



**HAL**  
open science

# Ultra-weak variational formulation for heterogeneous maxwell problem in the context of high performance computing

Sébastien Pernet, Margot Sirdey, Sébastien Tordeux

► **To cite this version:**

Sébastien Pernet, Margot Sirdey, Sébastien Tordeux. Ultra-weak variational formulation for heterogeneous maxwell problem in the context of high performance computing. ESAIM: Proceedings and Surveys, 2023, 75, pp.96-121. 10.1051/proc/202375096 . hal-03642116v2

**HAL Id: hal-03642116**

**<https://hal.science/hal-03642116v2>**

Submitted on 20 Sep 2023

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

# ULTRA-WEAK VARIATIONAL FORMULATION FOR HETEROGENEOUS MAXWELL PROBLEM IN THE CONTEXT OF HIGH PERFORMANCE COMPUTING

SÉBASTIEN PERNET<sup>1</sup>, MARGOT SIRDEY<sup>2</sup> AND SÉBASTIEN TORDEUX<sup>3</sup>

**Abstract.** Electromagnetic simulations on large domains require a huge memory consumption. Domain decomposition methods, based on Trefftz methods, could be an answer to this issue. In this paper, we associate to heterogeneous three-dimensional Maxwell equations one variational formulation which can be obtained either by upwind fluxes or Riemann traces. We associate to this variational formulation an iterative Trefftz Krylov solver. The poor conditioning due to the use of plane wave basis functions is bypassed thanks to a compression strategy. Moreover, the developed iterative solver is accelerated thanks to a left preconditioner. The considered numerical cases illustrate the performance of this basis reduction, which leads to the consideration of an industrial case of more than 750 millions of degrees of freedom.

## INTRODUCTION

Electromagnetic waves are present in a wide panel of applications, such as transports, high-technology or medicine. Their numerical modeling remains a challenging topic when considering complex industrial geometries. Some numerical methods can simulate accurately electromagnetic fields such as Finite Element Methods (FEM) [29, 37, 42], or Finite Difference (FD) [54]. However, they face two main issues: the numerical dispersion and the memory limitation.

First, numerical dispersion, see [1, 2, 35, 36], appears when the size of the domain is large with respect to the wavelength. In such cases, wave phenomena are oscillating and non dissipative. Then, the number of discretisation points per wavelength increases with the size of the domain to accurately approximate the wave phase. Many methods may partially counter this problem, such as high-order FEM [13, 44], polynomial Discontinuous Galerkin (DG) [14, 19, 22, 28, 33, 59], Boundary Element Methods (BEM) [6, 45, 53, 58], integral equation collocation method [5], or Trefftz type methods [31, 49]. The latter are at main interest in this paper. They take into account local basis functions adapted to the electromagnetic phenomena such that they reduce numerical dispersion. In three dimensions, these methods are associated to large linear systems which are extremely costly to invert when using LU factorization. It is then needed to choose the method depending on the studied geometry. Condensation methods, such as Hybrid DG (HDG) methods [12, 43, 46] have been introduced to reduce the number of degrees of freedom. Nevertheless, these optimisations are not sufficient to simulate electromagnetic waves on domains whose size is larger than 30 wavelengths. For BEM, this issue has been solved thanks to compression methods like Fast Multipole [16, 17] and Adaptive Cross Approximation

---

<sup>1</sup> DTIS, ONERA

<sup>2</sup> Makutu, EPC INRIA, e2s-UPPA, LMAP UMR CNRS 5142 - DTIS, ONERA

<sup>3</sup> Makutu, EPC INRIA, e2s-UPPA, LMAP UMR CNRS 5142

methods [4, 38], for example. But, implementation of BEM in the case of heterogeneous structures is a really hard task.

When dealing with larger heterogeneous geometries, the indefinite property of the matrix associated to most volumic methods prevents the use of algebraic iterative inversion. Many authors propose Domain Decomposition Methods (DDM), see [20, 57] for example, but DDM are often difficult to carry out in the context of High-Parallel Computing (HPC) architecture. The Ultra-Weak Variational Formulation (UWVF) methods [7, 10, 18, 34, 40, 56], developed by Cessenat-Després in [9], are conversely naturally adapted to domain decomposition since they involve definite positive matrices. Consequently, UWVF is easily implementable in HPC framework. Trefftz methods often use plane wave local basis functions, see [24, 26, 27, 41], but we could have chose spherical Bessel functions [39]. Practically, the advantage of using such basis results in an easy implementation. But plane wave basis functions difficultly approximate high-complexity electromagnetic waves such as point effects at corners or trapped waves in elaborated heterogeneous configurations. Moreover, they often provide ill-conditioned linear systems, due to rounding errors caused by plane waves linear-dependency, see [15]. This may conduct to a lack of accuracy for the numerical solution. To avoid such issue, previous works have developed strategies, for which we propose some improvements here, see [15, 40]. For homogeneous cases, the convergence of the method is well established, see [30].

Here, we aim at developing an iterative Trefftz solver using few memory and leading to an accurate numerical solution. The developed method is called the heterogeneous Cessenat-Després UWVF. This numerical method resorts to Maxwell equations to model electromagnetic waves. In the present paper we study a dimensionless Maxwell problem which needs to be constructed. We recall the general Maxwell formulas in absence of charges and currents and for an isotropic linear medium

$$\begin{cases} \nabla \cdot \mathbf{d} = 0, & \nabla \times \mathbf{e} = -\frac{\partial \mathbf{b}}{\partial t}, & \mathbf{d} = \varepsilon_0 \varepsilon_r \mathbf{e}, \\ \nabla \cdot \mathbf{b} = 0, & \nabla \times \mathbf{h} = \frac{\partial \mathbf{d}}{\partial t}, & \mathbf{b} = \mu_0 \mu_r \mathbf{h}, \end{cases}$$

where  $\varepsilon_0$  (*resp.*  $\varepsilon_r$ ) and  $\mu_0$  (*resp.*  $\mu_r$ ) are the permittivity (*resp.* relative permittivity) and the permeability (*resp.* relative permeability) of the vacuum (*resp.* of the medium). This system involves the electric and magnetic field intensities  $\mathbf{e}$  and  $\mathbf{h}$ , the electric displacement  $\mathbf{d}$  and the magnetic induction  $\mathbf{b}$ . We suppose they are all time-harmonic fields. Therefore, they can be represented by four complex valued normalised functions  $\mathbf{E}$ ,  $\mathbf{H}$ ,  $\mathbf{D}$  and  $\mathbf{B}$  which are associated to their normalisation amplitudes  $e_0$ ,  $d_0 = \varepsilon_0 e_0$ ,  $h_0 = \sqrt{\varepsilon_0/\mu_0} e_0$  and  $b_0 = \sqrt{\varepsilon_0 \mu_0} e_0$ ,

$$\begin{cases} \mathbf{e}(\mathbf{x}, t) = e_0 \mathcal{R}(e^{i\omega t} \mathbf{E}(\mathbf{x})), & \mathbf{h}(\mathbf{x}, t) = h_0 \mathcal{R}(e^{i\omega t} \mathbf{H}(\mathbf{x})), \\ \mathbf{d}(\mathbf{x}, t) = d_0 \mathcal{R}(e^{i\omega t} \mathbf{D}(\mathbf{x})), & \mathbf{b}(\mathbf{x}, t) = b_0 \mathcal{R}(e^{i\omega t} \mathbf{B}(\mathbf{x})), \end{cases}$$

where  $\omega$  is the angular frequency that accounts for time-harmonic dependency. We define the wavenumber  $k_0 = \omega/c_0$ , where  $c_0 = (\varepsilon_0 \mu_0)^{-\frac{1}{2}}$  is the velocity in the vacuum. We get the following Maxwell normalised problem defined on a connex Lipschitz domain  $\Omega \subset \mathbb{R}^3$

$$\nabla \times \mathbf{H} = ik_0 \varepsilon_r \mathbf{E} \quad \text{and} \quad \nabla \times \mathbf{E} = -ik_0 \mu_r \mathbf{H}, \quad \text{on } \Omega, \quad (1)$$

where both fields are in the space  $H(\nabla \times, \Omega)$  defined by

$$H(\nabla \times, \Omega) := \left\{ \mathbf{u} : \Omega \rightarrow \mathbb{C}^3, \int_{\Omega} |\mathbf{u}|^2 \, d\mathbf{x} < \infty, \int_{\Omega} |\nabla \times \mathbf{u}|^2 \, d\mathbf{x} < \infty \right\}.$$

The unknown of the problem is denoted by  $\mathbb{E} := (\mathbf{E}, \mathbf{H})$ . We complete by an impedance boundary condition defined on the domain boundary  $\partial\Omega$

$$\gamma_t \mathbf{E} + Z_{\partial\Omega} \mathbf{n} \times \gamma_t \mathbf{H} = \mathbf{g} \quad \text{on } \partial\Omega,$$

where we have the following definitions and notations:

- $\mathbf{n} \in \mathbb{R}^3$  is the outward unit normal to  $\partial\Omega$ ,
- $Z_{\partial\Omega} : \partial\Omega \rightarrow \mathbb{R}^+$  is a piecewise constant function, with a strictly non-positive real part,
- $\gamma_t \mathbf{E} = \mathbf{E} - (\mathbf{E} \cdot \mathbf{n}) \mathbf{n} = (\mathbf{n} \times \mathbf{E}) \times \mathbf{n}$  is the tangential component of the electric field,
- $\gamma_t \mathbf{H} = \mathbf{H} - (\mathbf{H} \cdot \mathbf{n}) \mathbf{n} = (\mathbf{n} \times \mathbf{H}) \times \mathbf{n}$  is the tangential component of the magnetic field,
- $\varepsilon_r$  and  $\mu_r$  are piecewise constant functions defined on a partition of  $\Omega$ . The material parameters are assumed to be constant over each mesh cell.
- $\mathbf{g} : \partial\Omega \rightarrow \mathbb{C}^3$  is a purely tangential field in the following functional space

$$L_t^2(\partial\Omega) := \left\{ \mathbf{u} \in (L^2(\partial\Omega))^3, \quad \mathbf{u} \cdot \mathbf{n} = 0 \right\}.$$

There exists a unique solution of problem (1), see [42]. The main goal of this paper is to develop a numerical method approximating accurately the solution to the Maxwell problem, with a low memory cost. In Section 1, we introduce the mesh associated to the domain  $\Omega$  and its properties. Afterwards, we design an heterogeneous plane wave mixed DG Trefftz variational formulation based on upwind fluxes; *mixed* in the sense that it involves both the electric and the magnetic fields. The well-posedness of the associated problem is proved thanks to a weak coercivity property. The well-posedness of the associated problem relies on the positivity of the variational formulation. In Section 2, we bring to the fore the equivalence between the mixed DG Trefftz formulation and the heterogeneous Cessenat-Després UWVF. Then, the Trefftz UWVF problem is discretised, leading to a matricial linear system, which can be inverted with LU solver. However, due to the memory challenge, we opt for an iterative Trefftz UWVF algorithm based on a singular regular decomposition of the matrix. In Section 4, we associate to the constructed iterative scheme a Krylov type method. The latter Krylov method, combined with a preconditioning strategy, is proved to be convergent thanks to Galerkin theory. Thereafter, we set up a compression strategy reducing the number of plane wave basis functions in the discrete Galerkin space, leading to a memory cost reduction. Finally, a global preconditioning aims at accelerating the Trefftz Krylov UWVF solver.

## 1. MIXED DISCONTINUOUS GALERKIN TREFFTZ FORMULATION

In this section, we devise a mixed DG Trefftz problem based on upwind fluxes. This goes through the introduction of Trefftz spaces and the reciprocity formula. Finally, we establish a weak coercivity property and show the well-posedness of the associated variational formulation.

### 1.1. Mesh properties and Trefftz spaces

We consider a three-dimensional mesh made of non-overlapping polyhedral elements  $T$  meshing the computational domain  $\Omega$ . Indeed, we choose a mesh that follows the partitions  $\Omega_p$ ,  $1 \leq p \leq P$ , of  $\Omega$ , in which  $\varepsilon_r$  and  $\mu_r$  are constant. It means that there exists a unique  $1 \leq p_0 \leq P$  such that  $T \subset \Omega_{p_0}$ , leading to piecewise constant functions  $\varepsilon_r$  and  $\mu_r$ . We denote by  $\mathcal{T}$  the set of elements  $T$  and by  $\mathcal{F}$  the set of faces  $F$  in  $\Omega$ . We define the following sets of faces:

- The set  $\mathcal{F}_{\text{int}}$  of interior faces

$$\mathcal{F}_{\text{int}} := \{ \partial T \cap \partial K : T, K \in \mathcal{T} \text{ with } T \neq K \text{ and } \text{area}(\partial T \cap \partial K) \neq 0 \},$$

where  $\text{area}(I)$  refers to the area of the face  $I$  (that is zero for edges and vertices).

- The set  $\mathcal{F}_{\text{ext}}$  of exterior faces

$$\mathcal{F}_{\text{ext}} := \{\partial T \cap \partial\Omega : T \in \mathcal{T} \text{ and } \text{area}(\partial T \cap \partial\Omega) \neq 0\}.$$

- The set  $\mathcal{F}_T$  of faces associated to an element  $T \in \mathcal{T}$

$$\mathcal{F}_T := \{F \in \mathcal{F}_{\text{int}} \cup \mathcal{F}_{\text{ext}} : \text{area}(F \cap \partial T) \neq 0\}.$$

We denote by  $\mathbb{X}_T$  the local Trefftz space defined for  $T \in \mathcal{T}$  as the set of functions

$$\mathbb{E}^T := (\mathbf{E}^T, \mathbf{H}^T) \in H(\nabla \times, T) \times H(\nabla \times, T) \text{ satisfying} \quad (2a)$$

$$\nabla \times \mathbf{H}^T = ik_0 \varepsilon_r \mathbf{E}^T \quad \text{and} \quad \nabla \times \mathbf{E}^T = -ik_0 \mu_r \mathbf{H}^T, \quad \text{on } T, \quad (2b)$$

$$\gamma_t \mathbf{E}^T \in L_t^2(\partial T) \text{ and } \gamma_t \mathbf{H}^T \in L_t^2(\partial T), \quad (2c)$$

where  $\mathbf{E}^T$  and  $\mathbf{H}^T$  are the restrictions of the electric field  $\mathbf{E}$  and the magnetic field  $\mathbf{H}$  on the element  $T$ .

**Remark 1.1.** *The elliptic regularity theory seems to restrict the numerical method of this paper to regular solutions, ie  $\mathbf{E}^T, \mathbf{H}^T \in H^{\frac{3}{2}}(T)$ . However, we did not face any difficulty to apply the present method to less regular solutions.*

The space  $\mathbb{X}_T$  is equipped with the inner product parametrised by  $\alpha > 0$  and  $\beta > 0$

$$(\mathbb{E}^T, \mathbb{E}'^T)_{\mathbb{X}_T} := \int_{\partial T} \left( \alpha \gamma_t \mathbf{E}^T \cdot \overline{\gamma_t \mathbf{E}'^T} + \beta \gamma_t \mathbf{H}^T \cdot \overline{\gamma_t \mathbf{H}'^T} \right) ds_{\mathbf{x}},$$

where  $ds_{\mathbf{x}}$  is the Lebesgue measure associated to surfaces of the three-dimensional mesh.

**Remark 1.2.** *Due to the unique continuation theorem, the pairing  $(\cdot, \cdot)_{\mathbb{X}_T}$  is an inner product which means*

$$(\mathbb{E}^T, \mathbb{E}^T)_{\mathbb{X}_T} = 0 \quad \implies \quad \gamma_t \mathbf{E}^T = 0 \text{ and } \gamma_t \mathbf{H}^T = 0 \text{ on } \partial T \quad \xrightarrow{(2b)} \quad \mathbb{E}^T = (\mathbf{E}^T, \mathbf{H}^T) = 0 \text{ in } T.$$

It leads to the global discontinuous Trefftz space  $\mathbb{X}_{\mathcal{T}}$  defined element by element by

$$\mathbb{X}_{\mathcal{T}} := \prod_{T \in \mathcal{T}} \mathbb{X}_T, \text{ equipped with the norm } (\mathbb{E}, \mathbb{E}')_{\mathbb{X}_{\mathcal{T}}} := \sum_{T \in \mathcal{T}} (\mathbb{E}^T, \mathbb{E}'^T)_{\mathbb{X}_T}.$$

**Remark 1.3.** *Every function  $\mathbb{E} = (\mathbb{E}^T)_{T \in \mathcal{T}} \in \mathbb{X}_{\mathcal{T}}$  is associated to a function defined on  $\Omega$  whose restriction to  $T$  is  $\mathbb{E}^T$ . Each  $\mathbb{E}^T = (\mathbf{E}^T, \mathbf{H}^T)$  satisfies (2b) and  $\mathbb{E}$  is generally discontinuous across faces.*

## 1.2. Trefftz reciprocity formula

We introduce the reciprocity formula, also called the virtual work formula, which involves the restrictions on  $T \in \mathcal{T}$  of the tangential component of the electric field and the tangential trace of the magnetic field

$$\gamma_t \mathbf{E}^T := (\mathbf{n}_T \times \mathbf{E}^T) \times \mathbf{n}_T \quad \text{and} \quad \gamma_{\times}^T \mathbf{H}^T := \mathbf{n}_T \times \gamma_t \mathbf{H}^T, \quad (3)$$

where  $\mathbf{n}_T \in \mathbb{R}^3$  is the outward unit normal to  $\partial T$ .

