathbb{H}_h}] \\ &= [P_h (I + B^{-1} T)^{-1} (P_h - I) (I + B^{-1} T)|_{\mathbb{H}_h}] + I|_{\mathbb{H}_h} \\ &= [P_h (I + B^{-1} T)^{-1} (P_h - I) B^{-1} T|_{\mathbb{H}_h}] + I|_{\mathbb{H}_h}, \end{aligned}$$

where  $\|P_h (I + B^{-1} T)^{-1} (P_h - I) B^{-1} T|_{\mathbb{H}_h}\| \leq c \|(P_h - I) B^{-1} T\| \rightarrow 0$ . Hence, the product of  $([P_h (I + B^{-1} T)^{-1} (P_h - I) B^{-1} T|_{\mathbb{H}_h}] + I|_{\mathbb{H}_h})^{-1}$  and  $[P_h (I + B^{-1} T)^{-1}|_{\mathbb{H}_h}]$  is the uniformly bounded inverse of  $P_h (I + B^{-1} T)|_{\mathbb{H}_h}$ .  $\square$

REMARK 5.6. The projection  $P_h$  in Lemma 5.5 can be replaced by operators which are not projections. If the  $P_h$  are orthogonal projections and if  $B_h = (P_h)^* B|_{\mathbb{H}_h}$ , then Lemma 5.5 reduces to the classical stability property of projection methods; see, e.g., [16, Theorem 13.7].

**5.3. Convergence analysis of the FEM.** To prove convergence of the FEM, we need the Hodge decomposition of the discrete functions in  $X_h$ . Define

$$\mathcal{S}_h := \{p_h \in H_{qp}^1(\Omega_b) : p_h|_K \in P_k \text{ for all } K \in \tau_h\}.$$

We have the discrete Hodge decomposition  $X_h = X_{h,0} \oplus X_{h,1}$  analogously to Lemma 4.3, where

$$\begin{aligned} X_{h,1} &:= \{\nabla p_h : p_h \in \mathcal{S}_h\} \subseteq X_1, \\ X_{h,0} &:= \left\{E_h \in X_h : 0 = \int_{\Omega} k^2(x) E_h \cdot \nabla p_h \, dx \text{ for all } \nabla p_h \in X_{h,1}\right\}. \end{aligned}$$

Unfortunately, it is not true that  $X_{h,0} \subset X_0$ . This causes difficulties in our convergence analysis. The following property of discrete compactness will help us to overcome these difficulties.

**DEFINITION 5.7.** *We say that the subspaces  $X_{h,0}$  have the discrete compactness property if, for any sequence  $E_{n,0} \in X_{h_n,0}$ ,  $n = 1, 2, \dots$ , such that  $\|E_{n,0}\|_X < c$  with some  $c$  independent of index  $n$ , there is an element  $E_0 \in X_0$  and a subsequence of  $E_{n,0}$  converging in  $L^2(\Omega)^3$  to  $E_0$ .*

**DEFINITION 5.8.** *Let  $\rho_K$  denote the diameter of the largest sphere inscribed in the tetrahedron  $K$ . We say that the partitions  $\tau_h$  are regular as  $h \rightarrow 0$  if there exist constants  $c, h_0 > 0$  such that  $\max_{K \in \tau_h} (h_K / \rho_K) \leq c$  for all  $h \in (0, h_0)$ .*

Analogously to [17, Theorems 7.17, 7.18, and 11.11], we can prove the following Lemma.

**LEMMA 5.9.** *Suppose that the partitions  $\tau_h$  are regular. Then the subspaces  $X_{h,0}$  possess the property of discrete compactness.*

Finally, the main convergence result is stated in the next theorem.

**THEOREM 5.10.** *Suppose that there only exists the trivial solution to the homogeneous variational equation (4.3) and that the partitions  $\tau_h$  of  $\Omega$  are regular. Then the finite element method (5.1) with Nédélec’s edge elements coupled to truncated Rayleigh series expansions converges.*

*Proof.* Define the discrete subspaces  $\mathbb{H}_{h,l} := X_{h,l} \times Y_{h,l}^+ \times Y_{h,l}^- \subset \mathbb{H}_h$  for  $l = 0, 1$ . Let  $P^{Y_{h,l}^\pm} : \mathbb{H} \rightarrow Y_{h,l}^\pm$ ,  $P^{X_{h,l}} : \mathbb{H} \rightarrow X_{h,l}$ , and  $P^{\mathbb{H}_{h,l}} : \mathbb{H} \rightarrow \mathbb{H}_{h,l}$  be orthogonal projections. Note that  $P^{X_{h,1}} \rightarrow P^{X_1}$  as  $h \rightarrow 0$ , where  $P^{X_1}$  is the orthogonal projection from  $\mathbb{H}$  to  $X_1$ . Indeed, for  $\nabla p \in X_1$ , the problem of finding  $\nabla p_h \in X_{h,1}$  such that  $\langle \nabla p_h, \nabla q_h \rangle = \langle \nabla p, \nabla q_h \rangle$  for all  $\nabla q_h \in X_{h,1}$  corresponds to the finite element scheme for the quasi-periodic boundary value problem of finding  $f \in H_{qp}^1(\Omega)$  such that

$$\Delta f = \Delta p \quad \text{in } \Omega, \quad e_3 \cdot \nabla f = e_3 \cdot \nabla p, \quad \text{on } \Gamma_b^\pm.$$

This boundary value problem only admits the unique quasiperiodic solution  $f = p$  if  $X$  does not contain constant functions. If  $X$  contains constants, i.e., if the direction of incidence is  $\hat{\theta} = (0, 0, -1)^\top$ , then the finite element scheme  $\langle \nabla p_h, \nabla q_h \rangle = \langle \nabla p, \nabla q_h \rangle$  can be considered in the factor space  $H_{qp}^1(\Omega)/\mathbb{C}$ . In any case, we have  $\nabla p_h \rightarrow \nabla p$  in  $L^2(\Omega)^3$  and  $P^{X_{h,1}} \rightarrow P^{X_1}$  as  $h \rightarrow 0$ . This together with  $P_h^X \rightarrow P^X$  implies the convergence  $P^{X_{h,0}} \rightarrow P^{X_0}$  as  $h \rightarrow 0$ . It is easy to see that  $P^{Y_{h,l}^\pm} \rightarrow P^{Y_l^\pm}$ .

Let the operator  $A : \mathbb{H} \rightarrow \mathbb{H}'$  be given as in (4.7). To prove convergence of the FEM, by Lemma 5.4, we only need to prove the stability of  $(P_h)^* A|_{\mathbb{H}_h}$ . For clarity, we divide our proof into five steps by introducing several auxiliary operators and then apply Lemma 5.5.



