



# A Modified Discontinuous Galerkin Method for Solving Helmholtz Problems

Mohamed Amara, Henri Calandra, Rabia Djellouli, Magdalena Grigoroscuta-Strugaru

## ► To cite this version:

Mohamed Amara, Henri Calandra, Rabia Djellouli, Magdalena Grigoroscuta-Strugaru. A Modified Discontinuous Galerkin Method for Solving Helmholtz Problems. [Research Report] RR-7050, INRIA. 2009, pp.30. inria-00421584v4

**HAL Id: inria-00421584**

**<https://inria.hal.science/inria-00421584v4>**

Submitted on 11 May 2011

**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.



INSTITUT NATIONAL DE RECHERCHE EN INFORMATIQUE ET EN AUTOMATIQUE

# *A Modified Discontinuous Galerkin Method for Solving Helmholtz Problems*

Mohamed Amara — Henri Calandra — Rabia Djellouli —

Magdalena Grigoroscuta-Strugaru

N° 7050

October 2009

---

A large, light blue stylized 'R' logo is positioned to the left of the text 'Rapport de recherche'.

*Rapport  
de recherche*



## A Modified Discontinuous Galerkin Method for Solving Helmholtz Problems

Mohamed Amara<sup>\*</sup>, Henri Calandra<sup>†</sup>, Rabia Djellouli<sup>‡</sup>,  
Magdalena Grigoroscuta-Strugaru<sup>§</sup>

Thème : Modélisation Avancée en GéophysIQUE 3D  
Équipe-Projet Magique-3D

Rapport de recherche n° 7050 — October 2009 — 31 pages

**Abstract:** A new solution methodology is proposed for solving Helmholtz problems in the mid- and high frequency regime. The proposed method falls in the category of the discontinuous Galerkin methods. The primal variable is obtained by solving in parallel a set of well posed local problems. The Lagrange multiplier is the solution of a global positive definite linear system. These two properties are the main features of the proposed method that distinguish it from the existing solution methodologies. Numerical results are presented to illustrate both the stability and the accuracy of the proposed method when applied for solving waveguide-type problems.

**Key-words:** Helmholtz equation, discontinuous Galerkin, plane waves, Lagrange multipliers, *positive definite matrix*, inf-sup condition, waveguide problems

<sup>\*</sup> LMA/CNRS UMR 5142, Université de Pau et des Pays de l'Adour, France  
INRIA Bordeaux Sud-Ouest Research Center, Team Project Magique-3D  
e-mail: mohamed.amara@univ-pau.fr

<sup>†</sup> TOTAL, Avenue Larribau, Pau, France  
e-mail: henri.calandra@total.com

<sup>‡</sup> Department of Mathematics, California State University Northridge and  
INRIA Bordeaux Sud-Ouest Research Center, Associate Team Project MAGIC, USA  
e-mail: rabia.djellouli@csun.edu

<sup>§</sup> INRIA Bordeaux Sud-Ouest Research Center, Team Project Magique-3D and  
LMA/CNRS UMR 5142, Université de Pau et des Pays de l'Adour, France  
e-mail: magdalena.grigoroscuta@inria.fr

# Une méthode modifiée de type Galerkin discontinu pour la résolution des problèmes de Helmholtz

**Résumé :** Nous proposons une nouvelle technique de résolution numérique des problèmes de Helmholtz en régime moyenne et haute fréquence. La méthode proposée est de type Galerkin discontinu. Elle diffère des techniques de résolution existantes essentiellement par les deux aspects suivants: (a) la variable primaire (le champ) est déterminée par la résolution en parallèle d'une classe des problèmes locaux bien posés au sens de Hadamard; (b) les multiplicateurs de Lagrange sont solutions d'un système linéaire global dont la matrice est définie positive. Nous présentons des résultats numériques pour illustrer la stabilité et la précision de la méthode proposée.

**Mots-clés :** équation de Helmholtz, Galerkin discontinu, ondes planes, multiplicateurs de Lagrange, *matrice définie positive*, condition inf-sup, guide d'ondes

## 1 Introduction

The Helmholtz equation belongs to the classical equations of mathematical physics that are well understood from a mathematical view point. However, the numerical approximation of the solution is still a challenging problem in spite the tremendous progress made during the last years (see, for example, the recent monograph [15] and the references therein). Indeed, the standard finite element method (FEM) is not well suited for solving Helmholtz problems in the mid- and high-frequency regime since highly oscillating solutions are not accurately approximated by piecewise polynomials. This phenomenon, related to the indefiniteness of the Helmholtz operator, is known as the pollution effect [3]. In order to maintain a certain level of accuracy, a mesh refinement is required and/or higher order FEM are used, leading to a prohibitive computational cost for high wavenumbers.

In response to this challenge, alternative techniques for alleviating the pollution effect were proposed. Numerous of these approaches use the plane waves, since they are expected to better approximate highly oscillating waves. Examples of such methods include the weak element method for Helmholtz equation [17], the Galerkin least-squares method [10], the partition of unity method [2], the residual free bubbles method [8], the least-squares method [16], the ultra-weak variational method [4] and recently, the discontinuous Galerkin method (DGM) designed by Farhat et al and presented in a series of papers [5, 6, 7]. In the latter method, the solution is approximated at the element mesh level using a superposition of plane waves which results in a discontinuous solution along interior boundaries of the mesh. The continuity is then restored in a weak sense through the use of Lagrange multipliers. The rectangular and quadrilateral elements constructed in [5, 6, 7] clearly outperform the standard Galerkin FEM. For example, for  $ka \geq 10$  and for a fixed level of accuracy, the so-called  $R$ -4-1 element reduces the total number of degrees of freedom (dofs) required by the  $Q$ 1-based finite element discretization for Helmholtz equation by a factor greater or equal to five. Similar results are obtained for the  $R$ -8-2a and  $R$ -8-2b elements when compared to the  $Q$ 2 element, and for  $Q$ -16-4 and  $Q$ -32-8 when compared to the  $Q$ 4 element. In spite of this impressive performance, the DGM has three important drawbacks. First, the method has to satisfy an *inf-sup* condition which is translated, in practice, as a compatibility requirement: the number of dofs of the Lagrange multiplier (corresponding to the dual variable) and of the field (the primal variable) cannot be chosen arbitrarily. The problem here is that there is no theoretical result on how to satisfy this compatibility requirement, except for the simple case of  $R$ -4-1 element (see [1]). Hence, for other elements, the existing choices are based on numerical experiments only. The second major issue with the DGM is that it becomes unstable as we refine the mesh. Such instabilities occur because of the singularity of the local problems and, to some extent, to the loss of the linear independence of the plane waves as the step size mesh discretization tends to zero. The latter affects dramatically the stability of the global system due to its ill-conditioning nature. Finally, the DGM exhibits a loss of accuracy for unstructured mesh [6].

We propose a new solution methodology for Helmholtz problems, that falls in the category of discontinuous Galerkin methods. The proposed formulation

distinguishes itself from existing procedures by the *well-posed* character of the local problems and by the resulting global system which is associated with a positive definite Hermitian matrix. More specifically, the computation domain is subdivided in quadrilateral- or triangular-shaped elements. The solution is approximated, at the element level, by a superposition of plane waves that are solution of the Helmholtz equation. The continuity of the solution at the interior interfaces of the elements is then enforced by Lagrange multipliers. Unlike the DGM, the proposed method does not require the continuity of the normal derivative. Consequently, Lagrange multipliers are introduced to restore in the weak sense the continuity of both the field and its normal derivative across interior boundaries of the mesh. Such choice leads to solving (a) local boundary value problems that are well posed in the sense of Hadamard [9] and (b) a global system, whose unknowns are the Lagrange multipliers. The Lagrange multiplier is the solution of a variational problem whose bilinear form can be written into two equivalent expressions. The approximation of both formulations leads to two linear systems corresponding to (a) a non-negative matrix and (b) a positive definite matrix. Note that the proposed technique is a two-step procedure where the local problems are first solved and then the Lagrange multipliers are evaluated. This two-step approach allows us to consider equally structured and unstructured meshes with either triangular- or quadrilateral-shaped elements. Since the proposed solution methodology resembles in some aspect the DGM, we will refer to it as mDGM (*modified Discontinuous Galerkin Method*).

The remainder of the paper is organized as follows. In Section 2, we introduce general notations and the model problem. Section 3 is devoted to the presentation of mDGM. In Section 4, we present the algebraic framework of the formulation. We compare in Section 5 the numerical performance of both methods: DGM and mDGM. The obtained results clearly indicate that mDGM outperforms DGM in terms of stability and accuracy. Finally, Section 6 concludes this paper.

## 2 Preliminaries

In this section, we introduce the model problem and specify the nomenclature and assumptions adopted throughout this paper.

### 2.1 The mathematical model

We consider the following class of waveguide-type problems:

$$(\text{BVP}) \begin{cases} -\Delta u - k^2 u = f & \text{in } \Omega, \\ \partial_n u = iku + g & \text{on } \partial\Omega \end{cases}$$

where  $\Omega \subset \mathbb{R}^2$  is an open bounded region, with smooth boundary  $\partial\Omega$ ,  $k$  is a positive number representing the wavenumber,  $\partial_n$  is the normal derivative and  $f$  and  $g$  are regular complex valued functions defined respectively on  $\Omega$  and  $\partial\Omega$ . The second equation of BVP is a representation of a class of non-homogeneous Robin boundary conditions, but other types of boundary condition can be considered.

Note that BVP is considered here for its simplicity since it allows us to compute analytically the solution  $u$  for a suitable choice of  $\Omega$ ,  $f$  and  $g$ . Such an expression of  $u$  is used when assessing the accuracy of mDGM.