**Proposition 1.1.** *For  $\mathbb{E} \in \mathbb{X}_{\mathcal{T}}$  and  $\mathbb{E}' \in \mathbb{X}_{\mathcal{T}}$ , the global reciprocity formula is*

$$r(\mathbb{E}, \mathbb{E}') := \sum_{T \in \mathcal{T}} r_T(\mathbb{E}, \mathbb{E}') \quad \text{with} \quad r_T(\mathbb{E}, \mathbb{E}') := \int_{\partial T} \left( \gamma_{\times}^T \mathbf{H}^T \cdot \overline{\gamma_t \mathbf{E}'^T} + \gamma_t \mathbf{E}^T \cdot \overline{\gamma_{\times}^T \mathbf{H}'^T} \right) ds_{\mathbf{x}} = 0. \quad (4)$$

*Proof.* The reciprocity formula takes the following form, with  $\mathbb{E}^T = (\mathbf{E}^T, \mathbf{H}^T)$  and  $\mathbb{E}'^T = (\mathbf{E}'^T, \mathbf{H}'^T)$  in  $\mathbb{X}_T$ ,

$$W_T = ik_0 \int_T \varepsilon_r \mathbf{E}^T \cdot \overline{\mathbf{E}'^T} + \mu_r \mathbf{H}^T \cdot \overline{\mathbf{H}'^T} dT, \quad \text{for } T \in \mathcal{T}.$$

Since both the solution  $\mathbb{E}^T \in \mathbb{X}_T$  and the test function  $\mathbb{E}'^T \in \mathbb{X}_T$  satisfy (2), we have

$$W_T = \int_T \nabla \times \mathbf{H}^T \cdot \overline{\mathbf{E}'^T} - \nabla \times \mathbf{E}^T \cdot \overline{\mathbf{H}'^T} dT = \int_T \mathbf{H}^T \cdot \overline{\nabla \times \mathbf{E}'^T} - \mathbf{E}^T \cdot \overline{\nabla \times \mathbf{H}'^T} dT.$$

By subtracting these two last expressions of  $W_T$ , we have

$$\int_T \nabla \times \mathbf{H}^T \cdot \overline{\mathbf{E}'^T} - \mathbf{H}^T \cdot \overline{\nabla \times \mathbf{E}'^T} dT + \int_T \mathbf{E}^T \cdot \overline{\nabla \times \mathbf{H}'^T} - \nabla \times \mathbf{E}^T \cdot \overline{\mathbf{H}'^T} dT = 0.$$

Due to the Stokes formula (see [42]), we get the local reciprocity formula

$$r_T(\mathbb{E}, \mathbb{E}') = \int_{\partial T} (\mathbf{n}_T \times \mathbf{H}^T) \cdot \overline{\mathbf{E}'^T} + \mathbf{E}^T \cdot (\mathbf{n}_T \times \overline{\mathbf{H}'^T}) ds_{\mathbf{x}} = 0, \quad \text{for } \mathbb{E} \text{ and } \mathbb{E}' \text{ in } \mathbb{X}_{\mathcal{T}}.$$

We end the proof by using definitions (3). □

**Remark 1.4.** We remark that  $\partial T$  can be decomposed into faces  $F \in \mathcal{F}_T$  leading to

$$\sum_{T \in \mathcal{T}} \int_{\partial T} f^T ds_{\mathbf{x}} = \sum_{T \in \mathcal{T}} \sum_{F \in \mathcal{F}_T} \int_F f^T ds_{\mathbf{x}}, \quad (5)$$

for any piecewise continuous regular function  $f$  whose restriction on one element  $T \in \mathcal{T}$  is denoted by  $f^T$ . The reciprocity formula (4) is then equivalent to

$$r_T(\mathbb{E}, \mathbb{E}') := \sum_{F \in \mathcal{F}_T} \int_F \left( \gamma_{\times}^T \mathbf{H}^T \cdot \overline{\gamma_t \mathbf{E}'^T} + \gamma_t \mathbf{E}^T \cdot \overline{\gamma_{\times}^T \mathbf{H}'^T} \right) ds_{\mathbf{x}} = 0.$$

### 1.3. Mixed Discontinuous Galerkin Trefftz variational formulation

Following the same approach than for polynomial mixed DG methods, see [22, 32], the mixed DG Trefftz formulation is deduced by inserting upwind fluxes, see [47], into the reciprocity formula (4).

**Problem 1.** Find  $\mathbb{E} \in \mathbb{X}_{\mathcal{T}}$  such that for all  $\mathbb{E}' \in \mathbb{X}_{\mathcal{T}}$

$$\sum_{T \in \mathcal{T}} \sum_{F \in \mathcal{F}_T} \int_F \left( \widehat{\gamma_t \mathbf{E}} \cdot \gamma_{\times}^T \overline{\mathbf{H}'^T} + \widehat{\gamma_{\times}^T \mathbf{H}} \cdot \overline{\gamma_t \mathbf{E}'^T} \right) ds_{\mathbf{x}} = 0, \quad (6)$$

with interior upwind fluxes for  $F \in \mathcal{F}_{\text{int}}$  separating two elements  $T$  and  $K$  defined as

$$\left( \widehat{\gamma_t \mathbf{E}} \right)_{|F} := \frac{Z_K}{Z_K + Z_T} \gamma_{\text{out}}^T \mathbb{E}^T + \frac{Z_T}{Z_K + Z_T} \gamma_{\text{out}}^K \mathbb{E}^K \quad \text{and} \quad \left( \widehat{\gamma_{\times}^T \mathbf{H}} \right)_{|F} := -\frac{1}{Z_K + Z_T} \gamma_{\text{out}}^T \mathbb{E}^T + \frac{1}{Z_K + Z_T} \gamma_{\text{out}}^K \mathbb{E}^K, \quad (7)$$

and boundary upwind fluxes for  $F \in \mathcal{F}_{\text{ext}}$  defined as

$$\left( \widehat{\gamma_t \mathbf{E}} \right)_{|F} := \frac{Z_{\partial\Omega}}{Z_{\partial\Omega} + Z_T} \gamma_{\text{out}}^T \mathbb{E}^T + \frac{Z_T}{Z_{\partial\Omega} + Z_T} \mathbf{g}^T \quad \text{and} \quad \left( \widehat{\gamma_{\times}^T \mathbf{H}} \right)_{|F} := -\frac{1}{Z_{\partial\Omega} + Z_T} \gamma_{\text{out}}^T \mathbb{E}^T + \frac{1}{Z_{\partial\Omega} + Z_T} \mathbf{g}^T, \quad (8)$$

where the general incoming and outgoing wave trace operators  $\gamma_{\text{in/out}}^T : \mathbb{X}_T \rightarrow L_t^2(\partial T)$  are defined by

$$\gamma_{\text{in}}^T \mathbb{E}^T := \gamma_t \mathbf{E}^T + Z_T \gamma_{\times}^T \mathbf{H}^T \quad \text{and} \quad \gamma_{\text{out}}^T \mathbb{E}^T := \gamma_t \mathbf{E}^T - Z_T \gamma_{\times}^T \mathbf{H}^T. \quad (9)$$

**Remark 1.5.** For this particular choice of numerical fluxes (7) and (8), the exact solution satisfies

$$\widehat{(\gamma_t \mathbf{E})}_{|F} = (\gamma_t \mathbf{E}^T)_{|F} = (\gamma_t \mathbf{E}^K)_{|F} \quad \text{and} \quad \widehat{(\gamma_{\times}^T \mathbf{H})}_{|F} = (\gamma_{\times}^T \mathbf{H}^T)_{|F} = -(\gamma_{\times}^K \mathbf{H}^K)_{|F}.$$

Therefore, the variational formulation of Problem 1 encodes the continuity of its solution.

**Remark 1.6.** An exterior face  $F \in \mathcal{F}_{\text{ext}}$  can be seen as a face interfacing an element  $T$  and a virtual element exterior to the domain. More precisely, denoting by  $Z_{\partial\Omega} = Z_K$  and  $\mathbf{g}^T = \gamma_{\text{out}}^K \mathbb{E}^K$ , (8) can be seen as (7).

Problem 1 can be interpreted as a variational problem and leads to the following proposition.

**Proposition 1.2.** Problem 1 can be written as

$$\text{Find } \mathbb{E} \in \mathbb{X}_{\mathcal{T}} \text{ such that for all } \mathbb{E}' \in \mathbb{X}_{\mathcal{T}} \text{ we have } a(\mathbb{E}, \mathbb{E}') = \ell(\mathbb{E}'), \quad (10)$$

where the sesquilinear form  $a$  and the linear form  $\ell$  are given by

$$a(\mathbb{E}, \mathbb{E}') := \sum_{T \in \mathcal{T}} \sum_{F \in \mathcal{F}_T} \int_F \left( a_{T,F}^{EE'} + a_{T,F}^{EH'} + a_{T,F}^{HE'} + a_{T,F}^{HH'} \right) ds_{\mathbf{x}}, \quad (11a)$$

$$\ell(\mathbb{E}') := \sum_{T \in \mathcal{T}} \sum_{F \in \mathcal{F}_T \cap \mathcal{F}_{\text{ext}}} \int_F \left( \ell_{T,F}^{E'} + \ell_{T,F}^{H'} \right) ds_{\mathbf{x}}, \quad (11b)$$

where we have if  $F \in \mathcal{F}_{\text{int}}$  separating two neighboring elements  $T$  and  $K$

$$\begin{cases} a_{T,F}^{EE'} := \frac{1}{Z_T + Z_K} (\gamma_t \mathbf{E}^T - \gamma_t \mathbf{E}^K) \cdot \overline{\gamma_t \mathbf{E}'^T}, & a_{T,F}^{EH'} := \frac{Z_T}{Z_T + Z_K} (\gamma_t \mathbf{E}^T - \gamma_t \mathbf{E}^K) \cdot \overline{\gamma_{\times}^T \mathbf{H}'^T}, \\ a_{T,F}^{HE'} := \frac{Z_K}{Z_T + Z_K} (\gamma_{\times}^T \mathbf{H}^T - \gamma_{\times}^T \mathbf{H}^K) \cdot \overline{\gamma_t \mathbf{E}'^T}, & a_{T,F}^{HH'} := \frac{Z_T Z_K}{Z_T + Z_K} (\gamma_{\times}^T \mathbf{H}^T - \gamma_{\times}^T \mathbf{H}^K) \cdot \overline{\gamma_{\times}^T \mathbf{H}'^T}, \end{cases} \quad (12a)$$

and if  $F \in \mathcal{F}_{\text{ext}}$

$$\begin{cases} a_{T,F}^{EE'} := \frac{1}{Z_T + Z_{\partial\Omega}} \gamma_t \mathbf{E}^T \cdot \overline{\gamma_t \mathbf{E}'^T}, & a_{T,F}^{EH'} := \frac{Z_T}{Z_T + Z_{\partial\Omega}} \gamma_t \mathbf{E}^T \cdot \overline{\gamma_{\times}^T \mathbf{H}'^T}, \\ a_{T,F}^{HE'} := \frac{Z_{\partial\Omega}}{Z_T + Z_{\partial\Omega}} \gamma_{\times}^T \mathbf{H}^T \cdot \overline{\gamma_t \mathbf{E}'^T}, & a_{T,F}^{HH'} := \frac{Z_T Z_{\partial\Omega}}{Z_T + Z_{\partial\Omega}} \gamma_{\times}^T \mathbf{H}^T \cdot \overline{\gamma_{\times}^T \mathbf{H}'^T}, \\ \ell_{T,F}^{E'} := \frac{1}{Z_T + Z_{\partial\Omega}} \mathbf{g}^T \cdot \overline{\gamma_t \mathbf{E}'^T}, & \ell_{T,F}^{H'} := \frac{Z_T}{Z_T + Z_{\partial\Omega}} \mathbf{g}^T \cdot \overline{\gamma_{\times}^T \mathbf{H}'^T}. \end{cases} \quad (12b)$$

*Proof.* Problem 1 reads

$$\sum_{T \in \mathcal{T}} \sum_{F \in \mathcal{F}_T} \int_F \mathcal{I}_{T,F} ds_{\mathbf{x}} = 0, \quad \text{with} \quad \mathcal{I}_{T,F} := \widehat{\gamma_t \mathbf{E}} \cdot \overline{\gamma_{\times}^T \mathbf{H}'^T} + \widehat{\gamma_{\times}^T \mathbf{H}} \cdot \overline{\gamma_t \mathbf{E}'^T}.$$

(a) If  $F \in \mathcal{F}_{\text{int}}$ , using the definition of the numerical traces (7), we have

$$\mathcal{I}_{T,F} = \left( \frac{Z_K}{Z_K + Z_T} \gamma_{\text{out}}^T \mathbb{E}^T + \frac{Z_T}{Z_T + Z_K} \gamma_{\text{out}}^K \mathbb{E}^K \right) \cdot \overline{\gamma_{\times}^T \mathbf{H}'^T} + \left( -\frac{1}{Z_K + Z_T} \gamma_{\text{out}}^T \mathbb{E}^T + \frac{1}{Z_K + Z_T} \gamma_{\text{out}}^K \mathbb{E}^K \right) \cdot \overline{\gamma_t \mathbf{E}'^T}.$$

Taking into account the definition of the outgoing trace operator in (9), we get

$$\begin{aligned} \mathcal{I}_{T,F} = & \left[ \frac{Z_K}{Z_K + Z_T} (\gamma_t \mathbf{E}^T - Z_T \gamma_\times^T \mathbf{H}^T) + \frac{Z_T}{Z_T + Z_K} (\gamma_t \mathbf{E}^K + Z_K \gamma_\times^T \mathbf{H}^K) \right] \cdot \overline{\gamma_\times^T \mathbf{H}^T} \\ & + \left[ -\frac{1}{Z_K + Z_T} (\gamma_t \mathbf{E}^T - Z_T \gamma_\times^T \mathbf{H}^T) + \frac{1}{Z_K + Z_T} (\gamma_t \mathbf{E}^K + Z_K \gamma_\times^T \mathbf{H}^K) \right] \cdot \overline{\gamma_t \mathbf{E}^T}. \end{aligned}$$

Expanding, we have

$$\mathcal{I}_{T,F} = \frac{(Z_K \gamma_t \mathbf{E}^T + Z_T \gamma_t \mathbf{E}^K)}{Z_K + Z_T} \cdot \overline{\gamma_\times^T \mathbf{H}^T} + \frac{(Z_T \gamma_\times^T \mathbf{H}^T + Z_K \gamma_\times^T \mathbf{H}^K)}{Z_K + Z_T} \cdot \overline{\gamma_t \mathbf{E}^T} - a_{T,F}^{EE'} - a_{T,F}^{HH'}.$$

We then remark that  $\frac{Z_K}{Z_K + Z_T} + \frac{Z_T}{Z_K + Z_T} = 1$  and  $\frac{Z_{\partial\Omega}}{Z_{\partial\Omega} + Z_T} + \frac{Z_T}{Z_{\partial\Omega} + Z_T} = 1$ . This leads to

$$\mathcal{I}_{T,F} = \gamma_t \mathbf{E}^T \cdot \overline{\gamma_\times^T \mathbf{H}^T} + \gamma_\times^T \mathbf{H}^T \cdot \overline{\gamma_t \mathbf{E}^T} - a_{T,F}^{EE'} - a_{T,F}^{EH'} - a_{T,F}^{HE'} - a_{T,F}^{HH'}. \quad (13)$$

(b) If  $F \in \mathcal{F}_{\text{ext}}$ , we have in the same way

$$\mathcal{I}_{T,F} = \gamma_t \mathbf{E}^T \cdot \overline{\gamma_\times^T \mathbf{H}^T} + \gamma_\times^T \mathbf{H}^T \cdot \overline{\gamma_t \mathbf{E}^T} + \ell_{T,F}^{E'} + \ell_{T,F}^{H'} - a_{T,F}^{EE'} - a_{T,F}^{EH'} - a_{T,F}^{HE'} - a_{T,F}^{HH'}. \quad (14)$$

Proposition 1.2 follows from (13), (14) and the reciprocity formula (4).  $\square$

#### 1.4. Positivity and injectivity properties

The well-posedness of Problem 1 relies on the positivity of the sesquilinear form  $a$  defined by (11a). It implies the study of the problem stability through interfaces  $F \in \mathcal{F}_{\text{int}}$  involving the jump of the tangential trace (*resp.* component) of the electric (*resp.* magnetic) field, respectively defined as

$$\llbracket \gamma_t \mathbf{E} \rrbracket_F := \gamma_t \mathbf{E}^T - \gamma_t \mathbf{E}^K \quad \text{and} \quad \llbracket \gamma_\times \mathbf{H} \rrbracket_F := \gamma_\times^T \mathbf{H}^T - \gamma_\times^T \mathbf{H}^K. \quad (15)$$

**Proposition 1.3.** *The sesquilinear form  $a$  is positive since  $\Re(a(\mathbb{E}, \mathbb{E})) \geq 0$ . We define the DG norm as*

$$\|\mathbb{E}\|_{\text{DG}} = \sqrt{\Re(a(\mathbb{E}, \mathbb{E}))}, \quad \text{for all } \mathbb{E} \in \mathbb{X}_{\mathcal{T}}, \quad (16)$$

for which we have the following properties

(i)  $\|\mathbb{E}\|_{\text{DG}}^2 = \|\mathbb{E}\|_{\text{int}}^2 + \|\mathbb{E}\|_{\text{ext}}^2$  with

$$\begin{aligned} \|\mathbb{E}\|_{\text{int}}^2 &= \sum_{F \in \mathcal{F}_{\text{int}}} \int_F \left( \frac{1}{Z_T + Z_K} \llbracket \gamma_t \mathbf{E} \rrbracket_F \cdot \overline{\llbracket \gamma_t \mathbf{E} \rrbracket_F} + \frac{Z_T Z_K}{Z_T + Z_K} \llbracket \gamma_\times \mathbf{H} \rrbracket_F \cdot \overline{\llbracket \gamma_\times \mathbf{H} \rrbracket_F} \right) \text{d}\mathbf{s}_{\mathbf{x}}, \\ \|\mathbb{E}\|_{\text{ext}}^2 &= \sum_{F \in \mathcal{F}_{\text{ext}}} \int_F \left( \frac{1}{Z_T + Z_{\partial\Omega}} \gamma_t \mathbf{E} \cdot \overline{\gamma_t \mathbf{E}} + \frac{Z_T Z_{\partial\Omega}}{Z_T + Z_{\partial\Omega}} \gamma_t \mathbf{H} \cdot \overline{\gamma_t \mathbf{H}} \right) \text{d}\mathbf{s}_{\mathbf{x}}. \end{aligned}$$

(ii)  $a(\mathbb{E}, \mathbb{E}) = 0 \implies \mathbb{E} \equiv 0$ .

The proof of Proposition 1.3 resorts to the following remark.

**Remark 1.7.** *Each face  $F \in \mathcal{F}_{\text{int}}$  is involved twice in the sum (5) of Remark 1.4. Thus, we can decompose the integral on  $\partial T$  by using a "face point of view"*

$$\sum_{T \in \mathcal{T}} \int_{\partial T} f^T \text{d}\mathbf{s}_{\mathbf{x}} = \sum_{F \in \mathcal{F}_{\text{ext}}} \int_F f^T \text{d}\mathbf{s}_{\mathbf{x}} + \sum_{F \in \mathcal{F}_{\text{int}}} \int_F f^T + f^K \text{d}\mathbf{s}_{\mathbf{x}},$$



where, in the right-hand side,  $T$  (resp.  $T$  and  $K$ ) is one element (resp. are two elements) with one face  $F$ .