Step 1. Introduce a new operator  $B_1: \mathbb{H}_1 \rightarrow \mathbb{H}'_1$  as

$$\begin{aligned} & \left\langle B_1(\nabla p, E_1^+, E_1^-), (\nabla \xi, V_1^+, V_1^-) \right\rangle \\ & := \left\langle A|_{\mathbb{H}_1}(\nabla p, E_1^+, E_1^-), (\nabla \xi, V_1^+, V_1^-) \right\rangle - \tilde{a}_4^+(E_1^+, V_1^+) - \tilde{a}_4^-(E_1^-, V_1^-) \\ & \quad - a_6^+((\nabla p, E_1^+, E_1^-), (\nabla \xi, V_1^+, V_1^-)) - a_6^-((\nabla p, E_1^+, E_1^-), (\nabla \xi, V_1^+, V_1^-)), \end{aligned}$$

where the sesquilinear forms  $\tilde{a}_4^\pm$  are given in (4.9). Obviously,  $-B_1$  is positively coercive over  $\mathbb{H}_1$ , i.e.,

$$\begin{aligned} & -\operatorname{Re} \left\langle B_1(\nabla p, E_1^+, E_1^-), (\nabla p, E_1^+, E_1^-) \right\rangle \\ & = a_1(\nabla p, \nabla p) + [a_4^+(E_1^+, E_1^+) + \tilde{a}_4^+(E_1^+, E_1^+)] + [-a_4^-(E_1^-, E_1^-) + \tilde{a}_4^-(E_1^-, E_1^-)] \\ & \geq c \left( \|\nabla p\|_{H(\operatorname{curl}, \Omega)}^2 + \|E_1^+\|_{H(\operatorname{curl}, D^+)}^2 + \|E_1^-\|_{H(\operatorname{curl}, D^-)}^2 \right), \end{aligned}$$

for some constant  $c > 0$ . Thus, the operators  $[(P_{h,1})^* B_1|_{\mathbb{H}_{h,1}}]$  are stable as the Galerkin approximations of  $B_1$ .

Define the operator  $B_0: Z \rightarrow Z'$ ,  $Z := \mathbb{H}_0 \times X_1$  by

$$\begin{aligned} & \left\langle B_0(E_0 + \nabla p, E_0^+, E_0^-), (V_0 + \nabla \xi, V_0^+, V_0^-) \right\rangle \\ & = -a_3^+(E_0^+, V_0^+) + \tilde{a}_3^+(E_0^+, V_0^+) + a_5^+(E_0, V_0^+) - \overline{a_5^+(V_0, E_0^+)} \\ & \quad - a_5^-(E_0, V_0^-) + \overline{a_5^-(V_0, E_0^-)} + a_3^-(E_0^-, V_0^-) + \tilde{a}_3^-(E_0^-, V_0^-) \\ & \quad + \int_{\Omega} [\operatorname{curl} E_0 \cdot \operatorname{curl} \bar{V}_0 + k^2(x) E_0 \cdot \bar{V}_0 + k^2(x) \nabla p \cdot \nabla \bar{\xi}] \, dx, \end{aligned}$$

with the sesquilinear forms  $\tilde{a}_3^\pm$  given in (4.8). From the proof of Theorem 4.1,  $B_0$  is positively coercive over  $Z$ , i.e.,

$$(5.4) \quad \begin{aligned} & \operatorname{Re} \left\langle B_0(E, E_0^+, E_0^-), (E, E_0^+, E_0^-) \right\rangle \\ & \geq c \left( \|E\|_{H(\operatorname{curl}, \Omega)}^2 + \|E_0^+\|_{H(\operatorname{curl}, D^+)}^2 + \|E_0^-\|_{H(\operatorname{curl}, D^-)}^2 \right), \end{aligned}$$

where  $E = E_0 + \nabla p$ . Consequently, the operators  $(P^{\mathbb{H}_{h,0}})^* B_0|_{\mathbb{H}_{h,0}}$  inherit the coercivity of  $B_0$  in (5.4). Note that, although  $\mathbb{H}_{h,0} \subset \mathbb{H}_0$  does not hold in general, we have  $\mathbb{H}_{h,0} \subset Z$ . Therefore, the operators  $(P^{\mathbb{H}_{h,0}})^* B_0|_{\mathbb{H}_{h,0}}: \mathbb{H}_{h,0} \rightarrow \mathbb{H}'_{h,0}$  are stable.

Next, we define the operators  $B: \mathbb{H} \rightarrow \mathbb{H}'$  and  $B_h: \mathbb{H}_h \rightarrow \mathbb{H}'_h$  as follows:

$$\begin{aligned} & \left\langle B(E, E^+, E^-), (E, E^+, E^-) \right\rangle \\ & := \left\langle B_0(E_0, E_0^+, E_0^-), (V_0, V_0^+, V_0^-) \right\rangle + \left\langle B_1(\nabla p, E_1^+, E_1^-), (\nabla p, E_1^+, E_1^-) \right\rangle, \\ & B_h(E_h, E_h^+, E_h^-) := \begin{bmatrix} (P^{\mathbb{H}_{h,0}})^* B_0|_{\mathbb{H}_{h,0}} & 0 \\ 0 & (P^{\mathbb{H}_{h,1}})^* B_1|_{\mathbb{H}_{h,1}} \end{bmatrix} \begin{bmatrix} (E_{h,0}, E_{h,0}^+, E_{h,0}^-) \\ (\nabla p_h, E_{h,1}^+, E_{h,1}^-) \end{bmatrix}. \end{aligned}$$

Obviously, the operators  $B_h$  are stable and the limit operator  $\lim_{h \rightarrow 0} B_h P_h$  is equal to  $B$ . If

we introduce the operators  $T_j : \mathbb{H} \rightarrow \mathbb{H}'$ ,  $j = 0, 1$ , by

$$\begin{aligned} \langle T_0(E, E^+, E^-), (V, V^+, V^-) \rangle &:= -2 \int_{\Omega} k^2(x) E \cdot \bar{V} dx, \\ \langle T_1(E, E^+, E^-), (V, V^+, V^-) \rangle &:= \\ &\quad -\tilde{a}_3^+(E_0^+, V_0^+) - \tilde{a}_3^-(E_0^-, V_0^-) + \tilde{a}_4^+(E_1^+, V_1^+) + \tilde{a}_4^-(E_1^-, V_1^-) \\ &\quad + a_5^+(E_0, V_1^+) - a_5^-(E_0, V_1^-) - \overline{a_5^+(V_0, E_1^+)} + \overline{a_5^-(V_0, E_1^-)} \\ &\quad + a_6^+((E, E^+, E^-), (V, V^+, V^-)) + a_6^-((E, E^+, E^-), (V, V^+, V^-)), \end{aligned}$$

then we arrive at

$$(5.5) \quad \begin{aligned} (P_h)^* A|_{\mathbb{H}_h} &= B_h + (P_h)^* T_1|_{\mathbb{H}_h} \\ &+ \begin{bmatrix} (P^{\mathbb{H}_{h,0}})^* T_0|_{\mathbb{H}_{h,0}} & (P_h^{\mathbb{H}_0})^* (A - T_1)|_{\mathbb{H}_{h,1}} \\ (P^{\mathbb{H}_{h,1}})^* (A - T_1)|_{\mathbb{H}_{h,0}} & 0 \end{bmatrix}. \end{aligned}$$