## 2.2 Nomenclature and assumptions

In what follows, we consider a regular triangulation  $\tau_h$  of  $\Omega$  into quadrilateral- or triangular-shaped subdomains  $K$  whose boundaries are denoted by  $\partial K$ . The step size mesh discretization is denoted by  $h$ . We introduce the space of the primal variable:

$$\mathcal{V} = \{v \in L^2(\Omega); v|_K \in H^1(K)\},$$

that we equip with the norm:

$$\|v\|_{\mathcal{V}} = \left( \sum_{K \in \tau_h} \|v^K\|_{1,K}^2 \right)^{\frac{1}{2}}, \quad \forall v \in \mathcal{V},$$

where  $\|\cdot\|_{1,K}$  is the  $H^1$ -norm on the element  $K$ . In addition, we introduce  $\|\cdot\|_{0,K}$  and  $|\cdot|_{1,K}$  to designate the  $L^2$ -norm and the  $H^1$ -seminorm respectively on the element  $K$ .

Note that  $\mathcal{V}$  contains functions that are discontinuous across interior boundaries since their regularity is only  $L^2(\Omega)$ . For any  $v \in \mathcal{V}$ , we define the jump across an interior edge  $e = \partial K \cap \partial K'$  by:

$$[v] = v^K - v^{K'}.$$

We introduce the space of the dual variable, corresponding here to Lagrange multipliers, by:

$$\mathcal{M} = \left\{ \mu \in \prod_{K \in \tau_h} L^2(\partial K); \mu = 0 \text{ on } \partial K \cap \partial \Omega \right\}$$

and we associate to  $\mathcal{M}$  the norm given by:

$$\|\mu\|_{\mathcal{M}} = \left( \sum_{K \in \tau_h} \|\mu^K\|_{0,\partial K}^2 \right)^{\frac{1}{2}}, \quad \forall \mu \in \mathcal{M},$$

where  $\mu^K$  designates the restriction of  $\mu$  to  $\partial K$ :  $\mu^K = \mu|_{\partial K}$  and  $\|\cdot\|_{0,\partial K}$  is the  $L^2$ -norm on  $\partial K$ . Moreover, for any function  $\mu \in \mathcal{M}$ , we define the jump across an interior edge  $e = \partial K \cap \partial K'$  by:

$$[[\mu]] = \mu^K + \mu^{K'}.$$

## 3 The continuous approach

The basic idea of mDGM is to evaluate  $u$ , the solution of BVP, using the following splitting:

$$u = \Phi(\lambda) + \varphi, \tag{1}$$

where  $\varphi$  and  $\Phi$  are elements of  $\mathcal{V}$  and  $\lambda \in \mathcal{M}$ .

To compute these three quantities, we proceed into two steps:



**Step 1 :** For all  $K \in \tau_h$  and  $\mu \in \mathcal{M}$ , we compute  $\varphi$  and  $\Phi(\mu)$ . This is achieved by solving local Helmholtz problems. This step is called the restriction procedure.

**Step 2 :** We determine  $\lambda \in \mathcal{M}$  by solving a global linear system to ensure the continuity in a weak sense of the solution  $u$  given by (1) and of the normal derivative of  $u$ . This step is called the optimization procedure.

### 3.1 Step 1: The restriction procedure

As stated earlier, this step is devoted to the computation of  $\varphi$  and  $\Phi(\mu)$ , for all  $\mu \in \mathcal{M}$ , by solving locally Helmholtz problems. More specifically, for all  $K \in \tau_h$ , we compute  $\varphi^K$  by solving the following boundary value problem:

$$(\text{BVP1}) \begin{cases} \text{Find } \varphi^K \in H^1(K) \text{ such that:} \\ -\Delta \varphi^K - k^2 \varphi^K = f & \text{in } K \\ \partial_n \varphi^K = i k \varphi^K + g & \text{on } \partial K \cap \partial \Omega \\ \partial_n \varphi^K = i \alpha \varphi^K & \text{on } \partial K \cap \dot{\Omega} \end{cases}.$$

Next, for all  $\mu \in \mathcal{V}$  and  $K \in \tau_h$ , we compute  $\Phi(\mu^K)$  by solving the boundary value problem given by:

$$(\text{BVP2}) \begin{cases} \text{Find } \Phi(\mu^K) \in H^1(K) \text{ such that:} \\ -\Delta \Phi(\mu^K) - k^2 \Phi(\mu^K) = 0 & \text{in } K \\ \partial_n \Phi(\mu^K) = i k \Phi(\mu^K) & \text{on } \partial K \cap \partial \Omega \\ \partial_n \Phi(\mu^K) = i \alpha \Phi(\mu^K) + \mu^K & \text{on } \partial K \cap \dot{\Omega} \end{cases},$$

with  $\alpha \in \mathbb{R}_+^*$ . Note that the presence of  $\alpha$  ensures the uniqueness of the solution of BVP 1 and BVP 2, as it will be shown later.

It is easy to verify that the variational formulation of both problems can be expressed in a compact form as follows:

$$\begin{cases} \text{Find } \Psi^K \in H^1(K) \text{ such that:} \\ a_K(\Psi^K, v^K) = L_K(v^K) \end{cases} \quad \forall v^K \in H^1(K) \quad (2)$$

where  $a_K(\cdot, \cdot)$  is a bilinear form given by:

$$a_K(v^K, w^K) = \int_K \nabla v^K \cdot \nabla \overline{w^K} dx - k^2 \int_K v^K \overline{w^K} dx - i \alpha \int_{\partial K \cap \dot{\Omega}} v^K \overline{w^K} ds - i k \int_{\partial K \cap \partial \Omega} v^K \overline{w^K} ds, \quad \forall v^K, w^K \in H^1(K) \quad (3)$$

and  $\Psi_K$  is either the solution of BVP1 or BVP2, i.e.:

$$\Psi^K = \begin{cases} \varphi^K & \text{for BVP1} \\ \Phi(\mu^K) & \text{for BVP2, } \forall \mu \in \mathcal{M}. \end{cases}$$

The right-hand side  $L_K(\cdot)$  is given by:

$$L_K(v^K) = \begin{cases} \int_K f v^K dx + \int_{\partial K \cap \partial \Omega} g \overline{v^K} ds & \text{for BVP1} \\ \int_{\partial K \cap \dot{\Omega}} \mu^K \overline{v^K} ds & \text{for BVP2} \end{cases}, \quad \forall v^K \in H^1(K). \quad (4)$$

Consequently, the solving of BVP1 and BVP2 requires solving *one* linear system with multiple right-hand side.

Note that the bilinear form  $a_K$  is neither Hermitian, nor symmetric. However, it is easy to check that  $a_K(\cdot, \cdot)$  is continuous on  $H^1(K) \times H^1(K)$  and satisfies the Gårding inequality in  $H^1(K)$  since

$$\Re a_K(v^K, v^K) + k^2 \|v^K\|_{0,K}^2 = |v|_{1,K}^2, \quad (5)$$

where  $\Re$  designates the real part. In addition, we have:

**Proposition 1.** For a fixed  $K \in \tau_h$ , the variational problem (2) admits a unique solution.

**Proof of Proposition 1.** Let  $K$  be a fixed element of  $\tau_h$ . From (5), it follows that the bilinear form  $a_K(\cdot, \cdot)$  satisfies the Fredholm alternative. Hence, the uniqueness ensures the existence of  $\Psi^K \in H^1(K)$ , solution of (14). To prove the uniqueness, we consider the homogeneous problem associated to the bilinear form  $a_K(\cdot, \cdot)$  and let  $w^K$  be its solution. We therefore have:

$$\begin{cases} \text{Find } w^K \in H^1(K) \text{ such that:} \\ a_K(w^K, v^K) = 0 \end{cases} \quad \forall v^K \in H^1(K) \quad (6)$$

In particular, for  $v^K = w^K$ , we have:

$$\alpha \|w^K\|_{0,\partial K \cap \hat{\Omega}}^2 + k \|w^K\|_{0,\partial K \cap \partial \Omega}^2 = 0.$$

Since  $\alpha > 0$ , we must have:

$$w^K = 0 \text{ on } \partial K \quad \text{and} \quad \partial_n w^K = 0 \text{ on } \partial K.$$

Therefore, using the continuation theorem [12, 20], we deduce that  $w^K = 0$  in  $K$  and the problem (2) has a unique solution.  $\square$

**Remark 1.** The presence of the Robin condition on  $\partial K \cap \hat{\Omega}$  with  $\alpha > 0$  is *crucial* to ensure that 0 is the only solution of the variational problem (6). Indeed, if  $\alpha = 0$  and  $K$  satisfies  $\partial K \cap \partial \Omega = \emptyset$ , then the resulting homogeneous Neumann boundary condition on  $\partial K$  is not sufficient to guarantee the uniqueness of the solution of (6), since  $k^2$  may become an interior eigenvalue.

Next, we define  $\varphi$  such that, for all element  $K$  in the mesh, the restriction of  $\varphi$  to  $K$  is  $\varphi^K$ , the solution of BVP1, i.e.  $\varphi|_K = \varphi^K$ . Similarly, for all element  $K$  and for all  $\mu$  in  $\mathcal{M}$ , we define  $\Phi(\mu)$  such that we have  $\Phi(\mu)|_K = \Phi(\mu^K)$ , where  $\Phi(\mu^K)$  is the solution of BVP2. Using the definition of  $\varphi$  and  $\Phi(\mu)$ ,  $\forall \mu \in \mathcal{M}$  we have:

$$\varphi \in \mathcal{V} \quad \text{and} \quad \Phi(\mu) \in \mathcal{V}, \quad \forall \mu \in \mathcal{M} \quad (7)$$

In summary, Step 1 allows us to compute, for all  $\mu$  in  $\mathcal{M}$ :

$$\varphi + \Phi(\mu) \in \mathcal{V} \quad (8)$$

by solving one variational problem given by (2) with different right-hand side given by (4). Step 1 can be viewed, to some extent, as a prediction step.