*Proof.* (i) Using (11a), we can write the real part of  $a$  as

$$\Re(a(\mathbb{E}, \mathbb{E})) = \sum_{T \in \mathcal{T}} \sum_{F \in \mathcal{F}_T} \int_F \Re \left( a_{T,F}^{EE} + a_{T,F}^{EH} + a_{T,F}^{HE} + a_{T,F}^{HH} \right) ds_{\mathbf{x}}.$$

(a) We first show that

$$\mathcal{I}^1 = \sum_{T \in \mathcal{T}} \sum_{F \in \mathcal{F}_T} \int_F \underbrace{\Re(a_{T,F}^{EH} + a_{T,F}^{HE})}_{:= \mathcal{I}_{T,F}^1} ds_{\mathbf{x}} = 0. \quad (17)$$

We evaluate  $\mathcal{I}_{T,F}^1$  by distinguishing the cases of an exterior face  $F \in \mathcal{F}_{\text{ext}}$  and an interior face  $F \in \mathcal{F}_{\text{int}}$ . For  $F \in \mathcal{F}_{\text{int}}$ , using the definitions of forms (12a) and taking into account that  $\Re(xy) = \Re(y\bar{x})$ , we have

$$\begin{aligned} \mathcal{I}_{T,F}^1 &= \Re \left( \frac{Z_T}{Z_T + Z_K} (\gamma_t \mathbf{E}^T - \gamma_t \mathbf{E}^K) \cdot \overline{\gamma_{\times}^T \mathbf{H}^T} + \frac{Z_K}{Z_T + Z_K} \gamma_t \mathbf{E}^T \cdot (\overline{\gamma_{\times}^T \mathbf{H}^T} - \overline{\gamma_{\times}^T \mathbf{H}^K}) \right), \\ &= \Re \left( \gamma_t \mathbf{E}^T \cdot \overline{\gamma_{\times}^T \mathbf{H}^T} - \frac{Z_T}{Z_T + Z_K} \gamma_t \mathbf{E}^K \cdot \overline{\gamma_{\times}^T \mathbf{H}^T} - \frac{Z_K}{Z_T + Z_K} \gamma_t \mathbf{E}^T \cdot \overline{\gamma_{\times}^T \mathbf{H}^K} \right). \end{aligned}$$

In a similar way, using (12b) for  $F \in \mathcal{F}_{\text{ext}}$ , we get

$$\mathcal{I}_{T,F}^1 = \Re \left( \frac{Z_T}{Z_T + Z_{\partial\Omega}} \gamma_t \mathbf{E}^T \cdot \overline{\gamma_{\times}^T \mathbf{H}^T} + \frac{Z_{\partial\Omega}}{Z_T + Z_{\partial\Omega}} \gamma_t \mathbf{E}^T \cdot \overline{\gamma_{\times}^T \mathbf{H}^T} \right) = \Re(\gamma_t \mathbf{E}^T \cdot \overline{\gamma_{\times}^T \mathbf{H}^T}).$$

Since  $\Re(r_T(\mathbb{E}, \mathbb{E})) = 2 \sum_{F \in \mathcal{F}_T} \int_F \Re(\gamma_t \mathbf{E}^T \cdot \overline{\gamma_{\times}^T \mathbf{H}^T})$ , see Remark 1.4, and summing over  $F \in \mathcal{F}_T$  we get

$$\sum_{F \in \mathcal{F}_T} \int_F \mathcal{I}_{T,F}^1 ds_{\mathbf{x}} = \frac{1}{2} \underbrace{\Re(r_T(\mathbb{E}, \mathbb{E}))}_{=0 \text{ see (4)}} - \sum_{F \in \mathcal{F}_T \cap \mathcal{F}_{\text{int}}} \int_F \Re \left( \frac{Z_T}{Z_T + Z_K} \gamma_t \mathbf{E}^K \cdot \overline{\gamma_{\times}^T \mathbf{H}^T} + \frac{Z_K}{Z_T + Z_K} \gamma_t \mathbf{E}^T \cdot \overline{\gamma_{\times}^T \mathbf{H}^K} \right) ds_{\mathbf{x}}.$$

Summing over elements, we then take the face point of view of Remark 1.7

$$\mathcal{I}^1 = - \sum_{F \in \mathcal{F}_{\text{int}}} \int_F \Re \left( \frac{Z_T}{Z_T + Z_K} (\gamma_t \mathbf{E}^K \cdot \overline{\gamma_{\times}^T \mathbf{H}^T} + \gamma_t \mathbf{E}^K \cdot \overline{\gamma_{\times}^K \mathbf{H}^T}) + \frac{Z_K}{Z_T + Z_K} (\gamma_t \mathbf{E}^T \cdot \overline{\gamma_{\times}^T \mathbf{H}^K} + \gamma_t \mathbf{E}^T \cdot \overline{\gamma_{\times}^K \mathbf{H}^K}) \right) ds_{\mathbf{x}}.$$

Remarking that  $\gamma_{\times}^T + \gamma_{\times}^K = 0$ , we finally have (17).

(b) It remains to evaluate

$$\mathcal{I}^2 = \sum_{T \in \mathcal{T}} \sum_{F \in \mathcal{F}_T} \int_F \Re(a_{T,F}^{EE}) ds_{\mathbf{x}} \quad \text{and} \quad \mathcal{I}^3 = \sum_{T \in \mathcal{T}} \sum_{F \in \mathcal{F}_T} \int_F \Re(a_{T,F}^{HH}) ds_{\mathbf{x}}.$$

We use again (12) and sum over elements, thus differentiating interior faces and exterior faces

$$\mathcal{I}^2 = \sum_{T \in \mathcal{T}} \sum_{F \in \mathcal{F}_T \cap \mathcal{F}_{\text{int}}} \int_F \Re \left( \frac{1}{Z_T + Z_K} (\gamma_t \mathbf{E}^T - \gamma_t \mathbf{E}^K) \cdot \overline{\gamma_t \mathbf{E}^T} \right) ds_{\mathbf{x}} + \sum_{T \in \mathcal{T}} \sum_{F \in \mathcal{F}_T \cap \mathcal{F}_{\text{ext}}} \int_F \Re \left( \frac{1}{Z_T + Z_{\partial\Omega}} \gamma_t \mathbf{E}^T \cdot \overline{\gamma_t \mathbf{E}^T} \right) ds_{\mathbf{x}}.$$

Assembling following a face point of view (see Remark 1.7), we have

$$\begin{aligned} \mathcal{I}^2 &= \sum_{F \in \mathcal{F}_{\text{int}}} \int_F \Re \left( \frac{1}{Z_T + Z_K} (\gamma_t \mathbf{E}^T - \gamma_t \mathbf{E}^K) \cdot \overline{\gamma_t \mathbf{E}^T} + \frac{1}{Z_K + Z_T} (\gamma_t \mathbf{E}^K - \gamma_t \mathbf{E}^T) \cdot \overline{\gamma_t \mathbf{E}^K} \right) ds_{\mathbf{x}} \\ &+ \sum_{F \in \mathcal{F}_{\text{ext}}} \int_F \Re \left( \frac{1}{Z_T + Z_{\partial\Omega}} \gamma_t \mathbf{E}^T \cdot \overline{\gamma_t \mathbf{E}^T} \right) ds_{\mathbf{x}}. \end{aligned}$$

Using the electric field jump definition in (15), we get the following formula for  $\mathcal{I}^2$ , which vanishes for  $F \in \mathcal{F}_{\text{int}}$ ,

$$\mathcal{I}^2 = \sum_{F \in \mathcal{F}_{\text{int}}} \int_F \frac{1}{Z_T + Z_K} \llbracket \gamma_t \mathbf{E} \rrbracket_F \cdot \overline{\llbracket \gamma_t \mathbf{E} \rrbracket_F} ds_{\mathbf{x}} + \sum_{F \in \mathcal{F}_{\text{ext}}} \int_F \frac{1}{Z_T + Z_{\partial\Omega}} \gamma_t \mathbf{E}^T \cdot \overline{\gamma_t \mathbf{E}^T} ds_{\mathbf{x}}.$$

With a similar reasoning and the definition of the magnetic field jump in (15), we also get

$$\mathcal{I}^3 = \sum_{F \in \mathcal{F}_{\text{int}}} \int_F \frac{Z_T Z_K}{Z_T + Z_K} \llbracket \gamma_{\times} \mathbf{H} \rrbracket_F \cdot \overline{\llbracket \gamma_{\times} \mathbf{H} \rrbracket_F} ds_{\mathbf{x}} + \sum_{F \in \mathcal{F}_{\text{ext}}} \int_F \frac{Z_T Z_{\partial\Omega}}{Z_T + Z_{\partial\Omega}} \gamma_{\times}^T \mathbf{H}^T \cdot \overline{\gamma_{\times}^T \mathbf{H}^T} ds_{\mathbf{x}}.$$

We have finally proved the proposition since  $\|\mathbb{E}\|_{\text{DG}}^2 = \Re(a(\mathbb{E}, \mathbb{E})) = \mathcal{I}^2 + \mathcal{I}^3$ .  $\square$

*Proof.* (ii) Let us prove the injectivity. Since  $a(\mathbb{E}, \mathbb{E}) = 0$ , we have  $\|\mathbb{E}\|_{\text{DG}}^2 = 0$  and therefore

$$\llbracket \gamma_t \mathbf{E} \rrbracket_F = \llbracket \gamma_{\times} \mathbf{H} \rrbracket_F = 0, \quad \text{on } F \in \mathcal{F}_{\text{int}} \quad \text{and} \quad \gamma_t \mathbf{E} = \gamma_{\times}^T \mathbf{H} = 0, \quad \text{on } F \in \mathcal{F}_{\text{ext}}. \quad (18)$$

We recall that if  $\mathbb{E}^T$  satisfies (2) and  $\exists F \in \mathcal{F}_T$  such that  $\gamma_t \mathbf{E}^T = 0$  and  $\gamma_{\times}^T \mathbf{H}^T = 0$  on  $F$ , then by unique continuation  $\mathbb{E}^T \equiv 0$ . Let  $\tilde{\mathcal{T}} \subset \mathcal{T}$  be the set of elements satisfying  $\tilde{\mathcal{T}} := \left\{ T \in \mathcal{T} \mid \mathbb{E}^T \equiv 0 \right\}$ .

We will show that  $\tilde{\mathcal{T}} = \mathcal{T}$ .

(a) The set  $\tilde{\mathcal{T}}$  is non-empty since  $\mathbb{E}^T$  is vanishing in any element  $T$  with one face  $F \subset \mathcal{F}_T$  included in the exterior boundary  $\mathcal{F}_{\text{ext}}$ , see (18).

(b) Let  $T \in \mathcal{T}$  and  $K \in \tilde{\mathcal{T}}$  with a common face  $F \in \mathcal{F}_T \cap \mathcal{F}_K$ . Let us show that  $T \in \tilde{\mathcal{T}}$ . Due to (18), we have

$$\gamma_t \mathbf{E}^T = \gamma_t \mathbf{E}^K = 0 \quad \text{and} \quad \gamma_{\times}^T \mathbf{H}^T = -\gamma_{\times}^K \mathbf{H}^K = 0 \quad \text{on } F \in \mathcal{F}_T \cap \mathcal{F}_K \quad \implies \quad \mathbb{E}^T = 0 \quad \text{in } T \quad \implies \quad T \in \tilde{\mathcal{T}}.$$

unique continuation

(c) It follows from the completeness of the neighborhood graph that  $\mathcal{T} = \tilde{\mathcal{T}}$ . Thus, we have  $\mathbb{E}^T = 0$ .  $\square$

The uniqueness of the solution to Problem 1 is proved in (ii) of Proposition 1.3. Its existence is ensured by the fact that the solution to (2), see [42], is also solution to Problem 1. However, it does not mean that the variational formulation (10), with the form  $a$  given by (11a), is well-posed for every  $\ell \in (\mathbb{X}_{\mathcal{T}})^*$ . This question remains an open problem.

**Remark 1.8.** *The point (ii) of Proposition 1.3 implies that any finite dimensional Galerkin approximation of Problem 1 is well-posed.*

## 2. CLASSICAL TREFFTZ SOLVER

The following section aims at devising a generalisation of the UWVF for heterogeneous media, based on Cessenat-Després trace operators (see [9]). Moreover, we show the equivalence between the Trefftz variational formulation based on upwind fluxes, see Problem 1, and the generalised Cessenat-Després UWVF.

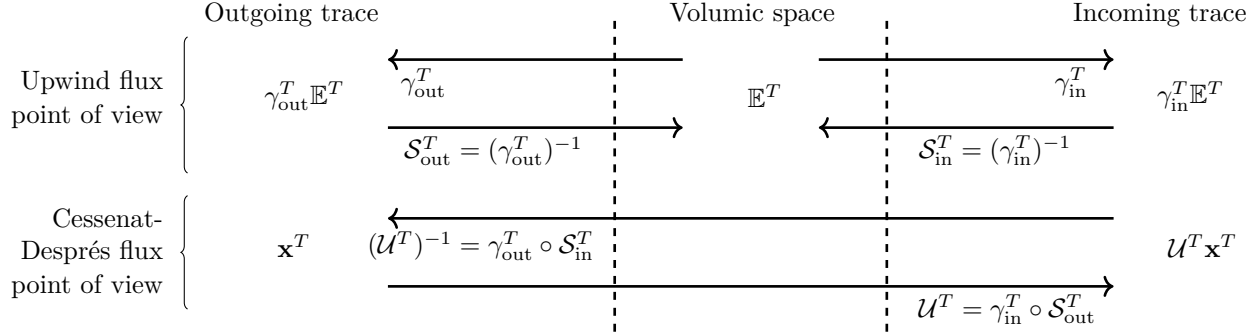


FIGURE 1. Comparison between the upwind flux and the Cessenat-Després flux points of view.

Let us first introduce, for each  $T \in \mathcal{T}$ , the Cessenat-Després operator  $\mathcal{U}^T$  associating to the outgoing trace  $\mathbf{x}^T$  of an electromagnetic field its incoming trace  $\mathcal{U}^T \mathbf{x}^T$  (see Figure 1)

$$\mathbf{x}^T := \gamma_{\text{out}}^T \mathbb{E}^T \stackrel{(9)}{=} \gamma_t \mathbf{E}^T - Z_T \gamma_{\times}^T \mathbf{H}^T \quad \text{and} \quad \mathcal{U}^T \mathbf{x}^T := \gamma_{\text{in}}^T \mathbb{E}^T \stackrel{(9)}{=} \gamma_t \mathbf{E}^T + Z_T \gamma_{\times}^T \mathbf{H}^T. \quad (19)$$

More precisely, the operator  $\mathcal{U}^T : L_t^2(\partial T) \rightarrow L_t^2(\partial T)$  is defined by  $\mathcal{U}^T := \gamma_{\text{in}}^T \circ \mathcal{S}_{\text{out}}^T$ , with the solution operator  $\mathcal{S}_{\text{out}}^T : L_t^2(\partial T) \rightarrow \mathbb{X}_T$ ,  $\mathbf{x}^T \mapsto \mathbb{E}^T$  satisfying (2) and  $\gamma_{\text{out}}^T \mathbb{E}^T = \mathbf{x}^T$ . Using above definitions, the electric tangential component and the magnetic tangential trace can then be deduced

$$\gamma_t \mathbf{E}^T = \frac{1}{2} (\mathcal{U}^T \mathbf{x}^T + \mathbf{x}^T) \quad \text{and} \quad \gamma_{\times}^T \mathbf{H}^T = \frac{1}{2Z_T} (\mathcal{U}^T \mathbf{x}^T - \mathbf{x}^T). \quad (20)$$

Problem 1 can be rephrased with Cessenat-Després notations. Indeed, by artificially adding and subtracting  $\mathbf{x}^T$  to upwind numerical fluxes (7) and (8), we have for  $F \in \mathcal{F}_T$

$$(\widehat{\gamma_t \mathbf{E}})_{|F} = \frac{1}{2} \left( (\widehat{\mathcal{U}^T \mathbf{x}})_{|F} + \mathbf{x}^T \right) \quad \text{and} \quad (\widehat{\gamma_{\times}^T \mathbf{H}})_{|F} = \frac{1}{2Z_T} \left( (\widehat{\mathcal{U}^T \mathbf{x}})_{|F} - \mathbf{x}^T \right), \quad (21)$$

where  $\widehat{\mathcal{U}^T \mathbf{x}}$  is the Cessenat-Després numerical flux associated to the element  $T$  defined by

$$(\widehat{\mathcal{U}^T \mathbf{x}})_{|F} := \frac{Z_K - Z_T}{Z_K + Z_T} \mathbf{x}^T + \frac{2Z_T}{Z_K + Z_T} \mathbf{x}^K \quad \text{for } F \in \mathcal{F}_{\text{int}}, \quad (22a)$$

$$(\widehat{\mathcal{U}^T \mathbf{x}})_{|F} := \frac{Z_{\partial\Omega} - Z_T}{Z_{\partial\Omega} + Z_T} \mathbf{x}^T + \frac{2Z_T}{Z_{\partial\Omega} + Z_T} \mathbf{g}^T \quad \text{for } F \in \mathcal{F}_{\text{ext}}. \quad (22b)$$

**Remark 2.1.** As for Remark 1.6, the formula (22b) can be seen as the interaction with a virtual element  $K$ , when setting  $Z_{\partial\Omega} = Z_K$  and  $\mathbf{g}^T = \mathbf{x}^K$ .

When replacing numerical fluxes (21) in Problem 1, we get for all  $\mathbf{x}' \in L_t^2(\partial\mathcal{T}) := \prod_{T \in \mathcal{T}} L_t^2(\partial T)$

$$\sum_{T \in \mathcal{T}} \int_{\partial T} Z_T^{-1} \left( \mathbf{x}^T \cdot \overline{\mathbf{x}'^T} - \widehat{\mathcal{U}^T \mathbf{x}} \cdot \overline{\widehat{\mathcal{U}^T \mathbf{x}'^T}} \right) ds_{\mathbf{x}} = 0.$$

This leads to an alternative formulation called the Cessenat-Després heterogeneous UWVF.

**Problem 2.** Find  $\mathbf{x} \in L_t^2(\partial\mathcal{T})$  such that for all  $\mathbf{x}' \in L_t^2(\partial\mathcal{T})$  we have

$$(\mathbf{x}, \mathbf{x}')_{L_t^2(\partial\mathcal{T})} - \left( \widehat{\mathcal{U}\mathbf{x}}, \mathcal{U}\mathbf{x}' \right)_{L_t^2(\partial\mathcal{T})} = 0, \quad (23)$$

where  $(\widehat{\mathcal{U}\mathbf{x}})^T := \widehat{\mathcal{U}^T \mathbf{x}}$ ,  $(\mathcal{U}\mathbf{x}')^T := \mathcal{U}^T \mathbf{x}'^T$  and the weighted scalar product on  $L_t^2(\partial\mathcal{T})$  is defined as

$$(\mathbf{x}, \mathbf{x}')_{L_t^2(\partial\mathcal{T})} := \sum_{T \in \mathcal{T}} (\mathbf{x}^T, \mathbf{x}'^T)_{L_t^2(\partial T)} = \sum_{T \in \mathcal{T}} \int_{\partial T} Z_T^{-1} \mathbf{x}^T \cdot \overline{\mathbf{x}'^T} \, d\mathbf{s}_{\mathbf{x}}. \quad (24)$$

As illustrated in Figure 1, a link is established between Problem 2 based on heterogeneous Cessenat-Després fluxes, see (19) and (21), and Problem 1 based on upwind fluxes, see (7) and (8).

**Theorem 2.1.** *Problem 1 is equivalent to Problem 2 in the following sense*

- If  $\mathbf{x}$  is solution to Problem 2 then  $\mathbb{E}$  defined by  $\mathbb{E}^T = \mathcal{S}_{\text{out}}^T \mathbf{x}^T$  is solution to Problem 1.
- If  $\mathbb{E}$  is solution to Problem 1 then  $\mathbf{x}$  defined by  $\mathbf{x}^T = \gamma_{\text{out}}^T \mathbb{E}^T$  is solution to Problem 2.

Moreover if  $\mathbf{x}$  is solution to Problem 2 then it is also solution to the variational formulation (10).

Taking into account (22), we obtain the variational formulation of Problem 2.