Step 2. It is easy to see that  $T_1$  is compact over  $\mathbb{H}$ , and the term  $(P_h)^* T_1|_{\mathbb{H}_h}$  can be treated by Lemma 5.5. Next we show that  $(P^{\mathbb{H}_{h,0}})^* T_0|_{\mathbb{H}_{h,0}}$  can be treated by Lemma 5.5 as well. Denote by  $\Pi$  the orthogonal projection from the space  $X$  into  $X_0$  with respect to the inner product  $\langle E, V \rangle_X = \int_{\Omega} \{\text{curl } E \cdot \overline{\text{curl } V} + E \cdot \bar{V}\} dx$ . Then,  $\Pi$  is also an orthogonal projection in the  $L^2(\Omega)^3$  sense. Moreover, by the proof of Lemma 4.3, the operator  $I - \Pi: X \rightarrow X_1$  is an orthogonal projection too. By the definitions of  $T_0$  and  $\Pi$ ,

$$(5.6) \quad \begin{aligned} [(P^{\mathbb{H}_{h,0}})^* T_0|_{\mathbb{H}_{h,0}}] P^{\mathbb{H}_{h,0}}|_{\mathbb{H}_h} &= (P_h)^* T_2|_{\mathbb{H}_h} + D_h^{(0)} + D_h^{(1)}, \\ D_h^{(0)} &:= -2(P^{X_{h,0}})^* [k^2(x)(I - \Pi)]|_{X_{h,0}} P^{X_{h,0}}|_{\mathbb{H}_h}, \\ D_h^{(1)} &:= -2(P_h)^* (P^{X_{h,0}} - P^{X_0})^* [k^2(x)\Pi]|_{X_{h,0}} P^{X_{h,0}}|_{\mathbb{H}_h} \\ &\quad - 2(P_h)^* (P^{X_0})^* [k^2(x)\Pi]|_{X_{h,0}} (P^{X_{h,0}} - P^{X_0})|_{\mathbb{H}_h}, \\ T_2 &:= -2(P^{X_0})^* [k^2(x)\Pi] P^{X_0}. \end{aligned}$$

Here  $T_2$  is compact due to Lemma 4.3. Again  $(P_h)^* T_2|_{\mathbb{H}_h}$  can be treated by Lemma 5.5 and it remains to show  $\|D_h^{(j)}\|_{\mathbb{H}_h \rightarrow \mathbb{H}'_h} \rightarrow 0$  for  $j = 0, 1$ .

The convergence  $\|D_h^{(1)}\| \rightarrow 0$  follows easily since  $[k^2(x)\Pi]: X \rightarrow X'$  is compact and since  $(P^{X_{h,0}} - P^{X_0}) \rightarrow 0$ . Consequently, it remains to prove that  $\|D_h^{(0)}\| \rightarrow 0$  as  $h \rightarrow 0$ . It suffices to show that  $\|(I - \Pi)|_{X_{h,0}}\|_{X_{h,0} \rightarrow X'} \rightarrow 0$  with  $h \rightarrow 0$ , i.e., that, for any sequence  $\|(I - \Pi)|_{X_{h_n,0}}\|_{X_{h_n,0} \rightarrow X'}$  with  $h_n \rightarrow 0$ , there is a subsequence tending to zero. Choose  $E_{h_n,0} \in X_{h_n,0}$  such that  $\|E_{h_n,0}\|_X = 1$  and  $\|(I - \Pi)E_{h_n,0}\|_{X'} = \|(I - \Pi)|_{X_{h_n,0}}\|_{X_{h_n,0} \rightarrow X'}$ . Recalling Lemma 5.9, without loss of generality we can assume the convergence  $E_{h_n,0} \rightarrow E_0 \in X_0$  in  $L^2(\Omega)^3$ . Since  $\Pi$  is bounded in  $L^2$ , we have  $(I - \Pi)E_{h_n,0} \rightarrow (I - \Pi)E_0 = 0$  in  $L^2(\Omega)^3$ . Noting that  $X \subseteq L^2(\Omega)^3$  and  $L^2(\Omega)^3 \subseteq X'$ , we finally conclude

$$\|(I - \Pi)E_{h_n,0}\|_{X'} \leq \|(I - \Pi)E_{h_n,0}\|_{L^2(\Omega)^3} \rightarrow \|(I - \Pi)E_0\|_{L^2(\Omega)^3} = 0.$$

This gives  $\|D_h^{(0)}\| \rightarrow 0$  as  $h \rightarrow 0$ .

Step 3. For  $E_h, V_h \in X_h$ , recall the decompositions

$$\begin{aligned} E_h &= E_{h,0} + \nabla p_h = \Pi(E_{h,0}) + (I - \Pi)(E_{h,0}) + \nabla p_h, \\ V_h &= V_{h,0} + \nabla \xi_h = \Pi(V_{h,0}) + (I - \Pi)(V_{h,0}) + \nabla \xi_h, \end{aligned}$$

with  $E_{h,0}, V_{h,0} \in X_{h,0}$  and  $\nabla p_h, \nabla \xi_h \in X_{h,1}$ . We set  $T := T_0 + T_1$  and claim that

$$(5.7) \quad (P_h)^* A|_{\mathbb{H}_h} = B_h + (P_h)^* T|_{\mathbb{H}_h} + D_h^{(0)} + D_h^{(1)} + D_h^{(2)},$$

where  $D_h^{(2)} : \mathbb{H}_h \rightarrow \mathbb{H}'_h$  is defined by

$$\begin{aligned} & \left\langle D_h^{(2)}(E_h, E_h^+, E_h^-), (V_h, V_h^+, V_h^-) \right\rangle \\ & := a_5^+ \left( (I - \Pi)(E_{h,0}), V_{h,1}^+ \right) - a_5^- \left( (I - \Pi)(E_{h,0}), V_{h,1}^- \right) \\ & \quad - \overline{a_5^+ \left( (I - \Pi)(V_{h,0}), E_{h,1}^+ \right)} + \overline{a_5^- \left( (I - \Pi)(V_{h,0}), E_{h,1}^- \right)}. \end{aligned}$$

In fact, the formulas (5.5) and (5.6) imply (5.7) if we can show that the operator  $D_h^{(2)}$  is the off-diagonal part of the matrix on the right-hand side of (5.5). Hence, it suffices to prove  $D_h^{(2)} = [(P^{\mathbb{H}_{h,0}})^*(A - T_1)|_{\mathbb{H}_{h,1}}]P^{\mathbb{H}_{h,1}} + [(P^{\mathbb{H}_{h,1}})^*(A - T_1)|_{\mathbb{H}_{h,0}}]P^{\mathbb{H}_{h,0}}$ . We conclude

$$\begin{aligned} & \left\langle (P^{\mathbb{H}_{h,1}})^*(A - T_1)|_{\mathbb{H}_{h,0}}(E_h, E_h^+, E_h^-), (V_h, V_h^+, V_h^-) \right\rangle \\ & = -a_1 \left( (I - \Pi)(E_{h,0}), \nabla \xi_h \right) + a_5^+ \left( (I - \Pi)(E_{h,0}), V_{h,1}^+ \right) \\ (5.8) \quad & \quad - a_5^- \left( (I - \Pi)(E_{h,0}), V_{h,1}^- \right) \\ & = a_5^+ \left( (I - \Pi)(E_{h,0}), V_{h,1}^+ \right) - a_5^- \left( (I - \Pi)(E_{h,0}), V_{h,1}^- \right), \end{aligned}$$

where we have used the identity

$$a_1 \left( (I - \Pi)(E_{h,0}), \nabla \xi_h \right) = \int_{\Omega} k^2(x) E_{h,0} \cdot \overline{\nabla \xi_h} dx - \int_{\Omega} k^2(x) \Pi(E_{h,0}) \cdot \overline{\nabla \xi_h} dx = 0.$$