### 3.2 Step 2: The optimization procedure

The objective here is to determine  $\lambda \in \mathcal{M}$  for which the function given by (8) is in  $H^1(\Omega)$ . This requirement can be viewed as a correction stage since we select the best-fit Lagrange multiplier  $\lambda$ . The determination of  $\lambda$  is accomplished by solving the following global variational problem:

$$(VF) \begin{cases} \text{Find } \lambda \in \mathcal{M} \text{ such that} \\ A(\lambda, \mu) = F(\mu), \quad \forall \mu \in \mathcal{M}, \end{cases} \quad (9)$$

where the bilinear form  $A(\cdot, \cdot)$  is given by:

$$\begin{aligned} A(\eta, \mu) = & \sum_{e \text{ interior edge}} \beta_e \int_e [\Phi(\eta)] [\overline{\Phi(\mu)}] ds \\ & + \sum_{\substack{e \text{ interior edge} \\ e = \partial K \cap \partial K'}} \gamma_e \int_e \left( \eta^K + \eta^{K'} + i\alpha \left( \Phi(\eta^K) + \Phi(\eta^{K'}) \right) \right) \\ & \left( \overline{\mu^K + \mu^{K'} + i\alpha \left( \Phi(\mu^K) + \Phi(\mu^{K'}) \right)} \right) ds \end{aligned} \quad (10)$$

and the linear form  $F(\cdot)$  is given by:

$$\begin{aligned} F(\mu) = & - \sum_{e \text{ interior edge}} \beta_e \int_e [\varphi] [\overline{\Phi(\mu)}] ds \\ & - i\alpha \sum_{\substack{e \text{ interior edge} \\ e = \partial K \cap \partial K'}} \gamma_e \int_e \left( \varphi^K + \varphi^{K'} \right) \\ & \left( \overline{\mu^K + \mu^{K'} + i\alpha \left( \Phi(\mu^K) + \Phi(\mu^{K'}) \right)} \right) ds. \end{aligned} \quad (11)$$

The parameters  $\beta_e$  and  $\gamma_e$  are two positive numbers that can be viewed as weight parameters. This problem expresses the continuity in the weak sense of the solution and its normal derivative. Note that the bilinear form  $A$  is Hermitian. Consequently, only half of the corresponding matrix will be stored.

**Remark 2.** Unlike the DGM, where only the primal variable  $u$  is discontinuous, the mDGM leads to the discontinuity of both variables: the primal variable  $u$  and the Lagrange multiplier  $\lambda$ , the dual variable. Consequently, the normal derivative of  $u$  is *discontinuous*. Alternatively, for numerical approximation purpose, we rewrite Eq. (10) as follows:

$$\begin{aligned} A(\eta, \mu) = & \sum_{e \text{ interior edge}} \beta_e \int_e [\Phi(\eta)] [\overline{\Phi(\mu)}] ds \\ & + \sum_{e \text{ interior edge}} \gamma_e \int_e [[\partial_n \Phi(\eta)]] [[\overline{\partial_n \Phi(\mu)}]] ds \\ & + \sum_{e \subset \partial \Omega} \omega_e \int_e (\partial_n \Phi(\eta) - i k \Phi(\eta)) (\overline{\partial_n \Phi(\mu) - i k \Phi(\mu)}) ds, \end{aligned} \quad (12)$$

where the weight parameter  $\omega_e$  is a positive number. Note that the second integral in Eq. (12) is equal to the second integral in Eq. (10), whereas the third

integral in Eq. (12) is in fact equal to 0. Consequently, the right-hand side given by (11) is modified depending on the use of Eq. (10) or (12).

The next result states the equivalence between solving BVP and solving the problem arising in the proposed two-step procedure.

**Theorem 1.**

- (i) Let  $u = \Phi(\lambda) + \varphi$ , where for all  $K$ ,  $\varphi^K$  is solution of BVP1 and  $\Phi(\lambda^K)$  is solution of BVP2 with  $\lambda$  solution of VF. Then  $u$  is the unique solution of BVP.
- (ii) Conversely, let  $u$  be the solution of BVP. For each  $K \in \tau_h$ , we define  $\lambda$  by:

$$\lambda^K = \begin{cases} 0 & \text{on } e \subset \partial K \cap \partial\Omega \\ \partial_n u^K - i\alpha u^K & \text{on } e \subset \partial K \cap \mathring{\Omega} \end{cases} \quad (13)$$

Let  $\varphi^K$  be the solution of BVP1 and  $\Phi(\lambda^K)$  the solution of BVP2. Then  $\lambda$  is solution of VF and  $u = \Phi(\lambda) + \varphi$ .

**Remark 3.** Eq. (13) indicates a clear distinction between mDGM and DGM, in which  $\lambda^K = \partial_n u^K$  and is continuous along the interior edges.

## 4 The algebraic approach

The implementation of mDGM requires first to introduce two finite-dimensional spaces  $\mathcal{V}_h$  and  $\mathcal{M}_h$  such that  $\mathcal{V}_h \subset \mathcal{V}$  and  $\mathcal{M}_h \subset \mathcal{M}$ . Similarly to the DGM formulation, we have considered in this paper spaces of plane waves functions. However, other shape functions satisfying the Helmholtz equation can also be considered. Moreover, unlike the DGM, mDGM allows - in principle - to choose the spaces  $\mathcal{V}_h$  and  $\mathcal{M}_h$  independently.

For any element  $K \in \tau_h$ , we denote by  $\mathcal{V}_h(K)$  (resp.  $\mathcal{M}_h(K)$ ) the set of functions of  $\mathcal{V}_h$  (resp.  $\mathcal{M}_h$ ) restricted to  $K$  (resp.  $\partial K$ ). Furthermore,  $n^K$  (resp.  $n^{\lambda^K}$ ) denotes the dimension of  $\mathcal{V}_h(K)$  (resp.  $\mathcal{M}_h(K)$ ). Last, the dimension of  $\mathcal{M}_h$ , which corresponds to the total number of dofs, is denoted by  $n^\lambda$ .

We show that when formulated in finite dimensional spaces, the proposed two-step procedure consists in solving linear algebraic systems in each step. Note that in Step 2, the resulting linear system to be solved depends on the approximation of the continuous formulations (10) and (12) respectively. As stated earlier, equations (10) and (12) are equivalent only at the continuous level. At the discontinuous level, the second and the third equation of BVP2 are satisfied in the weak sense.

### 4.1 Step 1: the restriction procedure

For an element  $K \in \tau_h$  and for any  $\mu_h^K \in \mathcal{M}_h(K)$ , we denote by  $\varphi_h^K \in \mathcal{V}_h(K)$  and  $\Phi_h(\mu_h^K) \in \mathcal{V}_h(K)$  the approximation of  $\varphi^K$  and  $\Phi(\mu_h^K)$  respectively. Similarly to the continuous formulation,  $\varphi_h$ ,  $\Phi_h(\mu_h)$  and  $\mu_h$  are given by:  $\varphi_h|_K = \varphi_h^K$ ,  $\Phi_h(\mu_h)|_K = \Phi_h(\mu_h^K)$  and  $\mu_h|_K = \mu_h^K$ , for any element  $K$  in the mesh.

To compute  $\varphi_h$  and  $\Phi_h(\mu_h)$ , for all  $K \in \tau_h$ , we set the variational problem (2) in the finite dimensional space  $\mathcal{V}_h(K)$ , that is:

$$\begin{cases} \text{Find } \Psi_h^K \in \mathcal{V}_h(K) \text{ such that:} \\ a_K(\Psi_h^K, v_h^K) = L_K(v_h^K), \quad \forall v_h^K \in \mathcal{V}_h(K) \end{cases} \quad (14)$$

where the forms  $a_K(\cdot, \cdot)$  and  $L_K(\cdot)$  are given by (3) and (4) respectively, and

$$\Psi_h^K = \begin{cases} \varphi_h^K & \text{for BVP1} \\ \Phi_h(\mu_h^K) & \text{for BVP2,} \quad \forall \mu_h \in \mathcal{M}_h. \end{cases} \quad (15)$$

Consequently, the variational problem (14)-(15) can be written in the following matrix form:

$$\left( \mathbf{K}^K - k^2 \mathbf{M}^K - i\alpha \mathbf{S}^{\partial K \cap \tilde{\Omega}} - ik \mathbf{S}^{\partial K \cap \partial \Omega} \right) \mathbf{X}^K = \text{rhs}, \quad (16)$$

where  $\mathbf{K}^K$  (resp.  $\mathbf{M}^K$ ) is the stiffness (resp. mass) matrix at the element level  $K$ .  $\mathbf{S}^{\partial K \cap \tilde{\Omega}}$  and  $\mathbf{S}^{\partial K \cap \partial \Omega}$  are mass-like matrices defined on  $\partial K \cap \tilde{\Omega}$  and  $\partial K \cap \partial \Omega$  respectively.  $\mathbf{X}^K$  is the vector in  $\mathbb{C}^{n^K}$  whose components are the values of  $\Psi_h^K$  in the basis of  $\mathcal{V}_h(K)$ .

The linear system (16) possesses the following properties:

- All the entries of the corresponding matrix can be evaluated analytically for plane waves shape functions.
- The linear system admits a unique solution, even when  $\partial K \cap \partial \Omega = \emptyset$ . Thanks to the positive number  $\alpha$  since the presence of the matrix  $\mathbf{S}^{\partial K \cap \tilde{\Omega}}$  guarantees the invertibility of the system. Note that this is not the case for the DGM, for which  $\mathbf{S}^{\partial K \cap \tilde{\Omega}}$  does not appear, leading to possibly a (weakly) singular system when  $\partial K \cap \partial \Omega = \emptyset$ .
- The corresponding matrix is neither Hermitian, nor symmetric. This cannot be viewed as a deficiency of the approach since the size of the system is *small* and thus can be solved easily using LU factorization. More specifically, the size of the system is  $n^K \times n^K$ , where  $n^K$  (the number of shape functions at the element level) does not exceed few hundreds.
- For an element  $K \in \tau_h$ , the number of rhs is  $n^{\lambda^K} + 1$ . We must point out that the obtained problems can be solved in parallel since they are independent from an element  $K$  to another.