**Proposition 2.1.** Find  $\mathbf{x} \in L_t^2(\partial\mathcal{T})$  such that for all  $\mathbf{x}' \in L_t^2(\partial\mathcal{T})$

$$\mathbf{a}(\mathbf{x}, \mathbf{x}') = \mathbf{l}(\mathbf{x}'), \quad \text{with} \quad \mathbf{a}(\mathbf{x}, \mathbf{x}') = (\mathbf{x}, \mathbf{x}')_{L_t^2(\partial\mathcal{T})} - \mathbf{k}(\mathbf{x}, \mathbf{x}'), \quad (25a)$$

where the sesquilinear form  $\mathbf{k} : L_t^2(\partial\mathcal{T}) \times L_t^2(\partial\mathcal{T}) \rightarrow \mathbb{C}$  and the antilinear form  $\mathbf{l} : L_t^2(\partial\mathcal{T}) \rightarrow \mathbb{C}$  are given by

$$\mathbf{k}(\mathbf{x}, \mathbf{x}') := \left( \Pi_{\mathcal{U}\mathbf{x}}, \mathcal{U}\mathbf{x}' \right)_{L_t^2(\partial\mathcal{T})} \quad \text{and} \quad \mathbf{l}(\mathbf{x}') := \left( \mathbf{g}_{\mathcal{U}}, \mathcal{U}\mathbf{x}' \right)_{L_t^2(\partial\mathcal{T})}, \quad (25b)$$

with  $\Pi_{\mathcal{U}} : L_t^2(\partial\mathcal{T}) \rightarrow L_t^2(\partial\mathcal{T})$  the flux operator and  $\mathbf{g}_{\mathcal{U}} \in L_t^2(\partial\mathcal{T})$  the second-member, both defined on each  $T \in \mathcal{T}$  and for all  $F \in \mathcal{F}_T$ , by

$$(\Pi_{\mathcal{U}}^T \mathbf{x})|_F := (\Pi_{\mathcal{U}\mathbf{x}})^T|_F = \begin{cases} \left( \widehat{\mathcal{U}^T \mathbf{x}} \right)|_F & \text{if } F \in \mathcal{F}_{\text{int}}, \\ \frac{Z_{\partial\Omega} - Z_T}{Z_{\partial\Omega} + Z_T} \mathbf{x}^T & \text{if } F \in \mathcal{F}_{\text{ext}}, \end{cases} \quad (\mathbf{g}_{\mathcal{U}}^T)|_F := \begin{cases} 0 & \text{if } F \in \mathcal{F}_{\text{int}}, \\ \frac{2Z_T}{Z_{\partial\Omega} + Z_T} \mathbf{g}^T & \text{if } F \in \mathcal{F}_{\text{ext}}. \end{cases} \quad (26)$$

The next proposition will assert that  $\mathbf{k}$  is contractant. Thus, (25) leads to a fixed point problem.

**Proposition 2.2.** *We have*

- The operator  $\Pi_{\mathcal{U}} : L_t^2(\partial\mathcal{T}) \rightarrow L_t^2(\partial\mathcal{T})$  satisfies  $\|\Pi_{\mathcal{U}}\|_{L_t^2(\partial\mathcal{T})} \leq 1$ .
- The operator  $\mathcal{U} : L_t^2(\partial\mathcal{T}) \rightarrow L_t^2(\partial\mathcal{T})$  satisfies  $\|\mathcal{U}\|_{L_t^2(\partial\mathcal{T})} = 1$ .
- The operator  $\mathbf{k} : L_t^2(\partial\mathcal{T}) \times L_t^2(\partial\mathcal{T}) \rightarrow \mathbb{C}$  is contractant, ie  $\|\mathbf{k}\|_{L_t^2(\partial\mathcal{T})} \leq 1$ ,

where, if  $V$  is an Hilbert space, the norm of the linear operator  $\mathcal{A} : V \rightarrow V$  and the sesquilinear form  $\mathbf{a} : V \times V \rightarrow \mathbb{C}$  are defined by

$$\|\mathcal{A}\|_V := \sup_{\mathbf{x} \in V \setminus \{0\}} \sup_{\mathbf{x}' \in V \setminus \{0\}} \frac{|(\mathcal{A}\mathbf{x}, \mathbf{x}')_V|}{\|\mathbf{x}\|_V \|\mathbf{x}'\|_V} \quad \text{and} \quad \|\mathbf{a}\|_V := \sup_{\mathbf{x} \in V \setminus \{0\}} \sup_{\mathbf{x}' \in V \setminus \{0\}} \frac{|\mathbf{a}(\mathbf{x}, \mathbf{x}')_V|}{\|\mathbf{x}\|_V \|\mathbf{x}'\|_V}. \quad (27)$$

*Proof.* (i) Taking the face point of view from Remark 1.7, we have for all  $\mathbf{x} \in L_t^2(\partial\mathcal{T})$

$$\sum_{T \in \mathcal{T}} \int_{\partial T} Z_T^{-1} |\Pi_{\mathcal{U}}^T \mathbf{x}|^2 \, ds_{\mathbf{x}} = \sum_{F \in \mathcal{F}} \mathcal{K}_F \text{ with } \mathcal{K}_F = \begin{cases} \int_F Z_T^{-1} |\Pi_{\mathcal{U}}^T \mathbf{x}|^2 + Z_K^{-1} |\Pi_{\mathcal{U}}^K \mathbf{x}|^2 \, ds_{\mathbf{x}} & \text{for } F \in \mathcal{F}_{\text{int}}, \\ \int_F Z_T^{-1} |\Pi_{\mathcal{U}}^T \mathbf{x}|^2 \, ds_{\mathbf{x}} & \text{for } F \in \mathcal{F}_{\text{ext}}. \end{cases}$$

For  $F \in \mathcal{F}_{\text{int}}$  interfacing two elements  $T$  and  $K$ , we have by (26)

$$\begin{cases} \mathcal{K}_F &= \int_F Z_T^{-1} \left( \frac{Z_K - Z_T}{Z_K + Z_T} \right)^2 |\mathbf{x}^T|^2 + \frac{2(Z_T - Z_K)}{(Z_K + Z_T)^2} (\mathbf{x}^T \cdot \overline{\mathbf{x}^K} + \mathbf{x}^K \cdot \overline{\mathbf{x}^T}) + \frac{4Z_T}{(Z_K + Z_T)^2} |\mathbf{x}^K|^2 \, ds_{\mathbf{x}} \\ &+ \int_F Z_K^{-1} \left( \frac{Z_T - Z_K}{Z_T + Z_K} \right)^2 |\mathbf{x}^K|^2 + \frac{2(Z_K - Z_T)}{(Z_T + Z_K)^2} (\mathbf{x}^K \cdot \overline{\mathbf{x}^T} + \mathbf{x}^T \cdot \overline{\mathbf{x}^K}) + \frac{4Z_K}{(Z_T + Z_K)^2} |\mathbf{x}^T|^2 \, ds_{\mathbf{x}}, \\ &= \int_F Z_T^{-1} |\mathbf{x}^T|^2 + Z_K^{-1} |\mathbf{x}^K|^2 \, ds_{\mathbf{x}}. \end{cases} \quad (28)$$

Remarking that  $\frac{Z_{\partial\Omega} - Z_T}{Z_{\partial\Omega} + Z_T} \leq 1$ , we obtain for  $F \in \mathcal{F}_{\text{ext}}$

$$\int_F Z_T^{-1} |\Pi_{\mathcal{U}}^T \mathbf{x}^T|^2 \, ds_{\mathbf{x}} = \int_F Z_T^{-1} \left( \frac{Z_{\partial\Omega} - Z_T}{Z_{\partial\Omega} + Z_T} \right)^2 |\mathbf{x}^T|^2 \, ds_{\mathbf{x}} = \int_F Z_T^{-1} |\mathbf{x}^T|^2 \, ds_{\mathbf{x}} - \int_F \frac{4Z_{\partial\Omega}}{(Z_{\partial\Omega} + Z_T)^2} |\mathbf{x}^T|^2 \, ds_{\mathbf{x}}. \quad (29)$$

Using the definition of the global scalar product (24), it follows

$$\|\Pi_{\mathcal{U}} \mathbf{x}\|_{L_t^2(\partial\mathcal{T})}^2 = \sum_{F \in \mathcal{F}} \mathcal{K}_F \text{ with } \mathcal{K}_F \leq \begin{cases} \int_F Z_T^{-1} |\mathbf{x}^T|^2 + Z_K^{-1} |\mathbf{x}^K|^2 \, ds_{\mathbf{x}} & \text{for } F \in \mathcal{F}_{\text{int}}, \\ \int_F Z_T^{-1} |\mathbf{x}^T|^2 \, ds_{\mathbf{x}} & \text{for } F \in \mathcal{F}_{\text{ext}}. \end{cases}$$

Applying Remark 1.7, we get an element assembling point of view leading to

$$\|\Pi_{\mathcal{U}} \mathbf{x}\|_{L_t^2(\partial\mathcal{T})}^2 \leq \sum_{T \in \mathcal{T}} \sum_{F \in \mathcal{F}_T} \int_F Z_T^{-1} |\mathbf{x}^T|^2 \, ds_{\mathbf{x}} = \|\mathbf{x}\|_{L_t^2(\partial\mathcal{T})}^2 \implies \|\Pi_{\mathcal{U}}\|_{L_t^2(\partial\mathcal{T})} \leq 1.$$

□

*Proof.* (ii) Recalling that we have the reciprocity formula (4), let us remark that

$$r_T(\mathbb{E}, \mathbb{E}') = \int_{\partial T} \left( \gamma_{\times}^T \mathbf{H}^T \cdot \overline{\gamma_t \mathbb{E}'^T} + \gamma_t \mathbf{E}^T \cdot \overline{\gamma_{\times}^T \mathbf{H}'^T} \right) \, ds_{\mathbf{x}} = \frac{\mathcal{I}_1 - \mathcal{I}_2}{2} = 0,$$

with  $\mathcal{I}_1$  and  $\mathcal{I}_2$  defined here as

$$\mathcal{I}_1 := \int_{\partial T} Z_T^{-1} \gamma_{\text{in}}^T \mathbb{E}^T \cdot \overline{\gamma_{\text{in}}^T \mathbb{E}'^T} \, ds_{\mathbf{x}} = \int_{\partial T} Z_T^{-1} (\gamma_t \mathbf{E}^T + Z_T \gamma_{\times}^T \mathbf{H}^T) \cdot \left( \overline{\gamma_t \mathbb{E}'^T} + Z_T \overline{\gamma_{\times}^T \mathbf{H}'^T} \right) \, ds_{\mathbf{x}},$$

$$\mathcal{I}_2 := \int_{\partial T} Z_T^{-1} \gamma_{\text{out}}^T \mathbb{E}^T \cdot \overline{\gamma_{\text{out}}^T \mathbb{E}'^T} \, ds_{\mathbf{x}} = \int_{\partial T} Z_T^{-1} (\gamma_t \mathbf{E}^T - Z_T \gamma_{\times}^T \mathbf{H}^T) \cdot \left( \overline{\gamma_t \mathbb{E}'^T} - Z_T \overline{\gamma_{\times}^T \mathbf{H}'^T} \right) \, ds_{\mathbf{x}}.$$

Consequently, for all  $\mathbb{E}^T \in \mathbb{X}_T$ ,  $\mathbb{E}'^T \in \mathbb{X}_T$  and for all  $\mathbf{x}^T \in L_t^2(\partial T)$ ,  $\mathbf{x}'^T \in L_t^2(\partial T)$ , we have

$$\int_{\partial T} Z_T^{-1} \gamma_{\text{out}}^T \mathbb{E}^T \cdot \overline{\gamma_{\text{out}}^T \mathbb{E}'^T} \, ds_{\mathbf{x}} = \int_{\partial T} Z_T^{-1} \gamma_{\text{in}}^T \mathbb{E}^T \cdot \overline{\gamma_{\text{in}}^T \mathbb{E}'^T} \, ds_{\mathbf{x}} = (\mathbb{E}^T, \mathbb{E}'^T)_{\mathbb{X}_T} = (\mathcal{S}_{\text{out}}^T \mathbf{x}^T, \mathcal{S}_{\text{out}}^T \mathbf{x}'^T)_{\mathbb{X}_T}.$$

Since  $\mathcal{S}_{\text{out}}^T : L_t^2(\partial T) \rightarrow \mathbb{X}_T$  is bijective, we notice that the space  $\mathbb{X}_T$  can be parametrised by  $L_t^2(\partial T)$ . Summing over elements  $T \in \mathcal{T}$ , it leads to  $\|\mathbf{x}\|_{L_t^2(\partial\mathcal{T})}^2 = \|\mathcal{U}\mathbf{x}\|_{L_t^2(\partial\mathcal{T})}^2$ . Therefore, we have  $\|\mathcal{U}\|_{L_t^2(\partial\mathcal{T})} = 1$ .  $\square$

*Proof.* (iii) Using the Cessenat-Després point of view, Problem 2 and the Cauchy-Schwarz inequality, we have

$$\mathbf{k}(\mathbf{x}, \mathbf{x}') = \left( \Pi_{\mathcal{U}}\mathbf{x}, \mathcal{U}\mathbf{x}' \right)_{L_t^2(\partial\mathcal{T})} \leq \|\Pi_{\mathcal{U}}\mathbf{x}\|_{L_t^2(\partial\mathcal{T})} \|\mathcal{U}\mathbf{x}'\|_{L_t^2(\partial\mathcal{T})} \stackrel{(i),(ii)}{\leq} \|\mathbf{x}\|_{L_t^2(\partial\mathcal{T})} \|\mathbf{x}'\|_{L_t^2(\partial\mathcal{T})}.$$

This leads to  $\|\mathbf{k}\|_{L_t^2(\partial\mathcal{T})} \leq 1$  and ends the proof.  $\square$

**Remark 2.2.** Property (ii) in Proposition 2.2 can be applied to the Cessenat-Després Problem 2. More precisely, we can replace the scalar product in the left-hand member of equation (23). We then we get another formulation for the heterogeneous UWVF

$$\text{Find } \mathbf{x} \in L_t^2(\partial\mathcal{T}) \text{ such that for all } \mathbf{x}' \in L_t^2(\partial\mathcal{T}), \quad \text{we have } \left( \mathcal{U}\mathbf{x} - \widehat{\mathcal{U}}\mathbf{x}, \mathcal{U}\mathbf{x}' \right)_{L_t^2(\partial\mathcal{T})} = 0.$$

This variational formulation encodes the continuity of the solution, ie  $\mathcal{U}^T \mathbf{x}^T = \widehat{(\mathcal{U}^T \mathbf{x})}$ .

Proposition 2.2 does not ensure that the operator  $\mathbf{k}$  is strictly contractant. Therefore, we have chosen to introduce an iterative Trefftz problem only at a discrete level.

### 3. ITERATIVE TREFFTZ SOLVER

An iterative Trefftz method based on Cessenat-Després numerical fluxes is introduced. The convergence is theoretically established in Proposition 3.1.

We define the global discrete Trefftz space  $\mathbb{X}_{\mathcal{T}}^h$ , that is a finite dimensional linear subspace of  $\mathbb{X}_{\mathcal{T}}$ , as

$$\mathbb{X}_{\mathcal{T}}^h := \left( \prod_{T \in \mathcal{T}} \mathbb{X}_T^h \right) \subset \mathbb{X}_{\mathcal{T}}.$$

The local Trefftz space  $\mathbb{X}_T^h := \text{span} \left\{ \mathbf{v}_T^\ell \in \mathbb{X}_T, \ell = 1, N_T \right\}$ , where the functions  $\mathbf{v}_T^\ell$  are electromagnetic plane waves, with direction  $\mathbf{d}_T^\ell$  and polarisation  $\mathbf{p}_T^\ell$ , which can be chosen as

$$\mathbf{v}_T^\ell := (\mathbf{E}^T, \mathbf{H}^T) \in \mathbb{X}_T \quad \text{with} \quad \mathbf{E}^T := \mathbf{p}_T^\ell e^{ik_0 \sqrt{\varepsilon_r \mu_r} \mathbf{d}_T^\ell \cdot \mathbf{x}} \quad \text{and} \quad \mathbf{H}^T := Z_T (\mathbf{p}_T^\ell \times \mathbf{d}_T^\ell) e^{ik_0 \sqrt{\varepsilon_r \mu_r} \mathbf{d}_T^\ell \cdot \mathbf{x}}, \quad (30)$$

with  $Z_T = \sqrt{\mu_r^T / \varepsilon_r^T}$  a normalised impedance defined on  $T$ . More precisely,  $\mathbf{d}_T^\ell \in \mathcal{D}$ , a finite discrete subspace of the unit sphere, and  $\mathbf{p}_T^\ell \in S_{\mathbf{d}}$ , an orthonormal basis of the two-dimensional linear subspace  $(\mathbf{d}_T^\ell)^\perp \subset \mathbb{R}^3$ . Many different choices for the definition of  $\mathcal{D}$  exist. Directions are given by the vertices of the triangular mesh surface of the unit sphere, or the surface mesh of a unit cube. For a complete plane waves approximation theory, one can refer to [32, 41].

Problem 2 can be discretised on the finite dimensional space  $\mathbb{Y}_{\mathcal{T}}^h \subset L_t^2(\partial\mathcal{T})$

$$\mathbb{Y}_{\mathcal{T}}^h := \prod_{T \in \mathcal{T}} \mathbb{Y}_T^h, \quad \text{with } \mathbb{Y}_T^h := \gamma_{\text{out}}^T \mathbb{X}_T^h.$$

More precisely, for each  $T \in \mathcal{T}$ , the local finite discrete space  $\mathbb{Y}_T^h$  is spanned by the incoming traces of plane waves in  $\mathbb{X}_T^h$ , which are defined by (30). Since  $\gamma_{\text{out}}^T$  is a bijective operator, the space  $\mathbb{Y}_T^h$  is of dimension  $N := \dim(\mathbb{Y}_T^h) = \dim(\mathbb{X}_T^h) = N_T$ . Consequently, the basis  $\mathbb{Y}_T^h$  is given by

$$\mathbb{Y}_T^h := \left\{ \mathbf{w}_T = \gamma_{\text{out}}^T \mathbf{v}_T \text{ such that } \mathbf{v}_T \in \mathbb{X}_T^h \right\} = \text{span}_{\ell=1, N} \left\{ \mathbf{w}_T^\ell \text{ such that } \mathbf{w}_T^\ell = \gamma_{\text{out}}^T \mathbf{v}_T^\ell \right\}.$$

The space  $\mathbb{Y}_{\mathcal{T}}^h$  has the dimension  $\#\text{dof} := N \#\text{elem}$ , with  $\#\text{elem} := \text{card}(\mathcal{T})$  the total number of elements in the mesh. Any element  $\mathbf{x} \in \mathbb{Y}_{\mathcal{T}}^h$  will be represented by a complex column vector  $[\mathbf{x}]$  of dimension  $\#\text{dof}$ . Their components  $[\mathbf{x}]_{\text{iglob}}$  are the amplitudes of the plane wave traces  $\mathbf{w}_i^\ell \in \mathbb{Y}_T^h$  in the  $i^{\text{th}}$  element  $T$ . It leads to

$$\mathbf{x} = \sum_{\text{iglob}=1}^{\#\text{dof}} [\mathbf{x}]_{\text{iglob}} \mathbf{w}^{\text{iglob}} = \sum_{i=1}^{\#\text{elem}} \sum_{\ell=1}^N [\mathbf{x}]_i^\ell \mathbf{w}_i^\ell, \quad \text{with } \text{iglob} := (i-1)N + \ell, \quad \text{and} \quad \begin{cases} [\mathbf{x}]_{\text{iglob}} := [\mathbf{x}]_i^\ell, \\ \mathbf{w}^{\text{iglob}} := \mathbf{w}_i^\ell. \end{cases} \quad (31)$$

The discrete direct Cessenat-Després heterogeneous UWVF is formulated as follows.