Analogously, it can be seen that

$$\begin{aligned} & \left\langle (P^{\mathbb{H}_{h,0}})^*(A - T_1)|_{\mathbb{H}_{h,1}}(E_h, E_h^+, E_h^-), (V_h, V_h^+, V_h^-) \right\rangle \\ (5.9) \quad & = -a_5^+ \left( (I - \Pi)(V_{h,0}), E_{h,1}^+ \right) + \overline{a_5^- \left( (I - \Pi)(V_{h,0}), E_{h,1}^- \right)}. \end{aligned}$$

Equations (5.8) and (5.9) imply that

$$[(P^{\mathbb{H}_{h,0}})^*(A - T_1)|_{\mathbb{H}_{h,1}}]P^{\mathbb{H}_{h,1}} + [(P^{\mathbb{H}_{h,1}})^*(A - T_1)|_{\mathbb{H}_{h,0}}]P^{\mathbb{H}_{h,0}}$$

coincides with  $D_h^{(2)}$ . Formula (5.7) is thus proven.

Step 4. We prove  $\|D_h^{(2)}\| \rightarrow 0$ . First we derive that  $\|(P^{\mathbb{H}_{h,1}})^*(A - T_1)|_{\mathbb{H}_{h,0}}\| \rightarrow 0$ . By (5.8), we choose functions  $E_{h,0}$  and  $V_{h,1}^\pm$  with  $\|E_{h,0}\|_{H(\text{curl}, \Omega)} = 1$  and  $\|V_{h,1}^\pm\|_{H(\text{curl}, D^\pm)} = 1$  such that

$$\|(P^{\mathbb{H}_{h,1}})^*(A - T_1)|_{\mathbb{H}_{h,0}}\| = a_5^+ (\nabla q_h, V_{h,1}^+) - a_5^- (\nabla q_h, V_{h,1}^-), \quad \nabla q_h := (I - \Pi)E_{h,0}.$$

Using the definition of  $a_5^\pm$ , we get

$$\begin{aligned} & a_5^+ (\nabla q_h, V_{h,1}^+) = - \int_{\Gamma_b^+} e_3 \times \nabla q_h \cdot \overline{\text{curl } V_{h,1}^+} ds, \\ (5.10) \quad & |a_5^+ (\nabla q_h, V_{h,1}^+)| \leq \|e_3 \times \nabla q_h\|_{H_t^{-1/2}(\Gamma_b^+)} \|\text{curl } V_{h,1}^+\|_{H_t^{1/2}(\Gamma_b^+)}. \end{aligned}$$

On the one hand, we have that for any  $a \in \mathbb{C}$ ,

$$\begin{aligned} \|e_3 \times \nabla q_h\|_{H_t^{-1/2}(\Gamma_b^+)} &= \|\nabla_{\Gamma_b^+}(q_h + a)\|_{H_t^{-1/2}(\Gamma_b^+)} \leq c \|q_h + a\|_{H^1(\Omega)} \\ &\leq c \|q_h + a\|_{H^1(\Omega)}, \end{aligned}$$

where  $\nabla_{\Gamma_b^+}$  denotes the surface gradient operator over  $\Gamma_b^+$ . Hence,

$$(5.11) \quad \|e_3 \times \nabla q_h\|_{H_t^{-1/2}(\Gamma_b^+)} \leq c \inf_{a \in \mathbb{C}} \|q_h + a\|_{H^1(\Omega)} \leq c \|\nabla q_h\|_{L^2(\Omega)^3}.$$

On the other hand, for  $V_{h,1}^+ = \sum_{n: |n| \leq C/h} c_n U_{n,1}^+ \in Y_1^+$ , there holds

$$(5.12) \quad \begin{aligned} \|\operatorname{curl} V_{h,1}^+\|_{H_t^{1/2}(\Gamma_b^+)} &\leq \|\operatorname{curl} V_{h,1}^+\|_{H^1(D^+)^3} \\ &= \left\| \sum_{n: |n| \leq C/h} c_n \operatorname{curl} U_{n,1}^+ \right\|_{H^1(D^+)^3}. \end{aligned}$$

In view of the identity  $\operatorname{curl} U_{n,1}^+ = -i(k^+)^2 / \sqrt{|\alpha_n|^2 + |\beta_n^+|^2} U_{n,0}^+$  (see [14, Lemma 3.1]) and the relation  $\sqrt{|\alpha_n|^2 + |\beta_n^+|^2} = \mathcal{O}(|n|)$  as  $|n| \rightarrow \infty$ , we get

$$\begin{aligned} \left\| \sum_{n: |n| \leq C/h} c_n \operatorname{curl} U_{n,1}^+ \right\|_{H^1(D^+)^3} &\leq c \left( \sum_{n: |n| \leq C/h} |c_n|^2 |n|^{-2} \|U_{n,0}^+\|_{H^1(D^+)^3} \right)^{1/2} \\ &\leq c \left( \sum_{n: |n| \leq C/h} |c_n|^2 \|U_{n,0}^+\|_{L^2(D^+)^3} \right)^{1/2} \\ &\leq c \left( \sum_{n: |n| \leq C/h} |c_n|^2 \|U_{n,1}^+\|_{L^2(D^+)^3} \right)^{1/2}, \end{aligned}$$

where the last two equalities follow from the estimates derived in the proof of [14, Lemma 4.5]. Recalling (5.12) and the representation of  $V_{h,1}^+$  as an expansion with respect to the basis functions  $U_{n,1}^+$ , we obtain

$$(5.13) \quad \|\operatorname{curl} V_{h,1}^+\|_{H_t^{1/2}(\Gamma_b^+)} \leq c \|V_{h,1}^+\|_{L^2(D^+)^3} \leq c \|V_{h,1}^+\|_{H(\operatorname{curl}, D^+)} \leq c.$$