## 4.2 Step 2: The optimization procedure

In this step, we set the global problem VF in finite dimension. We have:

$$\begin{cases} \text{Find } \lambda_h \in \mathcal{M}_h \text{ such that:} \\ A_h(\lambda_h, \mu_h) = F_h(\mu_h), \quad \forall \mu_h \in \mathcal{M}_h \end{cases} \quad (17)$$

where the forms  $A_h(\cdot, \cdot)$  and  $F_h(\cdot)$  are obtained from  $A(\cdot, \cdot)$  and  $F(\cdot)$  respectively by replacing  $\varphi$  with  $\varphi_h$  and  $\Phi(\mu_h)$  with  $\Phi_h(\mu_h)$ , for  $\mu_h \in \mathcal{M}_h$ . Hence, solving the variational problem (17) comes to solve the following linear algebraic system:

$$\mathbf{A}\mathbf{\Lambda} = \mathbf{b}, \quad (18)$$

where the entries of the matrix  $\mathbf{A}$  and of the vector  $\mathbf{b}$  depend on the expression of the continuous bilinear form  $A$ . More specifically, when using the bilinear form given by (10), the entries of the matrix  $\mathbf{A}$  and of the vector  $\mathbf{b}$  are given by:

$$\begin{aligned} \mathbf{A}_{lm} = & \sum_{e \text{ - interior edge}} \beta_e \int_e [\Phi_h(\mu_m)] [\overline{\Phi_h(\mu_l)}] ds \\ & + \sum_{\substack{e \text{ - interior edge} \\ e = \partial K \cap \partial K'}} \gamma_e \int_e \left( \mu_m^K + \mu_m^{K'} + i\alpha \left( \Phi_h(\mu_m^K) + \Phi_h(\mu_m^{K'}) \right) \right) \\ & \left( \overline{\mu_l^K + \mu_l^{K'} + i\alpha \left( \Phi_h(\mu_l^K) + \Phi_h(\mu_l^{K'}) \right)} \right) ds \end{aligned} \quad (19)$$

and

$$\begin{aligned} \mathbf{b}_l = & - \sum_{e \text{ - interior edge}} \beta_e \int_e [\varphi_h] [\overline{\Phi_h(\mu_l)}] ds \\ & - i\alpha \sum_{\substack{e \text{ - interior edge} \\ e = \partial K \cap \partial K'}} \gamma_e \int_e \left( \varphi_h^K + \varphi_h^{K'} \right) \\ & \left( \overline{\mu_l^K + \mu_l^{K'} + i\alpha \left( \Phi_h(\mu_l^K) + \Phi_h(\mu_l^{K'}) \right)} \right) ds \end{aligned} \quad (20)$$

for  $1 \leq l, m \leq n^\lambda$ .

On the other hand, when the bilinear form  $A$  is given by (12), the entries of the corresponding matrix, denoted by  $\tilde{\mathbf{A}}$ , are defined by:

$$\begin{aligned} \tilde{\mathbf{A}}_{lm} = & \sum_{e \text{ - interior edge}} \beta_e \int_e [\Phi_h(\mu_m)] [\overline{\Phi_h(\mu_l)}] ds \\ & + \sum_{e \text{ - interior edge}} \gamma_e \int_e [[\partial_n \Phi_h(\mu_m)]] [\overline{[\partial_n \Phi_h(\mu_l)]}] ds \\ & + \sum_{e \subset \partial\Omega} \omega_e \int_e (\partial_n \Phi_h(\mu_m) - i k \Phi_h(\mu_m)) (\overline{\partial_n \Phi_h(\mu_l) - i k \Phi_h(\mu_l)}) ds, \end{aligned} \quad (21)$$

and consequently, the right-hand side, denoted by  $\tilde{\mathbf{b}}$ , is given by:

$$\begin{aligned} \tilde{\mathbf{b}}_l = & - \sum_{e \text{ - interior edge}} \beta_e \int_e [\varphi_h] [\overline{\Phi_h(\mu_l)}] ds \\ & - \sum_{e \text{ - interior edge}} \gamma_e \int_e [[\partial_n \varphi_h]] [\overline{[\partial_n \Phi_h(\mu_l)]}] ds \\ & - \sum_{e \subset \partial\Omega} \omega_e \int_e (\partial_n \varphi_h - i k \varphi_h - g) (\overline{\partial_n \Phi_h(\mu_l) - i k \Phi_h(\mu_l)}) ds \end{aligned} \quad (22)$$

for  $1 \leq l, m \leq n^\lambda$ .

The unknown  $\mathbf{\Lambda}$  is a vector in  $\mathbb{C}^{n^\lambda}$  whose components are the values of  $\lambda_h$  in the basis of  $\mathcal{M}_h$ .

Hence, from a numerical point of view, two approaches are possible at Step 2 for determining the Lagrange multiplier:

- Approach 1: solving the linear system (19)-(20).
- Approach 2: solving the linear system (21)-(22).

Note that the matrices  $\mathbf{A}$  and  $\tilde{\mathbf{A}}$  are both *Hermitian*. Next, we show that the system in Approach 1 is positive semi-definite, whereas the one in Approach 2 is positive definite. Therefore, Approach 2 is expected to be more stable.

**Proposition 3.**

(i) The matrix  $\mathbf{A}$  is positive semi-definite. Moreover, for  $\Phi_h(\mu_h^K)$  satisfying, in the *strong* sense, the following equations:

$$\forall K \in \tau_h, \forall \mu_h \in \mathcal{M}_h, \quad \partial_n \Phi_h(\mu_h^K) = i\alpha \Phi_h(\mu_h^K) + \mu_h^K \text{ on } \partial K \cap \tilde{\Omega}, \quad (23)$$

and

$$\forall K \in \tau_h, \forall \mu_h \in \mathcal{M}_h, \quad \partial_n \Phi_h(\mu_h^K) = i k \Phi_h(\mu_h^K) \text{ on } \partial K \cap \partial\Omega, \quad (24)$$

$\mathbf{A}$  is a positive-definite matrix.

(ii) The matrix  $\tilde{\mathbf{A}}$  is positive definite.

**Proof of Proposition 3.** Let  $\{\mu_j\}_{1 \leq j \leq n^\lambda}$  be a basis of  $\mathcal{M}_h$  and let  $\mathbf{y} = {}^t[y_1, y_2, \dots, y_{n^\lambda}] \in \mathbb{C}^{n^\lambda}$  be an ordinary vector. We set:

$$\eta_h = \sum_{1 \leq l \leq n^\lambda} y_l \mu_l. \quad (25)$$

Therefore, it is easy to verify that:

$$\Phi_h(\eta_h) = \sum_{1 \leq l \leq n^\lambda} y_l \Phi_h(\mu_l). \quad (26)$$

Consequently, in each  $K \in \tau_h$ ,  $\Phi_h(\eta_h)$  satisfies:

$$a_K(\Phi_h(\eta_h^K), v_h^K) = \int_{\partial K \cap \tilde{\Omega}} \eta_h^K \overline{v_h^K} ds \quad \forall v_h^K \in \mathcal{V}_h(K). \quad (27)$$

(i) Using the definition of  $\eta_h$ , we have:

$$\mathbf{y}^* \mathbf{A} \mathbf{y} = \sum_{e \text{--interior edge}} \left( \beta_e \|\Phi_h(\eta_h)\|_{L^2(e)}^2 + \gamma_e \|[\eta_h + i\alpha \Phi_h(\eta_h)]\|_{L^2(e)}^2 \right). \quad (28)$$

Hence,

$$\mathbf{y}^* \mathbf{A} \mathbf{y} \geq 0, \quad \forall \mathbf{y} \in \mathbb{C}^{n^\lambda}$$

that is  $\mathbf{A}$  is a positive semi-definite matrix.

Next, we assume that condition (23) is satisfied. Let  $\mathbf{y}$  be a vector in  $\mathbb{C}^{n^\lambda}$  such that  $\mathbf{y}^* \mathbf{A} \mathbf{y} = 0$ . Then, it follows from Eq. (28) that:

$$\sum_{e \text{--interior edge}} \left( \beta_e \|\Phi_h(\eta_h)\|_{L^2(e)}^2 + \gamma_e \|[\eta_h + i\alpha \Phi_h(\eta_h)]\|_{L^2(e)}^2 \right) = 0. \quad (29)$$

Since for any interior edge  $e$ ,  $\beta_e > 0$  and  $\gamma_e > 0$ , we have:

$$\| [\Phi_h(\eta_h)] \|_{L^2(e)} = 0 \text{ and } \| [[\eta_h + i\alpha\Phi_h(\eta_h)]] \|_{L^2(e)} = 0 \text{ on all interior edges.} \quad (30)$$

Therefore,

$$[\Phi_h(\eta_h)] = 0 \text{ on all interior edges}$$

and, using Eq. (23), we have also:

$$[[\partial_n \Phi_h(\eta_h)]] = 0 \text{ on all interior edges.}$$

Consequently,  $\Phi_h(\eta_h) \in H^1(\Omega)$  and using (24), we have:

$$\begin{cases} -\Delta \Phi_h(\eta_h) - k^2 \Phi_h(\eta_h) &= 0 & \text{in } \Omega \\ \partial_n \Phi_h(\eta_h) &= i k \Phi_h(\eta_h) & \text{on } \partial\Omega. \end{cases} \quad (31)$$

Hence,  $\Phi_h(\eta_h) = 0$  because the boundary value problem (31) admits a unique solution. It follows from Eq. (23) that  $\eta_h^K = 0, \forall K \in \tau_h$ . Consequently,  $\eta_h = 0$  and therefore,  $y_l = 0$  for all  $1 \leq l \leq n^\lambda$ . Thus,  $\mathbf{y}^* \mathbf{A} \mathbf{y} > 0, \forall \mathbf{y} \in \mathbb{C}^{n^\lambda} \setminus \{0\}$ , that is  $\mathbf{A}$  is a positive definite matrix.

(ii) Observe that the matrix  $\tilde{\mathbf{A}}$  satisfies:

$$\begin{aligned} \mathbf{y}^* \tilde{\mathbf{A}} \mathbf{y} &= \sum_{e \text{--interior edge}} \left( \beta_e \| [\Phi_h(\eta_h)] \|_{L^2(e)}^2 + \gamma_e \| [[\partial_n \Phi_h(\eta_h)]] \|_{L^2(e)}^2 \right) \\ &+ \sum_{e \subset \partial\Omega} \omega_e \| \partial_n \Phi_h(\eta_h) - i k \Phi_h(\eta_h) \|_{L^2(e)}^2 \geq 0, \end{aligned}$$

that is  $\tilde{\mathbf{A}}$  is a positive semi-definite matrix.