**Problem 3.** Find  $\mathbf{x}^h \in \mathbb{Y}_{\mathcal{T}}^h$  such that for all  $\mathbf{x}' \in \mathbb{Y}_{\mathcal{T}}^h$ , we have

$$\mathbf{a}(\mathbf{x}^h, \mathbf{x}') = \mathbf{l}(\mathbf{x}') \iff \mathbf{A}[\mathbf{x}] = \mathbf{F}, \quad (32)$$

with  $\mathbf{A} \in \mathbb{C}^{\#\text{dof} \times \#\text{dof}}$  and  $\mathbf{F} \in \mathbb{C}^{\#\text{dof}}$  defined component by component by

$$\mathbf{A}_{\text{iglob}, \text{jglob}} := \mathbf{a}(\mathbf{w}^{\text{iglob}}, \mathbf{w}^{\text{jglob}}) \quad \text{and} \quad \mathbf{F}_{\text{iglob}} := \mathbf{l}(\mathbf{w}^{\text{iglob}}) \quad \text{for } \text{iglob}, \text{jglob} = 1, \#\text{dof}. \quad (33)$$

**Remark 3.1.** The global matrix  $\mathbf{A}$  is composed of block matrices  $\mathbf{A}_{i,j} := \left( \mathbf{A}_{i,j}^{\ell,k} \right)_{\ell,k=1}^N \in \mathbb{C}^{N \times N}$  for  $i, j = 1, \#\text{elem}$  defined as

$$\mathbf{A}_{i,j}^{\ell,k} = \mathbf{a}(\mathbf{w}_j^k, \mathbf{w}_i^\ell) \quad \text{with } \ell, k = 1, N.$$

Proposition 2.2 leads to a singular regular decomposition of matrix  $\mathbf{A} = \mathbf{M} - \mathbf{N}$ . In the homogeneous case, it corresponds to a Jacobi-block decomposition, see [11]. We then obtain the discrete iterative UWVF problem.

**Problem 4.** Compute the sequence  $\mathbf{x}_n^{\text{jac}} \in \mathbb{Y}_{\mathcal{T}}^h$ ,  $n \in \mathbb{N}^*$ , from the recurrence

$$(\mathbf{x}_{n+1}^{\text{jac}}, \mathbf{x}')_{L_t^2(\partial\mathcal{T})} - \mathbf{k}(\mathbf{x}_n^{\text{jac}}, \mathbf{x}') = \mathbf{l}(\mathbf{x}'), \quad \forall \mathbf{x}' \in \mathbb{Y}_{\mathcal{T}}^h \iff \mathbf{M}[\mathbf{x}_{n+1}^{\text{jac}}] = \mathbf{N}[\mathbf{x}_n^{\text{jac}}] + \mathbf{F}, \quad \text{with } [\mathbf{x}_0^{\text{jac}}] = 0, \quad (34)$$

with  $\mathbf{M}/\mathbf{N} \in \mathbb{C}^{\#\text{dof} \times \#\text{dof}}$  defined by

$$\mathbf{M}_{\text{iglob}, \text{jglob}} := (\mathbf{w}^{\text{iglob}}, \mathbf{w}^{\text{jglob}})_{L_t^2(\partial\mathcal{T})} \quad \text{and} \quad \mathbf{N}_{\text{iglob}, \text{jglob}} := \mathbf{k}(\mathbf{w}^{\text{iglob}}, \mathbf{w}^{\text{jglob}}), \quad \text{for } \text{iglob}, \text{jglob} = 1, \#\text{dof}. \quad (35)$$

**Remark 3.2.** Support properties of DG basis functions lead to a hermitian positive block-diagonal matrix  $\mathbf{M}$

$$\mathbf{M}_{i,j}^{\ell,k} = (\mathbf{w}_i^k, \mathbf{w}_i^\ell)_{L_t^2(\partial\mathcal{T})} \delta_{i,j} \quad \text{for } i, j = 1, \#\text{elem} \quad \text{and} \quad \ell, k = 1, N,$$

which allows a fast direct inversion of the linear system (34).

To ensure the convergence of the discrete iterative Problem 4, the matrix  $\mathbf{M}^{-1}\mathbf{N}$  associated to the form  $\mathbf{k}$  needs to be strictly contractant. In other terms, its spectral radius, denoted by  $\rho(\mathbf{M}^{-1}\mathbf{N})$ , has to verify  $\rho(\mathbf{M}^{-1}\mathbf{N}) < 1$ . Due to the conformal nature of the discretisation, ie  $\mathbb{Y}_{\mathcal{T}}^h \subset L_t^2(\partial\mathcal{T})$ , we have

$$\rho(\mathbf{M}^{-1}\mathbf{N}) \leq \|\mathbf{k}\|_{\mathbb{Y}_{\mathcal{T}}^h} = \sup_{\mathbf{x} \in \mathbb{Y}_{\mathcal{T}}^h \setminus \{0\}} \sup_{\mathbf{x}' \in \mathbb{Y}_{\mathcal{T}}^h \setminus \{0\}} \frac{|\mathbf{k}(\mathbf{x}, \mathbf{x}')|}{\|\mathbf{x}\|_{\mathbb{Y}_{\mathcal{T}}^h} \|\mathbf{x}'\|_{\mathbb{Y}_{\mathcal{T}}^h}} \leq \|\mathbf{k}\|_{L_t^2(\partial\mathcal{T})} \leq 1, \quad \text{with } \|\mathbf{x}\|_{\mathbb{Y}_{\mathcal{T}}^h} = \|\mathbf{x}\|_{L_t^2(\partial\mathcal{T})}.$$

It then remains to exclude the eigenvalues on the unit circle.

**Proposition 3.1.** The matrix  $\mathbf{M}^{-1}\mathbf{N}$  is strictly contractant, ie  $\rho(\mathbf{M}^{-1}\mathbf{N}) < 1$ .

*Proof.* We suppose that it exists  $\mathbf{x} \in \mathbb{Y}_{\mathcal{T}}^h$  and  $\lambda \in \mathbb{C}$ , such that  $\mathbf{x} \neq 0$  and  $|\lambda| = 1$ . We have

$$\mathbf{k}(\mathbf{x}, \mathbf{x}') = \lambda (\mathbf{x}, \mathbf{x}')_{L_t^2(\partial\mathcal{T})}, \quad \forall \mathbf{x}' \in \mathbb{Y}_{\mathcal{T}}^h \iff \mathbf{N}[\mathbf{x}] = \lambda \mathbf{M}[\mathbf{x}], \quad \text{with } \mathbf{x} \text{ represented by } [\mathbf{x}] \text{ through (31).}$$

It leads to

$$\|\Pi_{\mathcal{U}\mathbf{x}} - \lambda\mathcal{U}\mathbf{x}\|_{L_i^2(\partial\mathcal{T})}^2 = \|\Pi_{\mathcal{U}\mathbf{x}}\|_{L_i^2(\partial\mathcal{T})}^2 - \lambda(\mathcal{U}\mathbf{x}, \Pi_{\mathcal{U}\mathbf{x}})_{L_i^2(\partial\mathcal{T})} - \bar{\lambda}(\Pi_{\mathcal{U}\mathbf{x}}, \mathcal{U}\mathbf{x})_{L_i^2(\partial\mathcal{T})} + |\lambda|^2 \|\mathcal{U}\mathbf{x}\|_{L_i^2(\partial\mathcal{T})}^2.$$

Thanks to Proposition 2.2 and the definition of  $\mathbf{k}$ , see (25b), we obtain

$$\|\Pi_{\mathcal{U}\mathbf{x}} - \lambda\mathcal{U}\mathbf{x}\|_{L_i^2(\partial\mathcal{T})} \leq \|\mathbf{x}\|_{L_i^2(\partial\mathcal{T})}^2 - \lambda\bar{\lambda} \|\mathbf{x}\|_{L_i^2(\partial\mathcal{T})} - \bar{\lambda}\lambda \|\mathbf{x}\|_{L_i^2(\partial\mathcal{T})} + |\lambda|^2 \|\mathcal{U}\mathbf{x}\|_{L_i^2(\partial\mathcal{T})}^2 \Big|_{|\lambda|=1} \leq 0 \implies \Pi_{\mathcal{U}\mathbf{x}} = \lambda\mathcal{U}\mathbf{x}. \quad (36)$$

Then, we act by contradiction. Let us show that  $\mathbf{x} = 0$ . Let  $\tilde{\mathcal{T}} := \{T \in \mathcal{T} \mid \mathbf{x}^T \equiv 0\} \subset \mathcal{T}$ . It remains to prove that  $\tilde{\mathcal{T}} = \mathcal{T}$  to obtain Proposition 3.1.

(a) Let us show that  $\tilde{\mathcal{T}}$  is non-empty. From (36), we have  $\|\Pi_{\mathcal{U}\mathbf{x}}\|_{L_i^2(\partial\mathcal{T})} = \|\mathcal{U}\mathbf{x}\|_{L_i^2(\partial\mathcal{T})} = \|\mathbf{x}\|_{L_i^2(\partial\mathcal{T})}$  since  $\mathcal{U}$  is unitary from Proposition 2.2. From (28) and (29), it follows that

$$0 = \|\mathbf{x}\|_{L_i^2(\partial\mathcal{T})}^2 - \|\Pi_{\mathcal{U}\mathbf{x}}\|_{L_i^2(\partial\mathcal{T})}^2 = \sum_{T \in \mathcal{T}} \sum_{F \in \mathcal{F}_T \cap \mathcal{F}_{\text{ext}}} \int_F \underbrace{\frac{4Z_{\partial\Omega}}{(Z_{\partial\Omega} + Z_T)^2}}_{>0} |\mathbf{x}^T|^2 ds_{\mathbf{x}} \implies \mathbf{x}^T = 0 \quad \text{on } F \in \mathcal{F}_T \cap \mathcal{F}_{\text{ext}}.$$

Moreover, recalling the definition (26) and due to (36), we also have

$$\mathcal{U}^T \mathbf{x}^T = \frac{1}{\lambda} \frac{Z_{\partial\Omega} - Z_T}{Z_{\partial\Omega} + Z_T} \mathbf{x}^T = 0 \quad \text{on } F \in \mathcal{F}_T \cap \mathcal{F}_{\text{ext}}.$$

Combining the two last results and using (20), we obtain  $\gamma_t \mathbf{E}^T = \gamma_{\times} \mathbf{H}^T = 0$  on  $F$ , leading to  $\mathbb{E}^T \equiv 0$  and  $\mathbf{x}^T \equiv 0$  in  $T$  by the unique continuation.

(b) Let  $T \in \mathcal{T}$  and  $K \in \tilde{\mathcal{T}}$  with a common face  $F \in \mathcal{F}_T \cap \mathcal{F}_K$ . Let us show that  $T \in \tilde{\mathcal{T}}$ . Since  $K \in \tilde{\mathcal{T}}$ , we have  $\mathbf{x}^K = 0$  and  $\mathcal{U}^K \mathbf{x}^K = 0$ . From (36), we have  $\Pi_{\mathcal{U}\mathbf{x}}^T = \lambda \mathcal{U}^T \mathbf{x}^T$  and  $\Pi_{\mathcal{U}\mathbf{x}}^K = \lambda \mathcal{U}^K \mathbf{x}^K$  on  $F$ . From (22a) and (26), we obtain

$$\frac{Z_T - Z_K}{Z_K + Z_T} \mathbf{x}^K + \frac{2Z_K}{Z_K + Z_T} \mathbf{x}^T = \lambda \mathcal{U}^K \mathbf{x}^K \quad \text{and} \quad \frac{Z_K - Z_T}{Z_K + Z_T} \mathbf{x}^T + \frac{2Z_T}{Z_K + Z_T} \mathbf{x}^K = \lambda \mathcal{U}^T \mathbf{x}^T,$$

leading to  $\mathbf{x}^T = 0$  and  $\mathcal{U}^T \mathbf{x}^T = 0$  on  $F$ , since  $Z_{T/K} > 0$  and  $\lambda \neq 0$ . Thus,  $T \in \tilde{\mathcal{T}}$ . Due to (20), we also obtain  $\gamma_t \mathbf{E}^T = \gamma_{\times} \mathbf{H}^T = 0$  on  $F$ , leading to  $\mathbb{E}^T \equiv 0$  and  $\mathbf{x}^T \equiv 0$  in  $T$  by the unique continuation theorem.

(c) From the completeness and the connexity of the neighborhood graph, we then have  $\mathcal{T} = \tilde{\mathcal{T}}$ .  $\square$

**Remark 3.3.** A weaker contraction result has been proved in the context of an under-relaxation iterative method by Cessenat-Després (see [10]). This result takes the form of  $\rho((1 - \beta)\mathbf{I} + \beta \mathbf{M}^{-1}\mathbf{N}) < 1$  for all  $\beta \in ]0, 1[$ . In particular, the Proposition 3.1 refines this result by including the limit case  $\beta = 1$ .

**Remark 3.4.** This proof ensures the uniqueness of the solution to the variational formulation (32). Combined with the rank-nullity theorem, it leads to the well-posedness of Problem 4 which is finite dimensional.

This result makes it possible to use a Jacobi like iterative method to solve the linear system. However, the spectral radius of the iterative matrix  $\mathbf{M}^{-1}\mathbf{N}$  can be very close of 1 and can significantly slow down the iterative algorithm as illustrated in Table 1 and Fig. 2.

Improvements proposed in Section 4 will show that it is related to the Cessenat-Després method.

$\mathcal{D}_{\Omega}$	8	20	40	80
$\rho(\mathbf{M}^{-1}\mathbf{N})$	1 - 0.0064	1 - 0.0022	1 - 0.0014	1 - 0.00054

TABLE 1. Comparison of  $\rho(\mathbf{M}^{-1}\mathbf{N})$  thanks to the power method, on different  $\mathcal{D}_{\Omega}$  in wavelength.



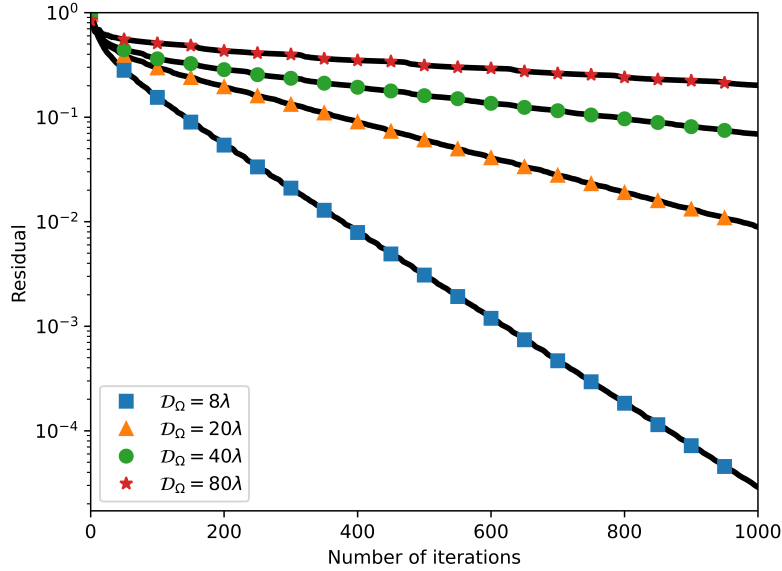


FIGURE 2. Comparison of the convergence rate of the iterative Cessenat-Despres method in function of the size of the domain.

**Remark 3.5.** Unless otherwise specified, the following configuration is considered for numerical experiments

- All distances or lengths are given in wavelengths, ie  $k = 2\pi$ .
- The domain is of size  $\mathcal{D}_\Omega$  and is decomposed into elements  $T$  which are cubes of size  $\mathcal{D}_T = 1$ .
- The reflexion coefficient on  $\partial\Omega$  is denoted by  $R_{\partial\Omega}$ . In particular, we set  $Z_{\partial\Omega} = (1 - R_{\partial\Omega})/(1 + R_{\partial\Omega})$  with  $R_{\partial\Omega} = 0.9$ .
- The number of plane wave basis functions is  $N = 52$ , that are homogeneously distributed in space.

#### 4. NEW ITERATIVE TREFFTZ SOLVER

In the present section, we propose alternatives to ensure the development of a robust Trefftz iterative method, even with the presence of rounding errors. First, we will guarantee its convergence thanks to a preconditioned Krylov method. Then, we resort to a reduction strategy to consider less basis functions providing a lower numerical cost. Finally, we set up a global preconditioner to improve the convergence rate of the method.

##### 4.1. Preconditioned iterative Trefftz solver

A large amount of Krylov type methods exist, see [50]. In this paper, we focus on one based on the positivity of the matrix  $\mathbf{A}$ , ie  $[\mathbf{x}]^* \mathbf{A} [\mathbf{x}] > 0$ . It is a variational method that resorts to Krylov spaces to reduce the dimension of the Galerkin space, see [51, 52]. It can be interpreted as an iterative numerical method which does not require the computation of the inverse of the matrix  $\mathbf{A}$ . When associated to a restart strategy, this method has a low memory cost and is appropriated to large numerical cases. Moreover, it is perfectly suited to real positive matrices, thus adapted to the inversion of Problem 3 which satisfies  $[\mathbf{x}]^* \mathbf{A} [\mathbf{x}] > 0$ . The Galerkin framework of the UWVF method leads to the introduction of the following iterative Krylov UWVF problem.

**Problem 5.** Setting  $\mathbf{x}_0 = 0$ , compute the sequence  $(\mathbf{x}_n)_{n \in \mathbb{N}^*}$  such that  $\mathbf{x}_n \in \mathbb{K}_n$  is the solution of

$$\text{Find } \mathbf{x}_n \in \mathbb{K}_n, \mathbf{a}(\mathbf{x}_n, \mathbf{x}') = \mathbf{l}(\mathbf{x}'), \forall \mathbf{x}' \in \mathbb{K}_n \iff \text{Find } [\mathbf{x}_n] \in [\mathbb{K}_n], [\mathbf{x}']^* \mathbf{A} [\mathbf{x}_n] = [\mathbf{x}']^* \mathbf{F}, \forall [\mathbf{x}'] \in [\mathbb{K}_n], \quad (37)$$

with  $\mathbb{K}_n$  a linear subspace of  $\mathbb{Y}_{\mathcal{T}}^h$  defined as:  $\mathbf{x}_n \in \mathbb{K}_n$  is represented by  $[\mathbf{x}_n] \in [\mathbb{K}_n]$  through the bijection defined by (31), where  $[\mathbb{K}_n]$  is the Krylov space associated to  $\mathbf{A}$  defined as

$$[\mathbb{K}_n] := \text{span}_{0 \leq k \leq n-1} (\mathbf{A}^k \mathbf{F}) = \text{span} \{ \mathbf{F}, \mathbf{A}\mathbf{F}, \mathbf{A}^2\mathbf{F}, \dots, \mathbf{A}^{n-1}\mathbf{F} \}. \quad (38)$$

The convergence theory of Problem 5 is established thanks to the Galerkin theory of the UWVF Problem 2.