Inserting the estimates (5.11) and (5.13) into (5.10) yields  $|a_5^+(\nabla q_h, V_{h,1}^+)| \leq c_5^+ \|\nabla q_h\|_{L^2(\Omega)^3}$  for some  $c_5^+ > 0$ , and analogously, there exists another non-negative constant  $c_5^-$  such that  $|a_5^-(\nabla q_h, V_{h,1}^-)| \leq c_5^- \|\nabla q_h\|_{L^2(\Omega)^3}$ . Thus, to prove  $\|(P^{\mathbb{H}h,1})^*(A - T_1)|_{\mathbb{H}_{h,0}}\| \rightarrow 0$ , we only need to verify  $\|\nabla q_h\|_{L^2(\Omega)^3} \rightarrow 0$  as  $h \rightarrow 0$ . However, we can choose  $E_{h,0}$  with  $\|E_{h,0}\|_X = 1$  such that  $\|(P^{\mathbb{H}h,1})^*(A - T_1)E_{h,0}\| = \|(P^{\mathbb{H}h,1})^*(A - T_1)|_{\mathbb{H}_{h,0}}\|$ . From the discrete compactness of the space  $X_{h,0}$  in Lemma 5.9, for any sequence  $E_{h_n,0}$ , we can always find a subsequence converging in  $L^2(\Omega)^3$  to an  $E_0 \in X_0$ . We denote this subsequence again by  $E_{h_n,0}$ . Then  $\|\nabla q_{h_n}\|_{L^2(\Omega)^3} = \|(I - \Pi)E_{h_n,0}\|_{L^2(\Omega)^3} \rightarrow \|(I - \Pi)E_0\|_{L^2(\Omega)^3} = 0$ . In other words, any sequence  $\|(P^{\mathbb{H}h_n,1})^*(A - T_1)E_{h_n,0}\|$  has a subsequence tending to zero. Consequently,  $\|(P^{\mathbb{H}h,1})^*(A - T_1)E_{h,0}\|$  converges to zero.

Arguing analogously, one can prove the convergence  $\|(P^{\mathbb{H}h,0})^*(A - T_1)|_{\mathbb{H}_{h,1}}\| \rightarrow 0$  as  $h \rightarrow 0$  via the identity (5.9). Hence, it holds that  $\|D_h^{(2)}\| \rightarrow 0$ .

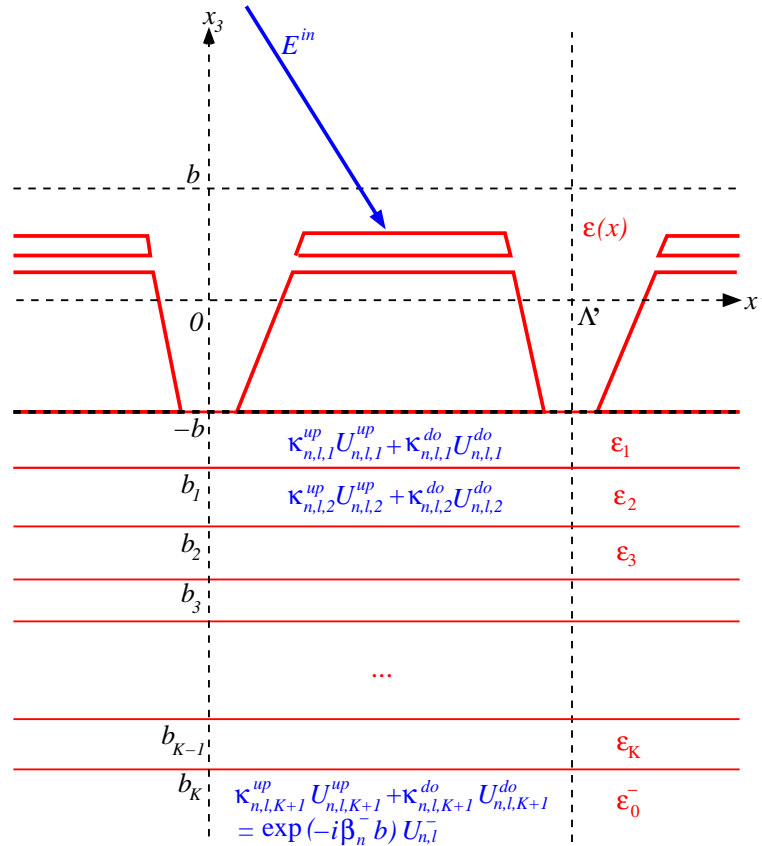


FIG. 6.1. Grating with multi-layer system.

Step 5. Setting  $D_h := D_h^{(0)} + D_h^{(1)} + D_h^{(2)}$ , equation (5.7) is the representation of Lemma 5.5. It can be concluded from Steps 1-4 that the operators  $B_h$  are stable,  $T$  is compact, and that  $D_h$  is only a small perturbation. By the uniqueness assumption in Theorem 5.10, we see from Theorem 3.1 that  $A$  is invertible. Now, applying Lemma 5.5 yields the stability of  $(P_h)^* A|_{\mathbb{H}_h}$ . The proof of the convergence of the FEM is thus completed.  $\square$

**6. Multi-layer system beneath the grating structure.** In many applications there is an adjacent multi-layer system beneath the lower face  $x_3 = -b$  of the grating. More precisely, as indicated in Figure 6.1, for a sequence  $b_k$ ,  $k = 0, \dots, K$ , of  $x_3$ -coordinates such that  $-b = b_0 > b_1 > \dots > b_K$ , the function  $\epsilon(x) + i\sigma(x)/\omega$  in the layer  $b_{k-1} > x_3 > b_k$  takes the constant value  $\epsilon_k$  with  $\text{Im } \epsilon_k \geq 0$  such that  $\text{Re } \epsilon_k > 0$  for  $\text{Im } \epsilon_k = 0$ . Of course, in the lower half space  $b_K > x_3$ , we suppose  $\epsilon(x) + i\sigma(x)/\omega = \epsilon_0^-$ .

For a variational formulation adapted to the multi-layer system, we need modified spaces  $Y_l^-$ ,  $l = 0, 1$ . Clearly, the tangential traces of  $E$  and  $\text{curl } E$  are continuous over the interfaces  $x_3 = b_k$ . Solving these transmission problems, each downward propagating mode  $E = \exp(-i\beta_n^- b) U_{n,l}^-$  in the half space  $b_K > x_3$  corresponds to an extended field  $E$  in  $b_0 > x_3$  such that  $E(x) = \kappa_{n,l,k}^{\text{up}} U_{n,l,k}^{\text{up}}(x) + \kappa_{n,l,k}^{\text{do}} U_{n,l,k}^{\text{do}}(x)$  for  $b_{k-1} > x_3 > b_k$ ,  $k = 1, 2, \dots, K$ ,

where  $\kappa_{n,l,k}^{\text{up}}, \kappa_{n,l,k}^{\text{do}} \in \mathbb{C}$  and

$$U_{n,0,k}^{\text{up}}(x) := e^{i[\alpha_n \cdot x' + \beta_{n,k}(x_3+b)]} \begin{cases} (0, -1, 0)^\top & \text{if } |\alpha_n|=0, \\ \frac{1+i(x_3+b)}{|\alpha_n|} (-\alpha_n^{(2)}, \alpha_n^{(1)}, 0)^\top & \text{if } \beta_{n,k}=0, \\ \frac{1}{|\alpha_n|} (-\alpha_n^{(2)}, \alpha_n^{(1)}, 0)^\top & \text{otherwise,} \end{cases}$$