Let  $\mathbf{y}$  be a vector in  $\mathbb{C}^{n^\lambda}$  such that  $\mathbf{y}^* \tilde{\mathbf{A}} \mathbf{y} = 0$ . Following the same reasoning as for (i), we have:

- $[\Phi_h(\eta_h)] = 0$  on all interior edges  $e$ . Hence,  $\Phi_h(\eta_h) \in H^1(\Omega)$ .
- $[[\partial_n \Phi_h(\eta_h)]] = 0$  on all interior edges. Using the fact that  $\Delta \Phi_h(\eta_h^K) \in L^2(K)$ , we deduce that  $\Delta \Phi_h(\eta_h) \in L^2(\Omega)$ .
- $\partial_n \Phi_h(\eta_h) = i k \Phi_h(\eta_h)$  on all the boundary edges of the domain.

Using the same argument as for (i), we conclude that  $y_l = 0$  for all  $1 \leq l \leq n^\lambda$ , and then  $\tilde{\mathbf{A}}$  is a positive definite matrix.  $\square$

## 5 Numerical investigation

In order to illustrate and assess the potential of mDGM for solving efficiently Helmholtz problems, we have performed numerical experiments using discrete spaces in which the shape functions are plane waves, as done in DGM [5, 6, 7]. More specifically,  $\mathcal{V}_h$  are the spaces introduced in [5]. Once the local space of shape functions,  $\mathcal{V}_h(K)$  is chosen, the Lagrange multiplier is approximated on



each edge using a subset or all set of shape functions that occur when evaluating  $\partial_n v_h^K - i\alpha v_h^K$ , for  $v_h^K \in \mathcal{V}_h(K)$ .

From now on, we suppose that  $\Omega$  is an  $a \times a$  square domain. We use a uniform partition of  $\Omega$  in rectangular-shaped elements  $K$ . The functions  $f$  and  $g$  are such that the exact solution  $u$  of BVP is a plane wave propagating in a direction  $\mathbf{d} = (\cos \theta, \sin \theta)$ . We vary the propagation angle  $\theta$  in the interval  $[0, 2\pi)$ . In order to compare the results obtained with mDGM to those delivered by DGM, we measure, for each propagating angle  $\theta$ , the relative error using the following modified  $H^1$  norm [5]:

$$\|v\|_{\widehat{H}^1} = \left( \sum_K \|v\|_{H^1(K)}^2 + \sum_{e-\text{interior edge}} \|[v]\|_{L^2(e)}^2 \right)^{\frac{1}{2}}, \quad \forall v \in \mathcal{V}. \quad (32)$$

Note that (32) is a modified  $H^1$  norm since it takes into account the  $H^1$  norm at the element level and the jump of the numerical solution along the interior interfaces of the mesh. We also use the *total* relative error, that is the *mean* value of the relative error obtained when  $\theta \in [0, 2\pi)$ .

For all numerical experiments, we have set  $\alpha = k$ , and  $\beta_e = 1$  and  $\gamma_e = h$  for all the interior edges  $e$ . The choice of these parameters results from our numerical investigation, given the lack of theoretical guidelines.

We present the results of two classes of numerical experiments: experiments using four plane waves per element, and experiments using eight plane waves per element. All the results are compared to the ones obtained with DGM.

In all numerical experiments, the linear systems (19)-(20) and (21)-(22) are solved using an LU decomposition method for sparse matrices designed by Pardiso [14, 18, 19].

### 5.1 Four plane waves per element

We equip each rectangular element with four plane waves positioned as indicated in Figure 1.

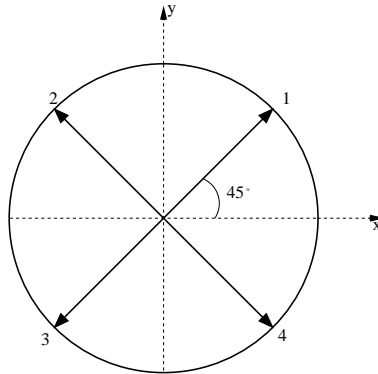


Figure 1: Four plane waves positioned at  $\theta_p = \pi/4 + (p-1)\pi/2$ ,  $\forall 1 \leq p \leq 4$

More specifically, for each  $K \in \tau_h$ , we consider:

$$\mathcal{V}_h(K) = \left\{ v_h^K = \sum_{1 \leq p \leq 4} e^{i k \theta_p \cdot \mathbf{x}} u_p, \theta_p = {}^t [\cos \theta_p, \sin \theta_p] \right. \\ \left. \theta_p = \pi/4 + (p-1)\pi/2, 1 \leq p \leq 4, u_p \in \mathbb{C} \right\}.$$

As stated earlier, the choice of the basis of  $\mathcal{M}_h$  is related to the computation

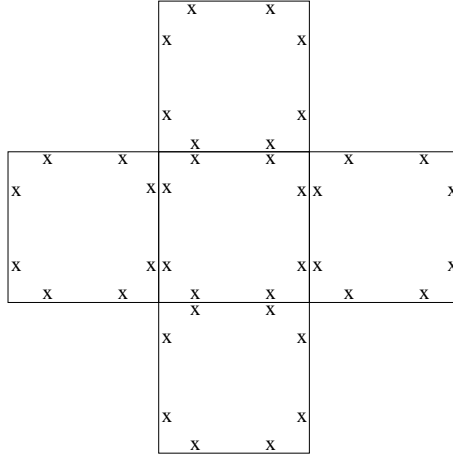


Figure 2: A 40 dofs stencil corresponding to 2 dofs per edge

of  $\partial_n v_h^K - i k v_h^K$  on the edges of the mesh, for  $v_h^K \in \mathcal{V}_h(K)$ . Observe that,  $\forall v_h^K \in \mathcal{V}_h(K)$ , on each interior edge  $e$  we have:

$$\partial_n v_h - i k v_h = \mu_1 e^{i k \frac{\sqrt{2}}{2} s} + \mu_2 e^{-i k \frac{\sqrt{2}}{2} s}, \quad (33)$$

where  $s$  represents the curvilinear abscissa and  $\mu_1, \mu_2 \in \mathbb{C}$ . Consequently, the Lagrange multiplier is approximated in the following discrete space:

$$\mathcal{M}_h = \left\{ \mu_h \in \mathcal{M}; \forall K \in \tau_h, \mu_h^K|_e = \mu_1^K e^{i k \frac{\sqrt{2}}{2} x} + \mu_2^K e^{-i k \frac{\sqrt{2}}{2} x} \text{ if } e \parallel \vec{x}, \right. \\ \left. \mu_h^K|_e = \mu_1^K e^{i k \frac{\sqrt{2}}{2} y} + \mu_2^K e^{-i k \frac{\sqrt{2}}{2} y} \text{ if } e \parallel \vec{y}, \mu_1, \mu_2 \in \mathbb{C} \right\}.$$

The spaces  $\mathcal{V}_h$  and  $\mathcal{M}_h$  defined above correspond to the so-called R-4-2 element in the nomenclature of DGM (see [5]). Note that for the DGM, considering two dofs per edge leads also to a complete approximation of the Lagrange multiplier. We must point out that, unlike the DGM, the Lagrange multipliers in mDGM are not continuous across interior boundaries. This is why on an edge shared by two elements, the Lagrange multipliers on each side of the edge are different. Consequently, the stencil of the matrix given by (19) is equal to 40 (see Figure 2). Note that in the DGM the Lagrange multipliers are equal on both sides of the edge, which leads to a stencil equal to 14.

The first experiments consist in comparing the error delivered by both numerical methods (DGM and mDGM) for different values of  $ka$ , while maintaining  $kh$  constant. More specifically, we consider  $ka = 10, 20, 30$  and we choose the step size of the mesh discretization  $h/a$  such that  $kh = \frac{1}{5}$ , which is about 30 elements per wavelength. The results are depicted in Figure 3 and Figure 4. These results indicate the following:

- The two methods deliver results with the same level of accuracy, as indicated in Figure 3: both curves are superposed.
- As expected, the relative error is  $\pi/2$  periodic (see Figure 3). On each period  $[(l-1)\pi/2, l\pi/2]$  (with  $l = 1, 2, 3, 4$ ), the error is symmetric with respect to the propagation angle  $(2l-1)\pi/4$ . Moreover, the error is minimal for  $\pi/4, 3\pi/4, 5\pi/4, 7\pi/4$  (less than 1%) and maximal for  $0, \pi/2, \pi, 3\pi/2$  (about 4%). This is due to the chosen basis, which includes the exact solution when the propagation angle is  $(2l-1)\pi/4$ , with  $l = 1, 2, 3, 4$  and to the fact that the Lagrange multiplier field contains all the functions needed to have a complete approximation.
- Figure 4 indicates that the R-4-2 element (for both methods) exhibits little pollution: increasing  $ka$ , while maintaining  $kh$  constant, leads to an increase in the relative error which is less than 0.5% at most (see Figure 4 at angles  $\theta = l\pi/2$ , with  $l = 0, 1, \dots, 8$ ).

Next, we compare the sensitivity of the total relative error (the mean value over the propagation angles) to the mesh size. The result depicted in Figure 5 is obtained for  $ka = 1$ . One can observe the following:

- For  $h/a > \frac{1}{100}$ , the errors delivered by the two methods are comparable. The two curves are on top of each other.
- For  $h/a < \frac{1}{100}$  mDGM outperforms DGM. As we refine the mesh ( $h/a < \frac{1}{100}$ ), DGM becomes unstable. Indeed, there is a dramatic loss in the accuracy of more than one order of magnitude. The error jumps from 0.09% (for  $h/a = \frac{1}{100}$ ) to 1.5% (for  $h/a = \frac{1}{190}$ ). The instability observed in DGM seems to be related to the severe ill conditioning of the local matrices. Observe that mDGM remains stable as we refine the mesh. The last point of the curve was obtained for  $h/a = \frac{1}{450}$ , the limit of our computing platform. The total relative error for this mesh size is 0.04%.

We must point out that the performance of mDGM in this case is not sensitive to the choice of the approach for solving the linear system in Step 2. Both approaches deliver results with the same level of accuracy.