**Proposition 4.1.** *Let  $\mathbf{x}^h \in \mathbb{Y}_{\mathcal{T}}^h$  and  $\mathbf{x}_n \in \mathbb{K}_n$  be the solutions of Problems 2 and 5, respectively. The convergence of the iterative Krylov UWVF method is ensured by*

$$\|\mathbf{x}^h - \mathbf{x}_n\|_{\text{DG}} \leq \sqrt{2} \min_{\mathbf{y} \in \mathbb{K}_n} \|\mathbf{x}^h - \mathbf{y}\|_{L_t^2(\partial\mathcal{T})}. \quad (39)$$

*Proof.* Let us recall that we have  $\|\mathbf{x}^h - \mathbf{x}_n\|_{\text{DG}}^2 = \Re(\mathbf{a}(\mathbf{x}^h - \mathbf{x}_n, \mathbf{x}^h - \mathbf{x}_n))$ . Since  $\mathbf{a}(\mathbf{x}^h - \mathbf{x}_n, \mathbf{x}_n - \mathbf{y}) = 0$  for all  $\mathbf{y} \in \mathbb{K}_n$ , we have

$$\|\mathbf{x}^h - \mathbf{x}_n\|_{\text{DG}}^2 = \Re(\mathbf{a}(\mathbf{x}^h - \mathbf{x}_n, \mathbf{x}^h - \mathbf{y})) = \Re([\mathbf{x}^h - \mathbf{y}]^* \mathbf{A}[\mathbf{x}^h - \mathbf{x}_n]).$$

Since  $\mathbf{M}$  is symmetric positive definite, we have

$$[\mathbf{x}^h - \mathbf{y}]^* \mathbf{A}[\mathbf{x}^h - \mathbf{x}_n] = [\mathbf{x}^h - \mathbf{y}]^* \mathbf{M}^{\frac{1}{2}} \mathbf{M}^{-\frac{1}{2}} \mathbf{A}[\mathbf{x}^h - \mathbf{x}_n] = (\mathbf{M}^{\frac{1}{2}}[\mathbf{x}^h - \mathbf{y}])^* \mathbf{M}^{-\frac{1}{2}} \mathbf{A}[\mathbf{x}^h - \mathbf{x}_n].$$

Then, we apply the Cauchy-Schwarz inequality and we get

$$\left\{ \begin{array}{l} \|\mathbf{x}^h - \mathbf{x}_n\|_{\text{DG}}^2 \leq \sqrt{(\mathbf{M}^{\frac{1}{2}}[\mathbf{x}^h - \mathbf{y}])^* \mathbf{M}^{\frac{1}{2}}[\mathbf{x}^h - \mathbf{y}]} \sqrt{(\mathbf{M}^{-\frac{1}{2}} \mathbf{A}[\mathbf{x}^h - \mathbf{x}_n])^* \mathbf{M}^{-\frac{1}{2}} \mathbf{A}[\mathbf{x}^h - \mathbf{x}_n]}, \\ \leq \sqrt{[\mathbf{x}^h - \mathbf{y}]^* \mathbf{M}[\mathbf{x}^h - \mathbf{y}]} \sqrt{(\mathbf{A}[\mathbf{x}^h - \mathbf{x}_n])^* \mathbf{M}^{-1} \mathbf{A}[\mathbf{x}^h - \mathbf{x}_n]}. \end{array} \right.$$

Let us now prove that, for all  $[\mathbf{x}] \in \mathbb{C}^{\#\text{dof}}$ ,  $(\mathbf{A}[\mathbf{x}])^* \mathbf{M}^{-1}(\mathbf{A}[\mathbf{x}]) \leq 2 \Re([\mathbf{x}]^* \mathbf{A}[\mathbf{x}])$ . Using the Cessenat-Després decomposition  $\mathbf{A} = \mathbf{M} - \mathbf{N}$ , we get

$$(\mathbf{A}[\mathbf{x}])^* \mathbf{M}^{-1}(\mathbf{A}[\mathbf{x}]) = [\mathbf{x}]^* \mathbf{M}[\mathbf{x}] - [\mathbf{x}]^* \mathbf{N}[\mathbf{x}] - [\mathbf{x}]^* \mathbf{N}^*[\mathbf{x}] + [\mathbf{x}]^* \mathbf{N}^* \mathbf{M}^{-1} \mathbf{N}[\mathbf{x}].$$

Let us take  $[\mathbf{z}] \in \mathbb{C}^{\#\text{dof}}$ , set as  $[\mathbf{z}] = \mathbf{M}^{-1} \mathbf{N}[\mathbf{x}]$ . It leads to  $[\mathbf{x}]^* \mathbf{N}^* \mathbf{M}^{-1} \mathbf{N}[\mathbf{x}] = [\mathbf{z}]^* \mathbf{N}[\mathbf{x}] = \mathbf{k}(\mathbf{x}, \mathbf{z})$ . Due to (iii) of Proposition 2.2, we have

$$[\mathbf{x}]^* \mathbf{N}^* \mathbf{M}^{-1} \mathbf{N}[\mathbf{x}] \leq \|\mathbf{x}\|_{L_t^2(\partial\mathcal{T})} \|\mathbf{z}\|_{L_t^2(\partial\mathcal{T})} = \sqrt{[\mathbf{z}]^* \mathbf{M}[\mathbf{z}]} \sqrt{[\mathbf{x}]^* \mathbf{M}[\mathbf{x}]} \leq \sqrt{[\mathbf{x}]^* \mathbf{N}^* \mathbf{M}^{-1} \mathbf{N}[\mathbf{x}]} \sqrt{[\mathbf{x}]^* \mathbf{M}[\mathbf{x}]}.$$

Therefore, we get  $[\mathbf{x}]^* \mathbf{N}^* \mathbf{M}^{-1} \mathbf{N}[\mathbf{x}] \leq [\mathbf{x}]^* \mathbf{M}[\mathbf{x}]$ , leading to

$$(\mathbf{A}[\mathbf{x}])^* \mathbf{M}^{-1}(\mathbf{A}[\mathbf{x}]) \leq 2 [\mathbf{x}]^* \mathbf{M}[\mathbf{x}] - 2 \Re([\mathbf{x}]^* \mathbf{N}[\mathbf{x}]) = 2 \Re([\mathbf{x}]^* \mathbf{A}[\mathbf{x}]) = 2 \|\mathbf{x}\|_{\text{DG}}^2.$$

The result follows from  $\|\mathbf{x}^h - \mathbf{x}_n\|_{\text{DG}} \leq \sqrt{2} \|\mathbf{x}^h - \mathbf{y}\|_{L_t^2(\partial\mathcal{T})}$ .  $\square$

We have not succeed to obtain an optimal estimate for the right-hand side of (39). Many convergence bounds of the Krylov residual exist in the literature, see [21]. But these bounds are generally pessimistic. Practically, Krylov solver often has a much better behaviour. Although we are aware of this aspect, we use the general convergence theory of Krylov methods, see [52], to estimate the right-hand side of (39).

**Proposition 4.2.** *Suppose that  $\mathbf{A}$  is diagonalisable, ie  $\mathbf{A} = \mathbf{X} \mathbf{D} \mathbf{X}^{-1}$ , we have for  $\mathbf{x}^h \in \mathbb{Y}_{\mathcal{T}}^h$*

$$\min_{\mathbf{y} \in \mathbb{K}_n} \|\mathbf{x}^h - \mathbf{y}\|_{L^2(\partial\mathcal{T})} \leq \min_{p \in \mathcal{P}_n^0} \max_{\lambda \in \sigma^{\mathbf{A}}} \left| \frac{p(\lambda)}{\lambda} \right| \sqrt{\kappa(\mathbf{M})} \kappa(\mathbf{X}) \|\mathbf{1}\|_{(\mathbb{Y}_{\mathcal{T}}^h)^*}, \quad (40)$$

where  $\mathcal{P}_n^0$  is the set of polynomials of degree  $\leq n$  satisfying  $p(0) = 1$ ,  $\sigma^{\mathbf{A}}$  is the spectrum of  $\mathbf{A}$ ,  $\kappa(\cdot)$  is the condition number and  $(\mathbb{Y}_{\mathcal{T}}^h)^*$  is the dual space of  $\mathbb{Y}_{\mathcal{T}}^h$ .

*Proof.* Let  $\mathbf{y} \in \mathbb{K}_n$ . Thus, it exists  $q$  a polynomial of degree  $n - 1$  such that  $\mathbf{y} = q(\mathbf{A})\mathbf{F}$ . Since  $\mathbf{A}$  is assumed to be diagonalisable, ie  $\mathbf{A} = \mathbf{X}\mathbf{D}\mathbf{X}^{-1}$ , we have

$$\mathbf{X}^{-1}([\mathbf{x}^h - \mathbf{y}]) = (\mathbf{D}^{-1} - q(\mathbf{D}))\mathbf{X}^{-1}\mathbf{F}.$$

Noting that the Euclidian norm of a diagonal matrix is the highest absolute value of its coefficients:

$$\|\mathbf{X}^{-1}[\mathbf{x}^h - \mathbf{y}]\| \leq \max_{\lambda \in \sigma^{\mathbf{A}}} \left| \frac{1}{\lambda} - q(\lambda) \right| \|\mathbf{X}^{-1}\mathbf{F}\|,$$

and using the estimate:

$$\|[\mathbf{x}^h - \mathbf{y}]\| \leq \|\mathbf{X}\| \|\mathbf{X}^{-1}[\mathbf{x}^h - \mathbf{y}]\|,$$

we can write

$$\|[\mathbf{x}^h - \mathbf{y}]\| \leq \max_{\lambda \in \sigma^{\mathbf{A}}} \left| \frac{1 - \lambda q(\lambda)}{\lambda} \right| \kappa(\mathbf{X}) \|\mathbf{F}\|,$$

with  $\kappa(\mathbf{X}) := \|\mathbf{X}\| \|\mathbf{X}^{-1}\|$  the condition number of  $\mathbf{X}$ .

We end the proof by using these following statements:

- the bijection  $q \in \mathcal{P}_{n-1} \mapsto p(x) = 1 - xq(x) \in \mathcal{P}_n^0$ ,
- $\forall \mathbf{z} \in \mathbb{Y}_{\mathcal{T}}^h$ ,

$$\|\mathbf{M}^{-1}\|^{-1} \|[\mathbf{z}]\|^2 \leq \|\mathbf{z}\|_{L^2(\partial\mathcal{T})}^2 = [\mathbf{z}]^* \mathbf{M} [\mathbf{z}] \leq \|\mathbf{M}\| \|[\mathbf{z}]\|^2.$$

□

The convergence rate depends on the spectrum of  $\mathbf{A}$  and draws a particular attention to eigenvalues close to 0, see in red in Figure 3. The fact that  $\mathbf{A}$  is diagonalisable is theoretically discutable, even if it appears to be in practice.

In the continuation, we introduce a preconditioned iterative Krylov UWVF method.

**Problem 6.** Setting  $\mathbf{x}_0^{\text{prec}} = 0$ , compute the sequence  $(\mathbf{x}_n^{\text{prec}})_{n \in \mathbb{N}^*}$  solution to: Find  $\mathbf{x}_n^{\text{prec}} \in \mathbb{K}_n^{\text{prec}}$ , or equivalently  $[\mathbf{x}_n^{\text{prec}}] \in [\mathbb{K}_n^{\text{prec}}]$ , such that we have

$$\mathbf{a}(\mathbf{x}_n^{\text{prec}}, \mathbf{x}') = \mathbf{l}(\mathbf{x}'), \quad \forall \mathbf{x}' \in \mathbb{K}_n^{\text{prec}} \iff [\mathbf{x}']^* \mathbf{A} [\mathbf{x}_n^{\text{prec}}] = [\mathbf{x}']^* \mathbf{F}, \quad \forall [\mathbf{x}'] \in [\mathbb{K}_n^{\text{prec}}], \quad (41)$$

where the finite dimensional space  $\mathbb{K}_n^{\text{prec}}$  is the image of  $[\mathbb{K}_n^{\text{prec}}]$  through (31). The space  $[\mathbb{K}_n^{\text{prec}}]$  is defined by

$$[\mathbb{K}_n^{\text{prec}}] := \mathbf{M}^{-\frac{1}{2}} [\widetilde{\mathbb{K}}_n^{\text{prec}}] \quad \text{with} \quad [\widetilde{\mathbb{K}}_n^{\text{prec}}] := \text{span}_{0 \leq k \leq n-1} \left( \widetilde{\mathbf{A}}^k \widetilde{\mathbf{F}} \right) = \text{span} \left\{ \widetilde{\mathbf{F}}, \widetilde{\mathbf{A}}\widetilde{\mathbf{F}}, \dots, \widetilde{\mathbf{A}}^{n-1}\widetilde{\mathbf{F}} \right\}, \quad (42)$$

where  $\widetilde{\mathbf{A}} = \mathbf{M}^{-\frac{1}{2}} \mathbf{A} \mathbf{M}^{-\frac{1}{2}}$  and  $\widetilde{\mathbf{F}} = \mathbf{M}^{-\frac{1}{2}} \mathbf{F}$ , leading to  $[\mathbb{K}_n^{\text{prec}}] = \text{span} \{ \mathbf{M}^{-1}\mathbf{F}, \dots, (\mathbf{M}^{-1}\mathbf{N})^{n-1} \mathbf{M}^{-1}\mathbf{F} \}$ .

We remark that  $[\mathbf{x}_n^{\text{jac}}] \in \mathbb{C}^{\#\text{dof}}$ , satisfying Problem 4, also belongs to  $[\mathbb{K}_n^{\text{prec}}]$ , ie  $\mathbf{x}_n^{\text{jac}} \in \mathbb{K}_n^{\text{prec}}$ . Thus, we obtain a convergence theorem for the preconditioned Krylov UWVF problem.

**Theorem 4.1.** *The preconditioned Krylov UWVF method converges minimally at the same rate than the iterative UWVF method. We have*

$$\|\mathbf{x}^h - \mathbf{x}_n^{\text{prec}}\|_{\text{DG}} \leq \sqrt{2} \|\mathbf{x}^h - \mathbf{x}_n^{\text{jac}}\|_{L_t^2(\partial\mathcal{T})} \xrightarrow{n \rightarrow +\infty} 0. \quad (43)$$

**Remark 4.1.** *From Proposition 3.1, the convergence rate of  $\mathbf{x}_n^{\text{jac}} \in \mathbb{Y}_{\mathcal{T}}^h$  to  $\mathbf{x}^h \in \mathbb{Y}_{\mathcal{T}}^h$  is  $\rho(\mathbf{M}^{-\frac{1}{2}} \mathbf{N} \mathbf{M}^{-\frac{1}{2}}) < 1$ .*

**Remark 4.2.** *The definition of  $\mathbb{K}_n^{\text{prec}}$  in (42) and the Theorem 4.1 put forward the advantage to perform a symmetric preconditioning of  $\mathbf{A}$  with  $\mathbf{M}^{-\frac{1}{2}}$ .*

The convergence rate of the preconditioned Krylov UWVF (*resp.* iterative UWVF) solver depends on the spectrum of  $\tilde{\mathbf{A}}$  (*resp.*  $\mathbf{A}$ ), see Figure 4 (*resp.* Figure 3). The spectrum of  $\tilde{\mathbf{A}}$  is contained into the unity circle and graphically justifies the use of the convergence Proposition 3.1 on Problem 6. We believe that the presence of fewer eigenvalues close to 0 results in a faster convergence of the Krylov solver, see Figure 4. Moreover,  $\tilde{\mathbf{A}}$  has a lower condition number than  $\mathbf{A}$ , see Table 2, which testifies its better convergence properties.

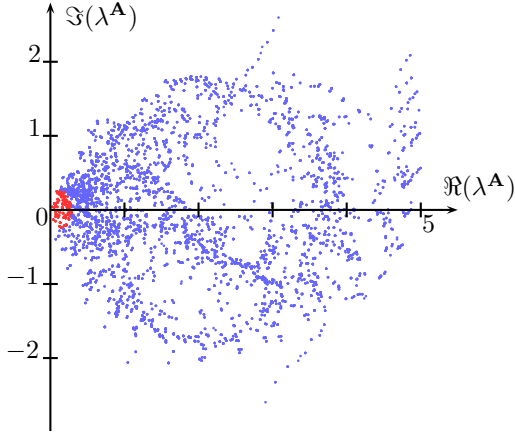


FIGURE 3. Real and complex parts, *resp.*  $\Re(\lambda^{\mathbf{A}})$  and  $\Im(\lambda^{\mathbf{A}})$ , of the spectrum  $\lambda^{\mathbf{A}}$  of matrix  $\mathbf{A}$  for  $\mathcal{D}_\Omega = 6$ .

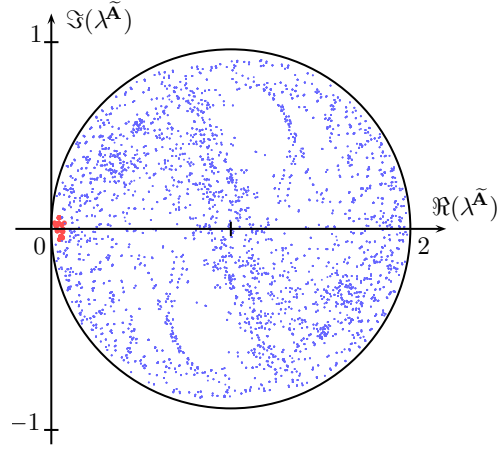


FIGURE 4. Real and complex parts, *resp.*  $\Re(\lambda^{\tilde{\mathbf{A}}})$  and  $\Im(\lambda^{\tilde{\mathbf{A}}})$ , of the spectrum  $\lambda^{\tilde{\mathbf{A}}}$  of matrix  $\tilde{\mathbf{A}}$  for  $\mathcal{D}_\Omega = 6$ .

$\mathcal{D}_\Omega$	1	2	3	4	5	6
$\kappa(\mathbf{A})$	18.5	23.4	45.1	60.7	88.8	113
$\kappa(\tilde{\mathbf{A}})$	4.82	10.9	25.6	29.9	59.8	60.9

TABLE 2. Comparison of the condition number of  $\mathbf{A}$  and  $\tilde{\mathbf{A}}$  according to  $\mathcal{D}_\Omega$  in wavelength.