$$U_{n,1,k}^{\text{up}}(x) := e^{i[\alpha_n \cdot x' + \beta_{n,k}(x_3+b)]} \begin{cases} (1, 0, 0)^\top & \text{if } |\alpha_n|=0, \\ \frac{1}{\sqrt{|\alpha_n|^2 + |\alpha_n|^4}} \left( -\alpha_n, |\alpha_n|^2(1+i(x_3+b)) \right)^\top & \text{if } \beta_{n,k}=0, \\ \frac{1}{|\alpha_n| \sqrt{|\alpha_n|^2 + |\beta_{n,k}|^2}} (-\beta_{n,k} \alpha_n, |\alpha_n|^2)^\top & \text{otherwise,} \end{cases}$$

$$U_{n,0,k}^{\text{do}}(x) := e^{i[\alpha_n \cdot x' - \beta_{n,k}(x_3+b)]} \begin{cases} (0, 1, 0)^\top & \text{if } |\alpha_n|=0, \\ \frac{1-i(x_3+b)}{|\alpha_n|} (\alpha_n^{(2)}, -\alpha_n^{(1)}, 0)^\top & \text{if } \beta_{n,k}=0, \\ \frac{1}{|\alpha_n|} (\alpha_n^{(2)}, -\alpha_n^{(1)}, 0)^\top & \text{otherwise,} \end{cases}$$

$$U_{n,1,k}^{\text{do}}(x) := e^{i[\alpha_n \cdot x' - \beta_{n,k}(x_3+b)]} \begin{cases} (-1, 0, 0)^\top & \text{if } |\alpha_n|=0, \\ \frac{1}{\sqrt{|\alpha_n|^2 + |\alpha_n|^4}} \left( \alpha_n, |\alpha_n|^2(1-i(x_3+b)) \right)^\top & \text{if } \beta_{n,k}=0, \\ \frac{1}{|\alpha_n| \sqrt{|\alpha_n|^2 + |\beta_{n,k}|^2}} (\beta_{n,k} \alpha_n, |\alpha_n|^2)^\top & \text{otherwise,} \end{cases}$$

$$\beta_{n,k} := \sqrt{\omega^2 \mu_0 \epsilon_k - |\alpha_n|^2}, \quad \beta_{n,K+1} := \beta_n^-.$$

Fix  $n$  and  $l$ . It is not hard to see (cf. [19, Section III.4]) that, for each linear combination of  $U_{n,l,K+1}^{\text{up}}$  and  $U_{n,l,K+1}^{\text{do}}$  in the half space  $x_3 < b_K$ , there exist unique linear combinations of the  $U_{n,l,k}^{\text{up}}$  and  $U_{n,l,k}^{\text{do}}$  in the layers  $b_k < x_3 < b_{k-1}$ ,  $k=1, \dots, K$ , such that the tangential traces over the interfaces  $x_3 = b_k$ ,  $k=1, \dots, K$ , of the functions and of their curls in the adjacent layers coincide. Similarly, to each linear combination of  $U_{n,l,1}^{\text{up}}$  and  $U_{n,l,1}^{\text{do}}$  in the layer  $b_1 < x_3 < b_0$  there exist unique linear combinations of the  $U_{n,l,k}^{\text{up}}$  and  $U_{n,l,k}^{\text{do}}$  in the layers  $b_k < x_3 < b_{k-1}$ ,  $k=2, \dots, K$ , and in the half space  $x_3 < b_K$  such that the tangential traces of the functions and of their curls in adjacent layers coincide. Hence, the coefficients  $\kappa_{n,l,k}^{\text{up}}, \kappa_{n,l,k}^{\text{do}}$  are uniquely determined. For instance, if all the  $\beta_{n,k}$  are non-zero and  $|\alpha_n| \neq 0$ , then

$$(6.1) \quad \begin{pmatrix} \kappa_{n,l,1}^{\text{up}} \\ \kappa_{n,l,1}^{\text{do}} \end{pmatrix} = \mathcal{M}_{n,l,1} \mathcal{M}_{n,l,2} \dots \mathcal{M}_{n,l,K} \begin{pmatrix} 0 \\ 1 \end{pmatrix},$$

$$\mathcal{M}_{n,0,k} := \begin{pmatrix} \frac{\beta_{n,k+1} + \beta_{n,k}}{2\beta_{n,k}} e^{i[\beta_{n,k+1} - \beta_{n,k}]b_k} & \frac{\beta_{n,k+1} - \beta_{n,k}}{2\beta_{n,k}} e^{-i[\beta_{n,k+1} + \beta_{n,k}]b_k} \\ \frac{\beta_{n,k+1} - \beta_{n,k}}{2\beta_{n,k}} e^{i[\beta_{n,k+1} + \beta_{n,k}]b_k} & \frac{\beta_{n,k+1} + \beta_{n,k}}{2\beta_{n,k}} e^{-i[\beta_{n,k+1} - \beta_{n,k}]b_k} \end{pmatrix},$$

$$\mathcal{M}_{n,1,k} := \sqrt{\frac{|\alpha_n|^2 + |\beta_{n,k}|^2}{|\alpha_n|^2 + |\beta_{n,k+1}|^2}} \begin{pmatrix} \left[ \frac{|\alpha_n|^2 + \beta_{n,k+1}^2}{|\alpha_n|^2 + \beta_{n,k}^2} + \frac{\beta_{n,k+1}}{\beta_{n,k}} \right] e^{i[\beta_{n,k+1} - \beta_{n,k}]b_k} & \left[ \frac{|\alpha_n|^2 + \beta_{n,k+1}^2}{|\alpha_n|^2 + \beta_{n,k}^2} - \frac{\beta_{n,k+1}}{\beta_{n,k}} \right] e^{-i[\beta_{n,k+1} + \beta_{n,k}]b_k} \\ \left[ \frac{|\alpha_n|^2 + \beta_{n,k+1}^2}{|\alpha_n|^2 + \beta_{n,k}^2} - \frac{\beta_{n,k+1}}{\beta_{n,k}} \right] e^{i[\beta_{n,k+1} + \beta_{n,k}]b_k} & \left[ \frac{|\alpha_n|^2 + \beta_{n,k+1}^2}{|\alpha_n|^2 + \beta_{n,k}^2} + \frac{\beta_{n,k+1}}{\beta_{n,k}} \right] e^{-i[\beta_{n,k+1} - \beta_{n,k}]b_k} \end{pmatrix}.$$

Note that the coefficients  $\kappa_{n,l,1}^{\text{up}}$  and  $\kappa_{n,l,1}^{\text{do}}$  can be computed by numerically stable algorithms; see, e.g., [19, Section III.6].