## 5.2 Eight plane waves per element

We approximate the primal variable using eight plane waves, positioned at:

$$\theta_p = (p-1)\pi/4, \quad \forall 1 \leq p \leq 8.$$

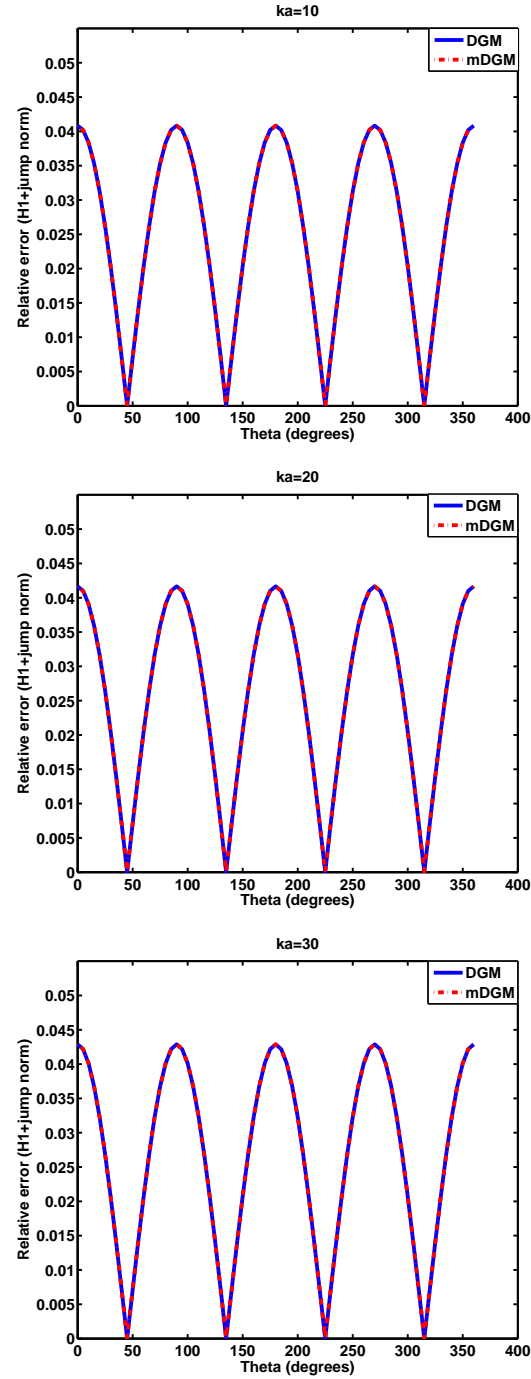


Figure 3: Performance of the two methods for  $kh=1/5$

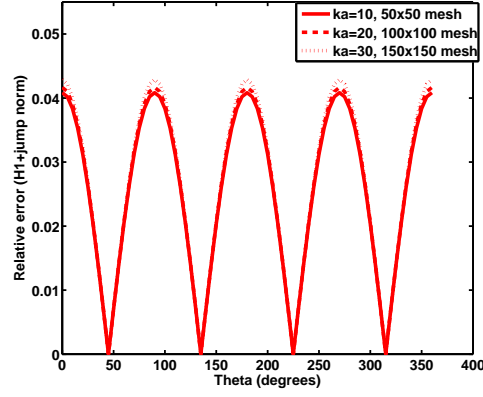


Figure 4: Pollution effect for the R-4-2 element

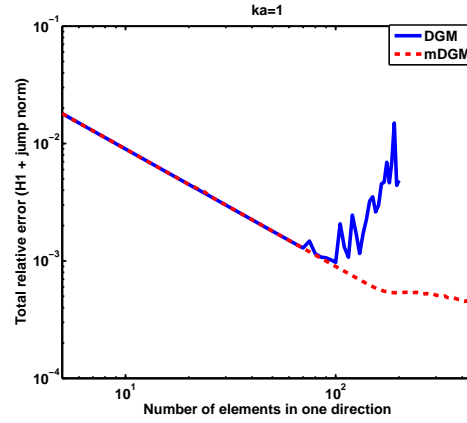


Figure 5: Convergence rates R-4-2 element

This choice corresponds to the following discrete space for the primal variable:

$$\mathcal{V}_h(K) = \left\{ v_h|_K = \sum_{1 \leq p \leq 8} e^{i k \theta_p \cdot x} u_p, \theta_p = {}^t [\cos \theta_p, \sin \theta_p], \right. \\ \left. \theta_p = (p-1)\pi/4, 1 \leq p \leq 8, u_p \in \mathbb{C} \right\}.$$

For an element  $v_h^K \in \mathcal{V}_h(K)$ , the full approximation of  $\partial_n v_h^K - i k v_h^K$  leads to five dofs per edge. More specifically in mDGM, the discrete space  $\mathcal{M}_h$  corresponding to the full approximation of the Lagrange multiplier is given by:

$$\begin{aligned} \mathcal{M}_h = & \left\{ \mu_h \in \mathcal{M}; \forall K \in \tau_h, \mu_h^K|_e = \mu_1^K + \mu_2^K e^{i k x} + \mu_3^K e^{-i k x} + \mu_4^K e^{i k \frac{\sqrt{2}}{2} x} \right. \\ & + \mu_5^K e^{-i k \frac{\sqrt{2}}{2} x} \text{ if } e \parallel \overrightarrow{x}, \mu_h^K|_e = \mu_1^K + \mu_2^K e^{i k y} + \mu_3^K e^{-i k y} \\ & \left. + \mu_4^K e^{i k \frac{\sqrt{2}}{2} y} + \mu_5^K e^{-i k \frac{\sqrt{2}}{2} y} \text{ if } e \parallel \overrightarrow{y}, \mu_1, \mu_2, \mu_3, \mu_4, \mu_5 \in \mathbb{C} \right\}. \end{aligned}$$

Following the nomenclature introduced in [5], such an approximation is called the *R-8-5* element. Note that mDGM can be implemented using less dofs per edge for the Lagrange multiplier. We recall that in DGM, the maximum number of dofs considered on each edge is three. Indeed, the computation of the normal derivative of the numerical solution leads to the following complete approximation:

$$\lambda_h = \mu_1 + \mu_2 e^{i k \frac{\sqrt{2}}{2} s} + \mu_3 e^{-i k \frac{\sqrt{2}}{2} s},$$

where  $s$  represents the curvilinear abscissa. This choice of approximation corresponds to the so-called *R-8-3* element.

We first present the results obtained in the case of Approach 1, than the ones obtained in the case of Approach 2.

### 5.2.1 Performance assessment in the case of Approach 1

Since the full approximation in DGM requires three dofs per edge, we first compare the performance of mDGM and DGM when using the *R-8-3* element. The result depicted in Figure 6 compares the relative error delivered by both methods, as a function of the propagation angle. This result is obtained for  $ka = 10$  and  $h/a = 1/20$ , that is  $kh = \frac{1}{2}$ , corresponding to about 12 elements per wavelength. It shows a clear superiority of mDGM over DGM. In addition, we have:

- DGM delivers the exact solution for  $l\pi/4$ , with  $l = 0, 1, \dots, 8$ . Note that in each of these cases, the exact solution is represented by one of the basis functions of the considered element and all the functions obtained when computing the normal derivative are in the Lagrange multiplier field. On the other hand, mDGM computes exactly the solution at angles  $\pi/4, 3\pi/4, 5\pi/4, 7\pi/4$  only, as it will be shown later (see Figure 8). This is not surprising since two dofs are removed from the full approximation of the Lagrange multiplier and therefore, for the plane waves propagating in the directions parallel to the axis, the approximation is not complete.
- The total relative error is about 0.091% for mDGM and is about 5% for DGM. This means that mDGM improves the accuracy by about one and a half order of magnitude.
- Observe that DGM *R-8-3* is an unstable element. Indeed, the error obtained for  $\theta = (2l - 1)\pi/8$  should be the same for all  $l = 1, 2, \dots, 8$  and symmetric with respect to  $0, \pi/4, 2\pi/4, \dots, 8\pi/4$ . Figure 6 shows that these values are not equal.

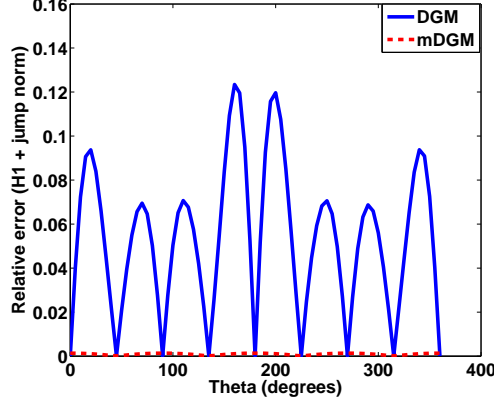


Figure 6: Performance of DGM and mDGM equipped with the  $R$ -8-3 element for  $ka=10$ ,  $h/a=1/20$

The instability observed in the DGM approach for  $R$ -8-3 element is not surprising since this element does not satisfy the numerical inf-sup condition required by DGM [11]. Almost everywhere in the mesh, for an element  $K$  there are twelve dofs for the Lagrange multiplier and only eight dofs for the primal variable. A dof must be removed from each edge. For this reason, two discrete spaces were suggested in [5] for the discrete dual variable, leading respectively to the so-called  $R$ -8-2a and  $R$ -8-2b elements. Since in the cited paper, the  $R$ -8-2b element was shown to deliver more accurate results than the  $R$ -8-2a element, we have compared the two methods when employing the  $R$ -8-2b element. We

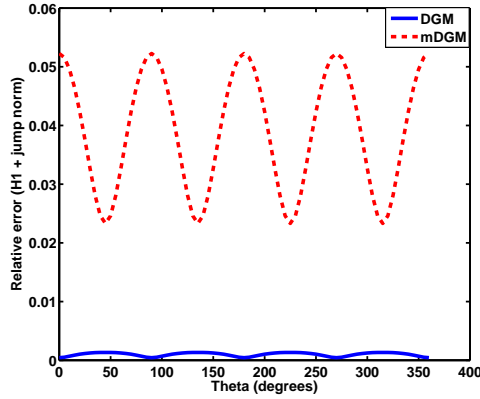


Figure 7: Performance of DGM and mDGM equipped with the  $R$ -8-2b element for  $ka=10$ ,  $h/a=1/20$

recall that the  $R$ -8-2b element corresponds to the following approximation of the Lagrange multiplier:

$$\lambda_h = \mu_1 e^{ik \frac{\sqrt{2}}{4} s} + \mu_2 e^{-ik \frac{\sqrt{2}}{4} s},$$

where  $s$  is the curvilinear abscissa. Similarly to the previous numerical experiment, we have set  $ka = 10$  and  $h/a = 1/20$ , which corresponds to  $kh = \frac{1}{2}$ . The result depicted in Figure 7 suggests the following:

- As expected, both methods preserve the symmetry of the error with respect to the propagation angles  $\theta = \pi/4, 3\pi/4, 5\pi/4, 7\pi/4$ .
- The total relative error obtained with the mDGM is about 5% while the one obtained with DGM is about 0.099%. This superiority of DGM over mDGM is most likely due to the poor approximation of the Lagrange multiplier in the mDGM (three out of five dofs are neglected), compared to the DGM, where only one dof out of three is neglected.