## 4.2. Basis reduction strategy

The preconditioning of the UWVF matrix  $\mathbf{A}$  ensures the convergence of the preconditioned Krylov solver, see Theorem 4.1. Previous work, see [3, 15, 40, 48], have shown that plane wave basis functions can be numerically linearly-dependant in the sense that  $\mathbf{M}$  has small eigenvalues. A rounding error on smallest eigenvalues of  $\mathbf{M}$  gives rise to a large error on  $\mathbf{M}^{-\frac{1}{2}}$  and impacts the convergence of the Krylov solver. Thus, we reduce the Galerkin space  $\mathbb{Y}_T^h$  by keeping only the eigenvectors associated to the highest eigenvalues. This will ensure a representation of  $\mathbf{x} \in \mathbb{Y}_T^h$  filtering the numerical noise in the plane wave basis, through (31). The diagonalisation of the hermitian matrix  $\mathbf{M}$  representing the  $L^2(\partial\mathcal{T})$  scalar product, see (35), is

$$\mathbf{M} = \mathbf{T}\mathbf{\Lambda}\mathbf{T}^* \quad \text{with } \mathbf{T} \in \mathbb{C}^{\#\text{dof} \times \#\text{dof}} \text{ and } \mathbf{\Lambda} \in \mathbb{C}^{\#\text{dof} \times \#\text{dof}},$$

where  $\mathbf{T}$  is the orthogonal eigenvector matrix and  $\mathbf{\Lambda}$  is the diagonal eigenvalue matrix satisfying

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{\#\text{dof}} \quad \text{with } \lambda_i = \mathbf{\Lambda}_{i,i}.$$

When performing the basis reduction, the orthogonal matrix then becomes rectangular

$$(\mathbf{T}_{\text{red}})_{i,j} = (\mathbf{T})_{i,j} \quad \text{for } i = 1, \#\text{dof} \text{ and } j = 1, \#\text{dof}_{\text{red}}, \quad \text{with } \lambda_{\#\text{dof}_{\text{red}}} \geq \lambda_1 \varepsilon > \lambda_{\#\text{dof}_{\text{red}}+1},$$

where the small threshold coefficient  $\varepsilon$  selects the largest eigenvalues. It leads to the following reduced iterative UWVF Galerkin problem. This problem involves less unknowns and is therefore less costly to solve numerically.

**Problem 7.** Let  $\mathbb{Y}_{\mathcal{T},\text{red}}^h$  be the linear subspace of  $\mathbb{Y}_{\mathcal{T}}^h$  parametrised by  $[\mathbf{x}_{\text{red}}] \in \mathbb{C}^{\#\text{dof}_{\text{red}}}$ . Any element  $\mathbf{x}_{\text{red}} \in \mathbb{Y}_{\mathcal{T},\text{red}}^h$  is given through

$$\mathbf{x} = \sum_{\text{iglob}=1}^{\#\text{dof}} [\mathbf{x}]_{\text{iglob}} \mathbf{w}^{\text{iglob}}, \quad \text{with } [\mathbf{x}] := \mathbf{T}_{\text{red}} \mathbf{\Lambda}_{\text{red}}^{-\frac{1}{2}} [\mathbf{x}_{\text{red}}]. \quad (44)$$

Find  $\mathbf{x}_{\text{red}} \in \mathbb{Y}_{\mathcal{T},\text{red}}^h$ , or equivalently  $[\mathbf{x}_{\text{red}}] \in \mathbb{C}^{\#\text{dof}_{\text{red}}}$ , such that we have

$$\mathbf{a}(\mathbf{x}_{\text{red}}, \mathbf{x}'_{\text{red}}) = \mathbf{l}(\mathbf{x}'_{\text{red}}), \quad \forall \mathbf{x}'_{\text{red}} \in \mathbb{Y}_{\mathcal{T},\text{red}}^h \iff [\mathbf{x}'_{\text{red}}]^* \mathbf{A}_{\text{red}} [\mathbf{x}_{\text{red}}] = [\mathbf{x}'_{\text{red}}]^* \mathbf{F}_{\text{red}}, \quad \forall [\mathbf{x}'_{\text{red}}] \in \mathbb{C}^{\#\text{dof}_{\text{red}}}, \quad (45)$$

with

$$\mathbf{A}_{\text{red}} := \mathbf{\Lambda}_{\text{red}}^{-\frac{1}{2}} \mathbf{T}_{\text{red}}^* \mathbf{A} \mathbf{T}_{\text{red}} \mathbf{\Lambda}_{\text{red}}^{-\frac{1}{2}} \quad \text{and} \quad \mathbf{F}_{\text{red}} := \mathbf{\Lambda}_{\text{red}}^{-\frac{1}{2}} \mathbf{T}_{\text{red}}^* \mathbf{F}. \quad (46)$$

Since  $\mathbb{Y}_{\mathcal{T},\text{red}}^h \subset \mathbb{Y}_{\mathcal{T}}^h$ , the convergence theory established in Proposition 3.1 remains true for the reduced system (45). Therefore, Theorem 4.1 leads to a convergent iterative reduced preconditioned Krylov UWVF method.

**Remark 4.3.** Let us notice that  $\mathbf{M}$  becomes the identity matrix of dimension  $\#\text{dof}_{\text{red}}$

$$\mathbf{I}_{\text{red}} := \mathbf{\Lambda}_{\text{red}}^{-\frac{1}{2}} \mathbf{T}_{\text{red}}^* \mathbf{M} \mathbf{T}_{\text{red}} \mathbf{\Lambda}_{\text{red}}^{-\frac{1}{2}},$$

which leads to a significant memory gain, particularly for large cases ie  $(D_{\Omega})^3 > 100^3$  in wavelength cube, since the iterative algorithm (34) becomes

$$[\mathbf{x}_{n+1}^{\text{jac}}] = \mathbf{N}_{\text{red}} [\mathbf{x}_n^{\text{jac}}] + \mathbf{F}_{\text{red}}, \quad \text{with } \mathbf{N}_{\text{red}} := \mathbf{\Lambda}_{\text{red}}^{-\frac{1}{2}} \mathbf{T}_{\text{red}}^* \mathbf{N} \mathbf{T}_{\text{red}} \mathbf{\Lambda}_{\text{red}}^{-\frac{1}{2}} \quad \text{and} \quad [\mathbf{x}_0^{\text{jac}}] = 0.$$

The number of basis functions after reduction, denoted by  $N_{\text{red}}$ , depends on the value of the threshold  $\varepsilon$  and on the size  $\mathcal{D}_T$  of each element, see Table 3. When the threshold  $\varepsilon$  is sufficiently small, ie  $\varepsilon \leq 10^{-4}$ , the convergence rate is the same, see Figure 7, leading to the same accuracy of the numerical solution visualised in Figure 6. As soon as  $\varepsilon \geq 10^{-3}$ , the convergence is faster in terms of iterations. But, the obtained numerical solution does not approximate the physical phenomenon anymore, see Figure 5 where  $\varepsilon = 10^{-2}$ . Consequently, to get a numerical solution of a given accuracy, the solver can use 16.5 GB instead of 275 GB without reduction, ie  $N_{\text{red}} = 48$  instead of  $N = 196$  basis functions, see Table 4. In any reduced configuration, the iterative or the direct solver is also faster in terms of time since we perform less costly matrix vector products, see Table 4.

**Remark 4.4.** Figures 5, 6, 7 and 12 are given for  $D_{\Omega} = 40$  and  $\mathcal{D}_T = 0.25$  in wavelength, and respect the configuration mentioned in Remark 3.5.

### 4.3. Global algebraic left-preconditioner

The non-local aspect of wave propagation encourages to set up a global preconditioner. A popular approach has been introduced in [55]. We propose an alternative which is based on the different subsets of faces of the cubic elements.

**Definition 4.1.** Let us consider an element  $T \in \mathcal{T}$ . We denote by

$\mathcal{D}_T \backslash \varepsilon$	$10^{-16}$	$10^{-15}$	$10^{-13}$	$10^{-11}$	$10^{-9}$	$10^{-7}$	$10^{-5}$	$10^{-4}$	$10^{-3}$	$10^{-2}$
0.25	180	175	154	126	96	70	48	36	30	16
0.5	196	196	190	186	174	132	96	84	70	48
1	196	196	196	196	196	190	180	174	148	114

TABLE 3. Values of  $N_{\text{red}}$  when reducing the basis of size  $N = 196$  according to  $\varepsilon$  and  $\mathcal{D}_T$  in wavelength.

$\varepsilon$	$10^{-16}$	$10^{-13}$	$10^{-11}$	$10^{-9}$	$10^{-7}$	$10^{-5}$	$10^{-4}$	$10^{-3}$	$10^{-2}$
$N_{\text{red}}$	180	154	126	96	70	48	36	30	16
Time (s)	4121	2936	1888	1085	496.1	254.9	141.8	99.68	25.61
Memory cost (Gb)	0.151	0.129	0.105	0.080	0.058	0.04	0.03	0.02	0.01

TABLE 4. Numerical data: basis reduction threshold  $\varepsilon$ , number of basis functions after reduction  $N_{\text{red}}$ , time to get the numerical solution, for the reduced Krylov solution, with a size of Krylov space set to 100, in function of the chosen basis reduction threshold  $\varepsilon$  for  $\mathcal{D}_\Omega = 5\lambda$ ,  $\mathcal{D}_T = 0.25$  and  $N = 196$ .

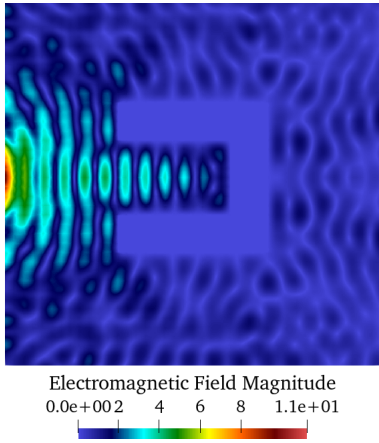


FIGURE 5. Magnitude of the electromagnetic field in a slice view of a three-dimensional cup for  $N_{\text{red}} = 24$ .

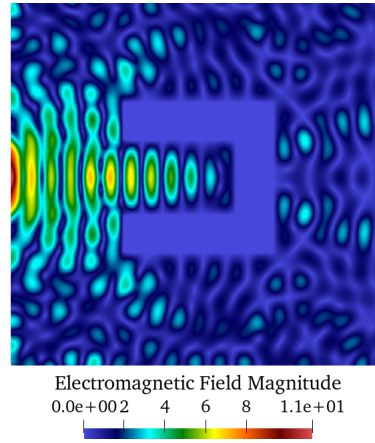


FIGURE 6. Magnitude of the electromagnetic field in a slice view of a three-dimensional cup for  $N_{\text{red}} = 48$ .

- (i)  $\mathcal{F}_x^T$  the set of the left and right faces of  $T$ ,
- (ii)  $\mathcal{F}_y^T$  the set of the front and back faces of  $T$ ,
- (iii)  $\mathcal{F}_z^T$  the set of the bottom and top faces of  $T$ .

The set  $\mathcal{T}$  is divided into one-dimensional subsets, either in  $x$ ,  $y$  or  $z$  direction, see Figure 9. These one-dimensional subdomains lead to three singular regular decompositions  $\mathbf{M} - \mathbf{N}$  of matrix  $\mathbf{A}$  (see [11]) which are associated to the following sesquilinear forms.

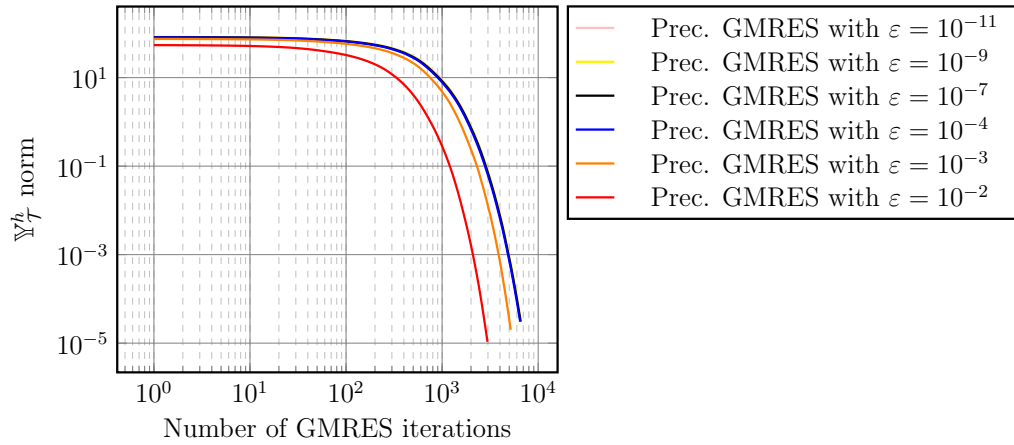


FIGURE 7. Comparison of the convergence rate of the preconditioned Krylov UWVF with basis reduction for different values of the threshold  $\varepsilon$ .

**Definition 4.2.** For all  $\mathbf{x} \in \mathbb{Y}_{\mathcal{T}}^h$  and for all  $\mathbf{x}' \in \mathbb{Y}_{\mathcal{T}}^h$ , we define the three sesquilinear forms  $\mathbf{k}^{\times/y/z}$  and their associated matrices  $\mathbf{N}^{\times/y/z} \in \mathbb{C}^{\#\text{dof} \times \#\text{dof}}$  as

$$\mathbf{k}^{\times/y/z}(\mathbf{x}, \mathbf{x}') := \sum_{T \in \mathcal{T}} \sum_{F \in \mathcal{F}_{\times/y/z}^T \cap \mathcal{F}_{\text{int}}} \left( \Pi_U \mathbf{x}, \mathcal{U}^T \mathbf{x}' \right)_{L_i^2(F)} = [\mathbf{x}']^* \mathbf{N}^{\times/y/z} [\mathbf{x}].$$

Moreover, we define the associated three regular matrices as  $\mathbf{M}^{\times/y/z} := \mathbf{A} + \mathbf{N}^{\times/y/z}$ .

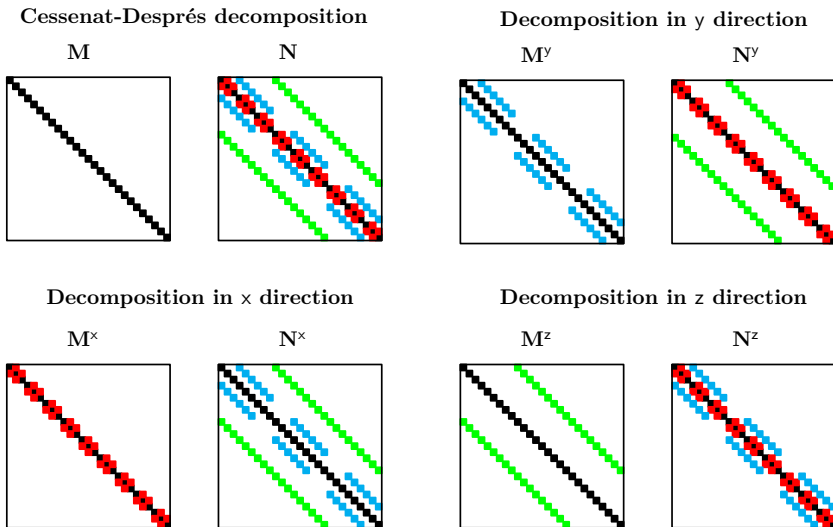


FIGURE 8. Matrices structures for the Cessenat-Després decomposition and the decompositions in  $x$ ,  $y$  or  $z$  direction, for  $\mathcal{D}_{\Omega} = 3$  in wavelength.

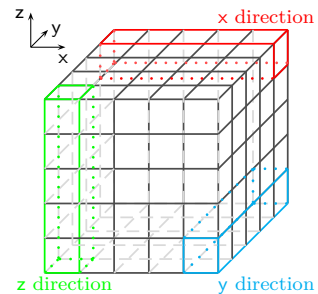


FIGURE 9. Global one-dimensional subdomains in a cube, for  $\mathcal{D}_{\Omega} = 5$  in wavelength.

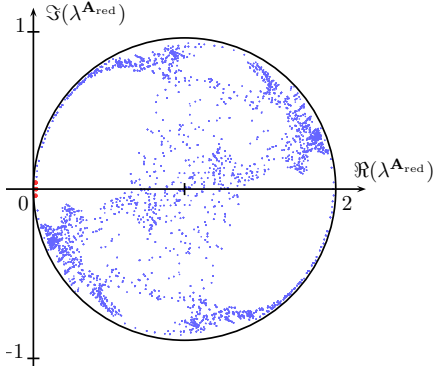


FIGURE 10. Real and complex parts, *resp.*  $\Re(\lambda^{\mathbf{A}_{\text{red}}})$  and  $\Im(\lambda^{\mathbf{A}_{\text{red}}})$ , of the spectrum  $\lambda^{\mathbf{A}_{\text{red}}}$  of  $\mathbf{A}_{\text{red}}$ , for  $\mathcal{D}_{\Omega} = 6$ .

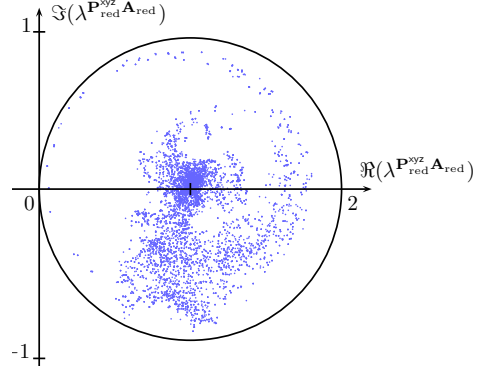


FIGURE 11. Real and complex parts, *resp.*  $\Re(\lambda^{\mathbf{P}_{\text{red}}^{\text{xyz}} \mathbf{A}_{\text{red}}})$  and  $\Im(\lambda^{\mathbf{P}_{\text{red}}^{\text{xyz}} \mathbf{A}_{\text{red}}})$ , of the spectrum  $\lambda^{\mathbf{P}_{\text{red}}^{\text{xyz}} \mathbf{A}_{\text{red}}}$  of  $\mathbf{P}_{\text{red}}^{\text{xyz}} \mathbf{A}_{\text{red}}$ , for  $\mathcal{D}_{\Omega} = 6$ .