Setting  $\tilde{U}_{n,l}^- := \kappa_{n,l,1}^{\text{do}} U_{n,l,1}^{\text{do}} + \kappa_{n,l,1}^{\text{up}} U_{n,l,1}^{\text{up}}$ , we define the modified spaces  $Y_l^-$  by (4.1) but with  $U_{n,l}^-$  replaced by  $\tilde{U}_{n,l}^-$ . Now the new variational formulation for the transmission problem is just (4.3) with a modified sesquilinear form (4.2) defined over  $\mathbb{H} := X \times Y^+ \times (Y_0^- \oplus Y_1^-)$  including the modified spaces  $Y_l^-$ . The modified sesquilinear form is the sum of (4.2) and the additional term

$$-\eta^- \sum_{l=0}^1 \sum_{n: e_3 \times \tilde{U}_{n,l}^- = 0} \left[ \int_{\Gamma_b^-} e_3 \times (E - E^-) \cdot (e_3 \times \bar{U}_{n,l}^-) \, ds \int_{\Gamma_b^-} (\text{curl } V^-) \cdot (\text{curl } \bar{U}_{n,l}^-) \, ds \right].$$

REMARK 6.1. All the results for the variational formulation and for the FEM coupled by the wave modes remain true for the case of multi-layer systems beneath the grating structure and the new variational form.

Indeed, we sketch the proof. From the definitions of the  $U_{n,l,1}^{\text{up}}$  and  $U_{n,l,1}^{\text{do}}$ , we observe that  $e_3 \times U_{n,l,1}^{\text{up}} = -e_3 \times U_{n,l,1}^{\text{do}}$  and  $(\text{curl } U_{n,l,1}^{\text{up}})_T = (\text{curl } U_{n,l,1}^{\text{do}})_T$  over the curve  $\Gamma_b^-$ . Consequently, the traces entering the sesquilinear forms satisfy

$$(6.2) \quad \begin{aligned} e_3 \times \tilde{U}_{n,l}^- &= [\kappa_{n,l,1}^{\text{do}} - \kappa_{n,l,1}^{\text{up}}] e_3 \times U_{n,l,1}^{\text{do}}, \\ (\text{curl } \tilde{U}_{n,l}^-)_T &= [\kappa_{n,l,1}^{\text{do}} + \kappa_{n,l,1}^{\text{up}}] (\text{curl } U_{n,l,1}^{\text{do}})_T. \end{aligned}$$

If  $\beta_{n,1} = 0$ , then  $[\kappa_{n,1,1}^{\text{do}} - \kappa_{n,1,1}^{\text{up}}] \neq 0$  since otherwise  $e_3 \times \tilde{U}_{n,1}^- = 0$ , which together with  $(\text{curl } \tilde{U}_{n,1}^-)_T = 0$  would be a contradiction to the one-to-one mapping between the linear combinations of wave modes mentioned above. This fact and the special choice of the additional term in the modified sesquilinear form guarantee (cf. [14, proof of Lemma 3.3]) the equivalence of the boundary value problem and the variational equation in the case of multi-layer systems.

The Fredholm property with index zero for the variational operator and convergence of the FEM coupled by wave modes follow from the fact that the operator corresponding to the modified variational form is a compact perturbation of that of the original form. To see this fact, we observe that

$$\beta_{n,k}/|n| \rightarrow i \quad \text{for } |n| \rightarrow \infty, \quad \text{and} \quad \beta_{n,k} - \beta_{n,k+1} = (k_k^2 - k_{k+1}^2)/(\beta_{n,k} + \beta_{n,k+1}) \sim |n|^{-1}$$

with  $k_k := \omega \sqrt{\epsilon_k \mu_0}$ . Consequently, equation (6.1) implies that  $\kappa_{n,l,1}^{\text{do}} \rightarrow 1$ ,  $\kappa_{n,l,1}^{\text{up}} \rightarrow 0$  and  $[\kappa_{n,l,1}^{\text{do}} \pm \kappa_{n,l,1}^{\text{up}}] \rightarrow 1$  for the factors in (6.2). In other words, the difference between the modified operator and the original one is the multiplication by operators represented with respect to the wave mode basis by the diagonal matrices  $([\kappa_{n,l,1}^{\text{do}} \pm \kappa_{n,l,1}^{\text{up}}] \delta_{n,n'})_{n,n'}$ . In view of

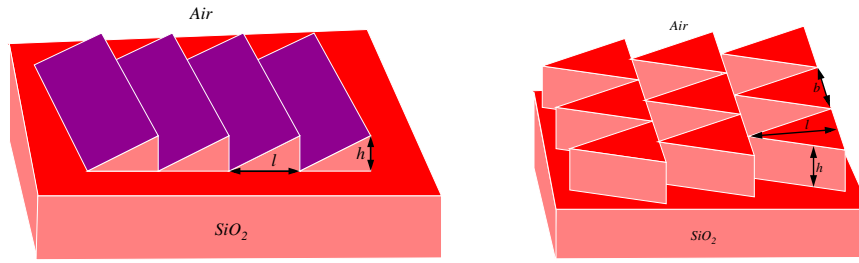


FIG. 7.1. Geometry of grating: left: echelle grating, right: blazes.

$$U_{n,l,1}^{\text{do}} = \exp(-i\beta_n^- b) U_{n,l}^- \text{ and}$$

$$\left\| \sum_{n \in \mathbb{Z}^2} \sum_{l=0}^1 c_{n,l} U_{n,l}^- \right\|_{H(\text{curl}, D^-)} \sim \left( \sum_{n \in \mathbb{Z}^2} \sum_{l=0}^1 e^{-2|n|b} \frac{1 + |n|^{2l}}{1 + |n|} |c_{n,l}|^2 \right)^{1/2}$$

(cf. [14, Lemma 3.1]), such a diagonal operator is a compact perturbation of the identity.

**7. Numerical example.** For a simple numerical test, we consider two profile gratings on the surface of a  $\text{SiO}_2$  body. The echelle grating (cf. the left of Fig. 7.1) is designed to deflect light into the direction specular with respect to the inclined upper faces. The idea of blazes (cf. the right of Fig. 7.1) with the width  $b$  less and the length  $l$  larger than the wavelength of light  $\lambda$ , is to provide a similar effective medium distribution and to function like an echelle grating. Hopefully, such blazes are of better stability (cf. [12]).

In Table 7.1 we compare the new 3D coupling algorithm (4.3) of Section 5.1 applied to the 2D echelle grating with the reliable results of the 2D FEM code solving the Helmholtz equation. The efficiencies

$$e_n^+ := \frac{\beta_n^+}{\beta_{(0,0)}^+} |E_n^+|^2, \quad e_n^- := \frac{(k^+)^2 \beta_n^-}{(k^-)^2 \beta_{(0,0)}^+} |E_n^-|^2$$

of the electric field solution are computed for wavelength  $\lambda = 500 \text{ nm}$ , period  $l = 10 \mu\text{m}$ , and height  $h = 0.5 \mu\text{m}$ . The grating is illuminated exactly from above under TE polarization. The FEM of Section 5.1 is applied with quadratic edge elements. The upper coupling modes  $n = (n_1, n_2)$  are restricted to  $|n_1| \leq 22$  and  $|n_2| \leq 2$ , the lower modes to  $|n_1| \leq 32$  and  $|n_2| \leq 2$ . Moreover, the coupling parameters  $\eta^\pm$  are set to zero. For the mesh-size tending to zero, the 3D results converge to those of the 2D simulation. Adding more coupling modes does not improve the accuracy.