Next, we enrich the approximation of the Lagrange multiplier in mDGM. We use the elements  $R-8-3$  and  $R-8-5$  and compare mDGM to DGM, equipped with the best element  $R-8-2b$ . The results reported in Figure 8 are obtained for  $ka = 10$  and  $h/a = 1/20$ . One can observe the following:

- There is a little improvement in the accuracy of the results delivered by  $R-8-3$  and  $R-8-5$  mDGM elements over DGM  $R-8-2b$  element. The total errors obtained with these elements are 0.091% and 0.048% respectively, whereas the one delivered by DGM equipped with the best element,  $R-8-2b$ , is 0.099%.

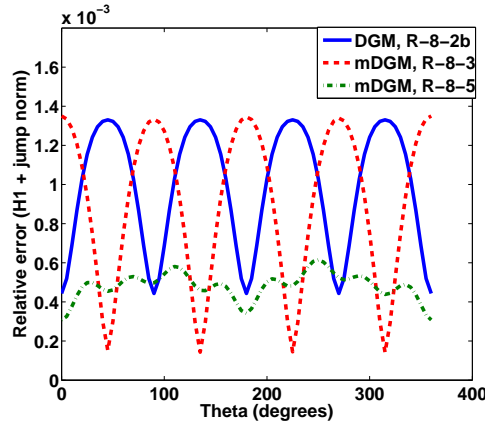


Figure 8: Performance of DGM equipped with the  $R-8-2b$  and of mDGM equipped with the  $R-8-3$  and  $R-8-5$  for  $ka=10$ ,  $h/a=1/20$

- The  $R-8-5$  element seems to be unstable, as illustrated by the loss of the symmetry at the propagation angles  $\theta = l\pi/4$ , with  $l = 0, 1, 2, \dots, 7$ , corresponding to the basis functions. This numerical instability is also noticeable when we compare  $R-8-5$  to  $R-8-3$ : for some propagation angles,  $R-8-3$  is more accurate than  $R-8-5$ , which is contrary to what is expected since using five dofs per edge leads to the full approximation of the Lagrange multiplier. We believe that the instability of the  $R-8-5$  element is due to the loss of the linear independence of the shape functions, but also to the fact that the matrix of the system is positive semi-definite.



Next, we investigate the sensitivity of the total relative error with respect to  $h$ , the step size of the mesh discretization. We consider the case of  $R$ -8-2b element and we set  $ka = 1$ . We evaluate the total relative error, as well as the smallest eigenvalue of the local system. The results are reported in Table I.

Table 1: Performance of the R-8-2b element for DGM and mDGM when  $ka=1$

| $h/a$ | DGM                  |  | mDGM                 |  |
|-------|----------------------|--|----------------------|--|
|       | Total relative error | The smallest eigenvalue                    | Total relative error | The smallest eigenvalue                    |
| 1/4   | 0.016%               | $5.7 \cdot 10^{-10} + 2.9 \cdot 10^{-18}i$ | 0.054%               | $5.7 \cdot 10^{-10} - 3.1 \cdot 10^{-11}i$ |
| 1/5   | 0.011%               | $9.5 \cdot 10^{-11} + 8.0 \cdot 10^{-19}i$ | 0.036%               | $9.5 \cdot 10^{-11} - 4.1 \cdot 10^{-12}i$ |
| 1/6   | 0.019%               | $2.2 \cdot 10^{-11} - 2.4 \cdot 10^{-18}i$ | 0.025%               | $2.2 \cdot 10^{-11} - 7.9 \cdot 10^{-13}i$ |
| 1/7   | 0.092%               | $6.5 \cdot 10^{-12} + 5.9 \cdot 10^{-19}i$ | 0.019%               | $6.4 \cdot 10^{-12} - 2.0 \cdot 10^{-13}i$ |
| 1/8   | 0.394%               | $2.2 \cdot 10^{-12} + 1.8 \cdot 10^{-18}i$ | 0.015%               | $2.2 \cdot 10^{-12} - 5.9 \cdot 10^{-14}i$ |
| 1/9   | 1.206%               | $8.6 \cdot 10^{-13} - 1.6 \cdot 10^{-18}i$ | 0.017%               | $8.6 \cdot 10^{-13} - 2.1 \cdot 10^{-14}i$ |
| 1/20  | 4.734%               | $1.5 \cdot 10^{-15} - 6.0 \cdot 10^{-21}i$ | 3.590%               | $1.5 \cdot 10^{-15} + 5.9 \cdot 10^{-17}i$ |
| 1/21  | 11.664%              | $9.8 \cdot 10^{-16} + 7.0 \cdot 10^{-18}i$ | 5.889%               | $9.6 \cdot 10^{-16} + 2.4 \cdot 10^{-17}i$ |
| 1/23  | 33.435%              | $4.7 \cdot 10^{-16} - 3.9 \cdot 10^{-18}i$ | 37.378%              | $4.7 \cdot 10^{-16} - 6.7 \cdot 10^{-17}i$ |

- The results reveal that in the DGM approach, the error decreases as long as  $h/a > \frac{1}{6}$ . Then, the error jumps from 0.01% to about 33% for  $kh = \frac{1}{23}$ . This is not surprising. Indeed, the local systems in the DGM are nearly singular and therefore extremely ill-conditioned when  $h$  becomes small, as indicated by the values corresponding to the smallest eigenvalues.
- The smallest error delivered by mDGM, which is about 0.01%, is obtained for  $kh = \frac{1}{8}$ . Then the error jumps to 37% for  $kh = \frac{1}{23}$ . This instability is, a priori, unexpected and very surprising since we have introduced the Robin-type boundary condition to address specifically this issue, as demonstrated in the case of the  $R$ -4-2 element. A quick look at Table I indicates that the local system corresponding to mDGM becomes nearly singular too. Hence, contrary to our goal, the presence of  $\alpha$  seems to be not sufficient to avoid the singularity of the local systems. We believe that the singularity of the local system in the mDGM formulation is due to the loss of the linear independence of the shape functions (eight plane waves) as  $h$  becomes small, as well as to the non-negativeness nature of the matrix  $\mathbf{A}$  of Approach 1.

The loss of the linear independence can be demonstrated as follows: let  $K$  be an element of the mesh. For simplicity, we assume  $K$  to be the square  $[0, h] \times [0, h]$ . Then, for a function  $w_h^K \in \mathcal{V}_h(K)$ , there exist  $c_1, c_2, \dots, c_8 \in \mathbb{C}$  such that:

$$w_h^K = \sum_{m=1}^8 c_m e^{ik\theta_m \cdot \mathbf{x}}.$$

Assume that:

$$a_K(w_h^K, w_h^K) = 0. \quad (34)$$

Consequently,  $w_h^K = 0$  on  $\partial K$ , which means:

$$\sum_{m=1}^8 c_m e^{i k \theta_m \cdot x} = 0 \quad \text{on } \partial K.$$

We write this equality on each of the four edges of  $K$ . For all  $x \in [0, h]$  and  $y \in [0, h]$ , we have:

$$\left\{ \begin{array}{l} c_1 + c_5 + (c_2 + c_4) e^{i k \frac{\sqrt{2}}{2} y} + c_3 e^{i k y} + (c_6 + c_8) e^{-i k \frac{\sqrt{2}}{2} y} + c_7 e^{-i k y} = 0, \\ c_3 + c_7 + (c_2 + c_8) e^{i k \frac{\sqrt{2}}{2} x} + c_1 e^{i k x} + (c_4 + c_6) e^{-i k \frac{\sqrt{2}}{2} x} + c_5 e^{-i k x} = 0, \\ c_1 e^{i k h} + c_5 e^{-i k h} + (c_2 e^{i k \frac{\sqrt{2}}{2} h} + c_4 e^{-i k \frac{\sqrt{2}}{2} h}) e^{i k \frac{\sqrt{2}}{2} y} + c_3 e^{i k y} + \\ \quad (c_6 e^{-i k \frac{\sqrt{2}}{2} h} + c_8 e^{i k \frac{\sqrt{2}}{2} h}) e^{-i k \frac{\sqrt{2}}{2} y} + c_7 e^{-i k y} = 0, \\ c_3 e^{i k h} + c_7 e^{-i k h} + (c_2 e^{i k \frac{\sqrt{2}}{2} h} + c_8 e^{-i k \frac{\sqrt{2}}{2} h}) e^{i k \frac{\sqrt{2}}{2} x} + c_1 e^{i k x} + \\ \quad (c_6 e^{-i k \frac{\sqrt{2}}{2} h} + c_4 e^{i k \frac{\sqrt{2}}{2} h}) e^{-i k \frac{\sqrt{2}}{2} x} + c_5 e^{-i k x} = 0. \end{array} \right.$$

Therefore, we deduce that:

$$\left\{ \begin{array}{l} c_1 = c_3 = c_5 = c_7 = 0 \\ c_4 = -c_2 \\ c_6 = c_2 \\ c_8 = -c_2 \\ c_2 \left( e^{i k \frac{\sqrt{2}}{2} h} - e^{-i k \frac{\sqrt{2}}{2} h} \right) = 0. \end{array} \right. \quad (35)$$

The problem here is that when  $h \rightarrow 0$ , we have:  $e^{i k \frac{\sqrt{2}}{2} h} \rightarrow 1$  and  $e^{-i k \frac{\sqrt{2}}{2} h} \rightarrow 1$  and hence, numerically speaking, it is not necessary to have  $c_2 = 0$  in order to have  $w_h^K = 0$  on  $\partial K$ . Consequently,  $w_h^K$  may not be equal to 0 when (34) is satisfied. This computation shows that when  $h$  tends to 0, the eight plane waves become linearly dependent, which leads to the singularity in the local matrix.