Then we get three iterative schemes, leading to left preconditioners

$$\mathbf{M}^{x/y/z}[\mathbf{x}_{n+1}] = \mathbf{F} + \mathbf{N}^{x/y/z}[\mathbf{x}_n].$$

By applying successively each of these one-dimensional preconditioners, we get a global preconditioner involving the three directions  $x, y$  and  $z$ . The latter associates  $\mathbf{F}$  to  $[\mathbf{y}]$  thanks to intermediate solution vectors  $[\mathbf{x}_{n+1}^0]$  and  $[\mathbf{x}_{n+1}^1]$  through the following iterative scheme

$$\begin{cases} \mathbf{M}^x[\mathbf{x}_{n+1}^0] = \mathbf{F} + \mathbf{N}^x \mathbf{F}, \\ \mathbf{M}^y[\mathbf{x}_{n+1}^1] = \mathbf{F} + \mathbf{N}^y[\mathbf{x}_{n+1}^0], \\ \mathbf{M}^z \mathbf{y} = \mathbf{F} + \mathbf{N}^z[\mathbf{x}_{n+1}^1]. \end{cases}$$

Thus, we obtain an iterative method of the form  $[\mathbf{x}_{n+1}] = \mathbf{P}^{\text{xyz}} \mathbf{F} + \mathbf{Q}^{\text{xyz}}[\mathbf{x}_n]$ . Taking  $[\mathbf{x}_0] = 0$ , it leads to an approximation of the solution to the problem  $\mathbf{A}[\mathbf{x}] = \mathbf{F}$ . Therefore, this iterative method defines a global preconditioner that can be applied to the non preconditioned Problem 5,

$$\mathbf{P}^{\text{xyz}} \mathbf{A}[\mathbf{x}] = \mathbf{P}^{\text{xyz}} \mathbf{F}, \quad \text{where} \quad \mathbf{P}^{\text{xyz}} := (\mathbf{M}^z)^{-1} - (\mathbf{M}^z)^{-1} \mathbf{N}^z (\mathbf{M}^y)^{-1} + (\mathbf{M}^z)^{-1} \mathbf{N}^z (\mathbf{M}^y)^{-1} \mathbf{N}^y (\mathbf{M}^x)^{-1},$$

where we choose to use left-preconditioning. We can resort to the same method for the preconditioned Problem 6 or the reduced preconditioned Problem 7. We do not provide any theory ensuring the convergence of the reduced left-preconditioned Krylov solver associated to  $\mathbf{P}_{\text{red}}^{\text{xyz}}$ . However, it is well-known that few isolated eigenvalues do not cause any problem to the Krylov solver. In the numerous numerical cases that we considered, the spectrum of  $\mathbf{P}_{\text{red}}^{\text{xyz}} \mathbf{A}_{\text{red}}$  is concentrated around 1 (see Figure 11), removing the small eigenvalues of  $\mathbf{A}_{\text{red}}$  (in red in Figure 10) which are slowing down the convergence (see the blue curve in Figure 12). The results shown in Figure 12 point out the smaller amount of Krylov iterations when using the reduced left-preconditioned matrix  $\mathbf{P}_{\text{red}}^{\text{xyz}} \mathbf{A}_{\text{red}}$ .

## CONCLUSION

In this paper, we have introduced an heterogeneous iterative Trefftz method solving three-dimensional Maxwell equations. As all iterative methods, the present solver does not need to store the LU factorisation of  $\mathbf{A}$  and can be matrix-free. The latter speed-up accelerates the convergence such that it takes a few minutes



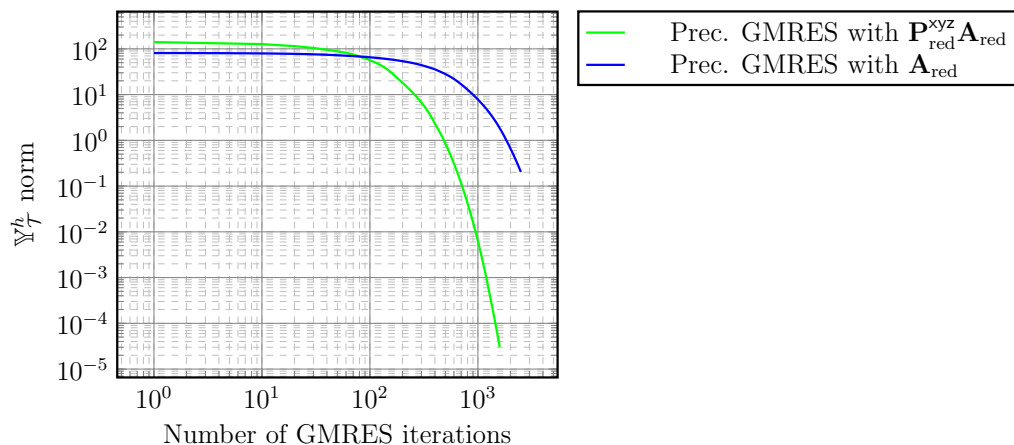


FIGURE 12. Comparison of the convergence rate of the preconditioned Krylov UWVF using either  $\mathbf{A}_{\text{red}}$  or  $\mathbf{P}_{\text{red}}^{\text{xyz}}\mathbf{A}_{\text{red}}$ .

to solve Maxwell equations on small heterogeneous domains, *i.e.* up to  $\mathcal{D}_{\Omega} = 40$ . This memory gain has been strengthened thanks to a basis reduction proposed in Section 4.2. This compression method turned out to be a successful development, reducing the memory cost of the heterogeneous preconditioned Krylov UWVF solver. It enables to consider large complex geometries, see Figure 13. Nevertheless, this basis reduction should be applied reasonably (*ie* a good choice is  $\varepsilon = 10^{-6}$ ) to use enough basis functions to describe the numerical solution. Moreover, the introduced Trefftz Krylov method is accelerated thanks to a new global preconditioner.

However, on some configurations, direct methods outperform iterative solvers since they can deal with many right-hand sides, see [23]. Judicious choices should then be made depending on the numerical case. But research on Krylov solvers efficiently dealing with multiple right-hand sides remains an active topic, see [8, 25]. It could then be preferable to use Krylov methods instead of direct solvers in the future.

To end this paper, let us consider an industrial application which consists of a radar echo on a boat. This large heterogeneous configuration is a boat surrounded by an homogeneous domain of size  $\mathcal{D}_{\Omega} = 24 \times 61 \times 154$  in wavelength and satisfying  $Z_{\partial\Omega} = 1$ . The boat has a perfect metal surface, such that the incident electromagnetic wave perfectly reflects when striking its top, see Figure 13. This large three-dimensional case is sped up thanks to a Cessenat-Després preconditioner. It also requires a restart strategy to fit into a 1 Terabyte Broadwell Intel Xeon CPU E5-2650 node. Our ways of checking the accuracy of this result are limited due to its dimensions, see Table 5. We ensure that the Krylov residual is  $10^{-3}$ . The large dimensions, and the heterogeneity of this numerical case both emphasize the robustness of the developed iterative Trefftz method.

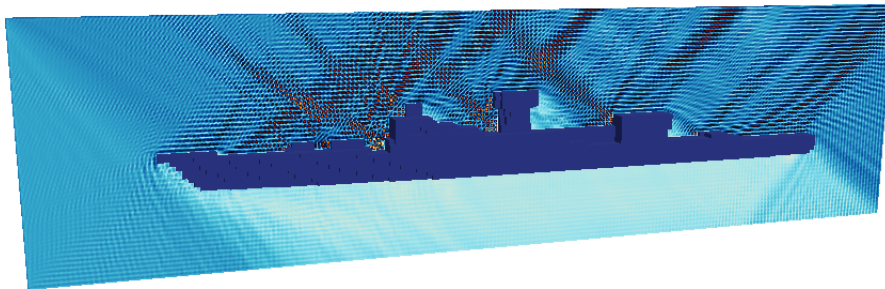


FIGURE 13. Visualisation of an electromagnetic wave striking the top of a boat.

#elem	$N_{\text{red}}$	$N$	$\#\text{dof}_{\text{red}}$	$\#\text{dof}$	Krylov iterations	Memory cost	Total duration
$> 14.4 \times 10^6$	46	52	$> 663 \times 10^6$	$> 750 \times 10^6$	473	$\approx 836$ GB	37.8 hours

TABLE 5. Numerical data and results associated to the boat case, see Figure 13.

## REFERENCES

- [1] M. Ainsworth. Dispersive properties of high-order Nédélec/edge Element approximation of the time-harmonic Maxwell equations. *Philosophical Transactions of the Royal Society of London. Series A: Mathematical, Physical and Engineering Sciences*, 362(1816):471–491, 2004.
- [2] M. Ainsworth, P. Monk, and W. Muniz. Dispersive and dissipative properties of Discontinuous Galerkin Finite Element Methods for the second-order wave equation. *Journal of Scientific Computing*, 27(1–3):5–40, 2006.
- [3] H. Barucq, A. Bendali, J. Diaz, and S Tordeux. Local strategies for improving the conditioning of the plane-wave Ultra Weak Variational Formulation. *Journal of Computational Physics*, 441:110449, September 2021.
- [4] M. Bebendorf. Approximation of Boundary Element matrices. *Numerische Mathematik*, 86(4):565–589, 2000.
- [5] O. P. Bruno and L. A. Kunyansky. A fast, high-order algorithm for the solution of surface scattering problems: basic implementation, tests, and applications. *Journal of Computational Physics*, 169(1):80–110, 2001.
- [6] A. Buffa and R. Hiptmair. Galerkin Boundary Element Methods for electromagnetic scattering. In *Topics in computational wave propagation*, pages 83–124. Springer, 2003.
- [7] A. Buffa and P. Monk. Error estimates for the Ultra Weak Variational Formulation of the Helmholtz equation. *Mathematical Modelling and Numerical Analysis*, 42(6):925–940, 2008.
- [8] H. Calandra, S. Gratton, R. Lago, X. Vasseur, and L. M. Carvalho. A modified block flexible GMRES method with deflation at each iteration for the solution of non-Hermitian linear systems with multiple right-hand sides. *SIAM Journal on Scientific Computing*, 35(5):S345–S367, 2013.
- [9] O. Cessenat. Application d’une nouvelle formulation variationnelle aux équations d’ondes harmoniques. *Problèmes d’Helmholtz 2D et de Maxwell 3D*. PhD thesis, University of Paris XI Dauphine, 1996.
- [10] O. Cessenat and B. Després. Application of an Ultra Weak Variational Formulation of elliptic PDE to the two-dimensional Helmholtz problem. *SIAM J. Num. Analysis*, 35(1):255–299, 1998.
- [11] P. G. Ciarlet. *Introduction à l’analyse numérique matricielle et à l’optimisation* PG Ciarlet. Masson, 1990.
- [12] B. Cockburn, J. Gopalakrishnan, and R. Lazarov. Unified Hybridization of Discontinuous Galerkin, mixed, and continuous Galerkin methods for second order elliptic problems. *SIAM Journal on Numerical Analysis*, 47(2):1319–1365, 2009.
- [13] G. Cohen. *Higher-order numerical methods for transient wave equations*. Springer Science & Business Media, 2001.
- [14] G. Cohen and S. Pernet. *Finite Element and Discontinuous Galerkin methods for transient wave equations*. Springer, 2017.
- [15] S. Congreve, J. Gedicke, and I. Perugia. Numerical investigation of the conditioning for plane wave Discontinuous Galerkin methods. In *European Conference on Numerical Mathematics and Advanced Applications*, pages 493–500. Springer, 2017.
- [16] E. Darve. The Fast Multipole Method: numerical implementation. *Journal of Computational Physics*, 160(1):195–240, 2000.
- [17] E. Darve and P. Havé. Efficient Fast Multipole Method for low-frequency scattering. *Journal of Computational Physics*, 197(1):341–363, 2004.
- [18] B. Després. Sur une Formulation Variationnelle Ultra Faible. *Comptes Rendus de l’Académie des Sciences, Série I* 318:939–944, 1994.
- [19] V. Dolean, H. Fol, S. Lanteri, and R. Perrusel. Solution of the time-harmonic Maxwell equations using Discontinuous Galerkin methods. *Journal of computational and applied mathematics*, 218(2):435–445, 2008.
- [20] V. Dolean, M. J. Gander, S. Lanteri, J.-F. Lee, and Z. Peng. Effective transmission conditions for Domain Decomposition Methods applied to the time-harmonic curl-curl Maxwell’s equations. *Journal of computational physics*, 280:232–247, 2015.
- [21] M. Embree. How descriptive are GMRES convergence bounds. *NA Report* 99, 8, 1999.
- [22] C. Farhat, R. Tezaur, and P. Weidemann-Goiran. Higher-order extensions of a Discontinuous Galerkin method for mid-frequency Helmholtz problems. *International journal for numerical methods in engineering*, 61(11):1938–1956, 2004.
- [23] F. Faucher and O. Scherzer. Adjoint-state method for Hybridizable Discontinuous Galerkin discretization, application to the inverse acoustic wave problem. *Computer Methods in Applied Mechanics and Engineering*, 372:113406, 2020.
- [24] G. Gabard. Discontinuous Galerkin methods with plane waves for time-harmonic problems. *Journal of Computational Physics*, 225:1961–1984, 2007.
- [25] L. Giraud, Y.-F. Jing, and Y.-F. Xiang. A block minimum residual norm subspace solver with partial convergence management for sequences of linear systems. *SIAM Journal on Matrix Analysis and Applications*, 2022.
- [26] C. Gittelsohn and R. Hiptmair. Dispersion analysis of plane wave discontinuous methods. *International Journal for Numerical Methods in Engineering*, 98(5):313–323, 2014.
- [27] C. Gittelsohn, R. Hiptmair, and I. Perugia. Plane wave Discontinuous Galerkin methods: analysis of the  $h$ -version. *Mathematical Modelling and Numerical Analysis*, 43:297–331, 2009.
- [28] J. S. Hesthaven and T. Warburton. *Nodal Discontinuous Galerkin methods: algorithms, analysis, and applications*. Springer Science & Business Media, 2007.

- [29] R. Hiptmair. Finite Elements in computational electromagnetism. *Acta Numerica*, 11:237–339, 2002.
- [30] R. Hiptmair, A. Moiola, and I. Perugia. Plane wave Discontinuous Galerkin methods for the 2-D Helmholtz equation: analysis of the p-version. *SIAM Journal on Numerical Analysis*, 49(1):264–284, 2011.
- [31] R. Hiptmair, A. Moiola, and I. Perugia. A survey of Trefftz methods for the Helmholtz equation. In *Building bridges: connections and challenges in modern approaches to numerical partial differential equations*, pages 237–279. Springer, 2016.
- [32] R. Hiptmair, A. Moiola, I. Perugia, and C. Schwab. Approximation by harmonic polynomials in star-shaped domains and exponential convergence of Trefftz *hp*-DGFEM. *Mathematical Modelling and Numerical Analysis*, 48:727–752, 2014.
- [33] P. Houston, I. Perugia, A. Schneebeli, and D. Schötzau. Interior Penalty method for the indefinite time-harmonic Maxwell equations. *Numerische Mathematik*, 100(3):485–518, 2005.
- [34] T. Huttunen, M. Malinen, and P. Monk. Solving Maxwell’s equations using the Ultra Weak Variational Formulation. *Journal of Computational Physics*, 223(2):731–758, 2007.
- [35] F. Ihlenburg and I. Babuska. Finite Element solution of the Helmholtz equation with high wave number – part i: the h-version of the FEM. *Computers Math. Applic.*, 30(9):9–37, 1995.
- [36] F. Ihlenburg and I. Babuska. Finite Element solution of the Helmholtz equation with high wave number – part ii: the h-p version of the FEM. *SIAM J. Numer. Anal.*, 34(1):315–358, 1997.
- [37] J. M. Jin. *The Finite Element Method in Electromagnetics, Second Edition*. John Wiley & Sons, New York, 2002.
- [38] S. Kurz, O. Rain, and S. Rjasanow. The adaptive cross-approximation technique for the 3-D Boundary Element Method. *IEEE transactions on Magnetics*, 38(2):421–424, 2002.
- [39] T. Luostari, T. Huttunen, and P. Monk. The Ultra Weak Variational Formulation using Bessel basis functions. *Communications in Computational Physics*, 11(2):400–414, 2012.
- [40] T. Luostari, T. Huttunen, and P. Monk. Improvements for the Ultra Weak Variational Formulation. *International Journal for Numerical Methods in Engineering*, 94(6):598–624, 2013.
- [41] A. Moiola, R. Hiptmair, and I. Perugia. Plane wave approximation of homogeneous Helmholtz solutions. *Z. Angew. Math. Phys.*, 62:809–837, 2011.
- [42] P. Monk. *Finite Element Methods for Maxwell’s Equations*. Numerical Analysis and Scientific Computations. Clarendon Press, 2003.
- [43] D. Moro, N. C. Nguyen, and J. Peraire. A Hybridized Discontinuous Petrov-Galerkin scheme for scalar conservation laws. *International journal for numerical methods in engineering*, 91(9):950–970, 2012.
- [44] J.-C. Nédélec. A new family of mixed Finite Elements in  $\mathbb{R}^3$ . *Numerische Mathematik*, 50(1):57–81, 1986.
- [45] J.-C. Nédélec. *Acoustic and electromagnetic equations: integral representations for harmonic problems*. Springer Science & Business Media, 2001.
- [46] N. C. Nguyen, J. Peraire, and B. Cockburn. Hybridizable Discontinuous Galerkin methods for the time-harmonic Maxwell’s equations. *Journal of Computational Physics*, 230(19):7151–7175, 2011.
- [47] S. Pernet, N. Serdiuk, M. Sirdey, and S. Tordeux. Discontinuous Galerkin method based on Riemann fluxes for the time domain Maxwell system, 2021.
- [48] E. Perrey-Debain. Plane wave decomposition in the unit disc: convergence estimates and computational aspects. *Journal of Computational and Applied Mathematics*, 193(1):140–156, 2006.
- [49] B. Pluymers, B. Hal, D. Vandepitte, and W. Desmet. Trefftz-based methods for time-harmonic acoustics. *Archives of Computational Methods in Engineering*, 14:343–381, 01 2007.
- [50] Y. Saad. Krylov subspace methods for solving large unsymmetric linear systems. *Mathematics of computation*, 37(155):105–126, 1981.
- [51] Y. Saad. *Iterative Methods for Sparse Linear Systems*. PWS Publishing Company, Boston, 1996.
- [52] Y. Saad and M. H. Schultz. GMRES: A Generalized Minimal RESidual algorithm for solving non-symmetric linear systems. *SIAM Journal on scientific and statistical computing*, 7(3):856–869, 1986.
- [53] S. A. Sauter and C. Schwab. *Boundary Element Methods*. Springer-Verlag, Berlin-Heidelberg, 2011.
- [54] A. Taflove. *Computational Electrodynamics. The Finite-Difference Time-Domain Method*. Artech house, Inc., Norwood, MA 02062, 1995.
- [55] A. Vion and C. Geuzaine. Double sweep preconditioner for optimized Schwarz methods applied to the Helmholtz problem. *J. Comput. Phys.*, 266:171–190, 2014.
- [56] D. Wang, R. Tezaur, J. Toivanen, and C. Ferhat. Overview of the Discontinuous Enrichment Method, the Ultra Weak Variational Formulation, and the Partition of Unity Method for the acoustic scattering in the medium frequency regime and performance comparisons. *International Journal for Numerical Methods in Engineering*, 89:403–417, 2012.
- [57] N. Zerbib. *Méthodes de Sous-Structuration et de Décomposition de Domaine pour la Résolution des Équations de Maxwell : Application au Rayonnement d’antenne dans un Environnement Complexe*. PhD thesis, National Institute for Applied Sciences (INSA), INSA Toulouse, 2006.
- [58] K. Zhao, M. N. Vouvakis, and J.-F. Lee. The adaptive cross approximation algorithm for accelerated method of moments computations of EMC problems. *IEEE transactions on electromagnetic compatibility*, 47(4):763–773, 2005.

- [59] L. Zhu, E. Burman, and H. Wu. Continuous Interior Penalty Finite Element Method for Helmholtz equation with high wave number: One dimensional analysis. preprint available at [arXiv:1211.1424](https://arxiv.org/abs/1211.1424).