Next we apply the same 3D algorithm to the blazes and compare the results with those obtained by the algorithm of Huber et al. (cf. [15]). Here the periods are chosen as  $\Lambda_1 = l = 10 \mu\text{m}$  and  $\Lambda_2 = b = \lambda/2$  and the other parameters like for the echelle grating. The resulting efficiencies coincide up to numerical errors; see Table 7.2.



TABLE 7.1

Computation of efficiencies for echelle grating. Comparison of FEM from Section 5.1 with two-dimensional FEM simulation.

meshsize	$e_{-2,0}^+$	$e_{0,0}^+$	$e_{1,0}^-$	$e_{2,0}^-$
125.0 nm	4.82	0.0027	43.23	3.78
62.5 nm	4.530	0.0022	45.0080	4.1289
31.2 nm	4.5039	0.0019	45.0559	4.1142
2D code	4.5025	0.0019	45.0630	4.1145

TABLE 7.2

Computation of efficiencies for blases. Comparison of the FEM from Section 5.1 (left numbers in column) with the FEM of [15] (right numbers).

meshsize	$e_{0,0}^+$	$e_{0,0}^+$	$e_{1,0}^+$	$e_{1,0}^+$
125.0 nm	2.8328	3.0985	0.1661	0.1661
62.5 nm	2.8172	2.8333	0.1918	0.1918
31.2 nm	2.8119	2.8136	0.1944	0.1944

meshsize	$e_{0,0}^-$	$e_{0,0}^-$	$e_{1,0}^-$	$e_{1,0}^-$
125.0 nm	75.2800	76.289	10.1503	10.1465
62.5 nm	75.5412	75.553	10.7248	10.7197
31.2 nm	75.4717	75.490	10.7787	10.7711

REFERENCES

[1] T. ABBOUD, *Formulation variationnelle des équations de Maxwell dans un réseau bipériodique de  $\mathbb{R}^3$* , C. R. Acad. Sci. Paris Sér. I Math., 317 (1993), pp. 245–248.

[2] H. AMMARI, *Uniqueness theorems for an inverse problem in a doubly periodic structure*, Inverse Problems, 11 (1995), pp. 823–833.

[3] H. AMMARI AND G. BAO, *Coupling of finite element and boundary element methods for the scattering by periodic chiral structures*, J. Comput. Math., 3 (2008), pp. 261–283.

[4] ———, *Maxwell’s equations in periodic chiral structures*, Math. Nachr., 251 (2003), pp. 3–18.

[5] H. AMMARI AND J. C. NÉDÉLEC, *Analysis of the diffraction from chiral gratings*, in Mathematical Modeling in Optical Science, G. Bao, L. Cowsar, and W. Masters, eds., Frontiers Appl. Math., 22, SIAM, Philadelphia, 2001, pp. 179–206.

[6] ———, *Coupling finite elements and integral equations to solve the Maxwell equations in a heterogeneous medium*, in Équations aux Dérivées Partielles et Applications, Éd. Sci. Méd. Elsevier, Gauthier-Villars, Paris, 1998, pp. 19–33.

[7] G. BAO, *Variational approximation of Maxwell’s equation in biperiodic structures*, SIAM J. Appl. Math., 57 (1997), pp. 364–381.

[8] G. BAO AND D. C. DOBSON, *On the scattering by biperiodic structures*, Proc. Amer. Math. Soc., 128 (2000), pp. 2715–2723.

[9] A. BUFFA, M. COSTABEL, AND D. SHEEN, *On traces for  $H(\text{curl}, \Omega)$  in Lipschitz domains*, J. Math. Anal. Appl., 276 (2002), pp. 847–867.

[10] X. CHEN AND A. FRIEDMAN, *Maxwell’s equations in a periodic structure*, Trans. Amer. Math. Soc., 323 (1991), pp. 465–507.

[11] D. C. DOBSON, *A variational method for electromagnetic diffraction in biperiodic structures*, RAIRO Modél. Math. Anal. Numér., 28 (1994), pp. 419–439.

[12] H. ELFSTRÖM, M. KUITTINEN, T. VALLIUS, B. H. KLEEMANN, J. RUOFF, AND R. ARNOLD, *Fabrication of blased gratings by area-coded effective medium structures*, Opt. Comm., 266 (2006), pp. 697–703.

[13] R. HIPTMAIER, *Coupling of finite elements and boundary elements in electromagnetic scattering*, SIAM J. Numer. Anal., 41 (2003), pp. 919–944.

[14] G. HU AND A. RATHSFELD, *Scattering of time-harmonic electromagnetic plane waves by perfectly conducting diffraction gratings*, IMA J. Appl. Math., in press, (2014), doi:10.1093/imamat/hxt054.

- [15] M. HUBER, J. SCHOEBERL, A. SINWEL, AND S. ZAGLMAYR, *Simulation of diffraction in periodic media with a coupled finite element and plane wave approach*, SIAM J. Sci. Comput., 31 (2009), pp. 1500–1517.
- [16] R. KRESS, *Linear Integral Equations*, Springer, New York, 1989.
- [17] P. MONK, *Finite Element Method for Maxwell's Equations*, Oxford University Press, Oxford, 2003.
- [18] J. C. NÉDÉLEC AND F. STARLING, *Integral equation methods in a quasi-periodic diffraction problem for the time-harmonic Maxwell equation*, SIAM J. Math. Anal., 22 (1991), pp. 1679–1701.
- [19] M. NEVIÈRE AND E. POPOV, *Light Propagation in Periodic Media*, Marcel Dekker, Basel, 2003.
- [20] J. NITSCHKE, *Über ein Variationsprinzip zur Lösung von Dirichlet-Problemen bei Verwendung von Teilräumen, die keinen Randbedingungen unterworfen sind*, Abh. Math. Sem. Univ. Hamburg, 36 (1970), pp. 9–15.
- [21] S. PRÖSSDORF AND B. SILBERMANN, *Numerical Analysis for Integral and Related Operator Equations*, Akademie-Verlag, Berlin, 1991.
- [22] ———, *Projektionsverfahren und die näherungsweise Lösung singulärer Gleichungen*, Teubner, Leipzig, 1977.
- [23] A. RATHSFELD, *Shape derivatives for the scattering by biperiodic gratings*, Appl. Numer. Math., 72 (2013), pp. 19–32.
- [24] A. SCHAEDEL, L. ZSCHIEDRICH, S. BURGER, R. KLOSE, AND F. SCHMIDT, *Domain decomposition method for Maxwell's equations: scattering of periodic structures*, J. Comput. Phys., 226 (2007), pp. 477–493.
- [25] G. SCHMIDT, *On the diffraction by biperiodic anisotropic structures*, Appl. Anal., 82 (2003), pp. 75–92.
- [26] ———, *Electromagnetic scattering by periodic structures*, J. Math. Sci. (N. Y.), 124 (2004), pp. 5390–5405.
- [27] R. STERNBERG, *Mortaring by a method of J. A. Nitsche*, in Computational Mechanics. New trends and Applications, S. Idelsohn, E. Onate, and E. Dvorkin, eds., CIMNE, Barcelona, Spain, 1988.