**Remark 4.** We have performed additional numerical experiments for higher frequencies and observed that both methods become unstable as we refine the mesh. DGM is stable as long as  $kh > \frac{1}{6}$ , while mDGM seems to remain stable longer ( $kh > \frac{1}{9}$ ). We must point out that the source for the instability is however different. DGM is unstable not only because of the singularity of the local systems, but also due to the loss of the linear independence of the plane waves as  $h$  becomes small. We believe that the resulting local system in DGM is also very sensitive to this loss.

On the other hand, mDGM restores the stability of the local problems and leads to a more stable formulation when the shape functions remain linearly independent, as it has been demonstrated in the case of the  $R$ -4-2 element (see Figure 5) and in the next experiments. However, the previous numerical results suggest that the loss of the linear independence affects dramatically the stability of mDGM when using Approach 1, that is when the linear system is non-negative only.

The following experiment reveals the behavior of mDGM and DGM when the

used shape functions remain linearly independent as we refine the mesh. This experiment consists in approximating the solution, at the element level, using seven plane waves, positioned at:

$$\theta_p = 2(p-1)\pi/7, \quad 1 \leq p \leq 7. \quad (36)$$

We have maintained the same two dofs per edge as in  $R$ -8-2b (see Section 5.2). Following the nomenclature introduced in [5], we will refer to this element as  $R$ -7-2. For  $ka = 1$ , we have compared the sensitivity of the total relative error to the mesh size for the mDGM  $R$ -8-2b and  $R$ -7-2 elements. The result depicted in Figure 9 illustrates the following:

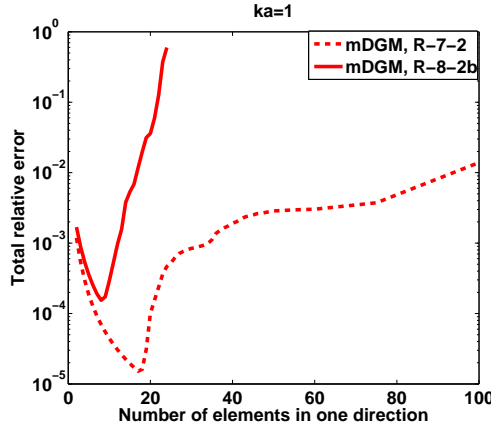


Figure 9: Convergence rates for mDGM  $R$ -8-2b and  $R$ -7-2 elements

- The error delivered by the  $R$ -7-2 element decreases as long as  $kh > \frac{1}{18}$ , unlike the  $R$ -8-2b element, which delivers the smallest error for  $h/a = \frac{1}{9}$ . This means that there is a reduction of factor 2 on the mesh size, while maintaining the stability.
- The  $R$ -7-2 element is shown to be more accurate. For each mesh size, the error delivered by this element is smaller than the one obtained with the  $R$ -8-2b element. Moreover, the most accurate approximation (about 0.001% for  $h/a = \frac{1}{18}$ ) obtained with  $R$ -7-2 is one order of magnitude lower than  $R$ -8-2b element (about 0.01% for  $h/a = \frac{1}{9}$ ).
- Last, note that the  $R$ -7-2 element remains stable while refining the mesh. For  $h/a = \frac{1}{100}$ , the total relative error is about 1%, unlike the  $R$ -8-2b element in which the error jumps from 0.01% (for  $h/a = \frac{1}{9}$ ) to 59% (for  $h/a = \frac{1}{24}$ ).

We have also compared the errors delivered by the mDGM  $R$ -7-2 element to the ones obtained with the DGM  $R$ -7-2 element. The result is reported in Figure 10. The following observations are noteworthy:

- The accuracy of the two methods is comparable for  $h/a > \frac{1}{8}$ . In this region the two curves are superposed.

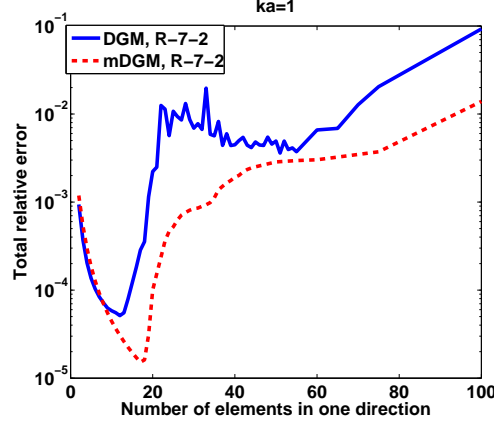


Figure 10: Convergence rates for mDGM R-7-2 and DGM R-7-2 elements

- The DGM  $R$ -7-2 element delivers the most accurate approximation (which is about 0.005%) for  $h/a = \frac{1}{12}$ . Observe that mDGM becomes unstable slightly later: the smallest error (about 0.001%) is obtained for  $h/a = \frac{1}{18}$ .
- Although for both methods we observe numerical instabilities as soon as  $h/a < \frac{1}{18}$ , mDGM is more accurate than DGM. For any mesh size the error delivered by mDGM is smaller than the one obtained with DGM. Moreover, for some mesh sizes, mDGM outperforms DGM by one order of magnitude. We believe that this is due to the local problems which are nearly singular in DGM.

The results depicted in Figure 9 seem to be surprising since one may expect that approximating the primal variable with eight plane waves leads to more accurate results than when using seven plane waves. Here, it seems that the linear independence when using seven plane waves is less sensitive (remains longer) to the mesh size  $h/a$ . Consequently, the linear system corresponding to Approach 1 is more stable in this case than when using eight plane waves.

In Table II we report the total relative error, as well as the smallest eigenvalue of the local system for the  $R$ -7-2 element for both DGM and mDGM methods, when  $ka = 1$ .

- As in the case of the  $R$ -8-2b element, the eigenvalues of the local matrices corresponding to each mesh size have the same real part in DGM and mDGM. The Robin-type condition used in mDGM leads to more important imaginary parts in mDGM.
- The values of the smallest eigenvalues and the comparison to the ones reported in Table I show that the linearly independence of the seven plane waves is less sensitive to the mesh refinement. Indeed, for example for  $h/a = \frac{1}{9}$  the real part of the smallest eigenvalue obtained when using eight plane waves at the element level is  $8.64 \cdot 10^{-13}$ . This is four orders of magnitude larger than the real part of the smallest eigenvalue of the matrix obtained with seven plane waves ( $4.00 \cdot 10^{-9}$ ).

Table 2: Performance of the R-7-2 element for DGM and mDGM when  $ka=1$

| $h/a$ | DGM                  |  | mDGM                 |  |
|-------|----------------------|--|----------------------|--|
|       | Total relative error | The smallest eigenvalue                    | Total relative error | The smallest eigenvalue                    |
| 1/4   | 0.020%               | $5.2 \cdot 10^{-07} - 3.0 \cdot 10^{-17}i$ | 0.030%               | $5.2 \cdot 10^{-07} - 3.2 \cdot 10^{-08}i$ |
| 1/5   | 0.014%               | $1.4 \cdot 10^{-07} - 4.8 \cdot 10^{-17}i$ | 0.019%               | $1.4 \cdot 10^{-07} - 6.6 \cdot 10^{-09}i$ |
| 1/6   | 0.010%               | $4.6 \cdot 10^{-08} - 8.0 \cdot 10^{-18}i$ | 0.012%               | $4.6 \cdot 10^{-08} - 1.9 \cdot 10^{-09}i$ |
| 1/7   | 0.008%               | $1.8 \cdot 10^{-08} + 2.0 \cdot 10^{-17}i$ | 0.009%               | $1.8 \cdot 10^{-08} - 6.3 \cdot 10^{-10}i$ |
| 1/8   | 0.007%               | $8.1 \cdot 10^{-09} + 3.9 \cdot 10^{-17}i$ | 0.007%               | $8.1 \cdot 10^{-09} - 2.5 \cdot 10^{-10}i$ |
| 1/9   | 0.006%               | $4.0 \cdot 10^{-09} + 2.5 \cdot 10^{-17}i$ | 0.005%               | $4.0 \cdot 10^{-09} - 1.1 \cdot 10^{-10}i$ |
| 1/12  | 0.005%               | $7.1 \cdot 10^{-10} + 6.0 \cdot 10^{-18}i$ | 0.003%               | $7.1 \cdot 10^{-10} - 1.5 \cdot 10^{-11}i$ |
| 1/18  | 0.035%               | $6.3 \cdot 10^{-11} + 7.7 \cdot 10^{-18}i$ | 0.001%               | $6.3 \cdot 10^{-11} - 8.5 \cdot 10^{-13}i$ |
| 1/20  | 0.222%               | $3.3 \cdot 10^{-11} + 3.1 \cdot 10^{-19}i$ | 0.010%               | $3.3 \cdot 10^{-11} - 4.1 \cdot 10^{-13}i$ |
| 1/25  | 1.075%               | $8.7 \cdot 10^{-12} - 2.0 \cdot 10^{-18}i$ | 0.052%               | $8.7 \cdot 10^{-12} + 2.4 \cdot 10^{-17}i$ |
| 1/34  | 0.589%               | $1.4 \cdot 10^{-12} + 4.4 \cdot 10^{-18}i$ | 0.099%               | $1.4 \cdot 10^{-12} - 9.9 \cdot 10^{-15}i$ |
| 1/50  | 0.493%               | $1.4 \cdot 10^{-13} - 1.1 \cdot 10^{-19}i$ | 0.285%               | $1.4 \cdot 10^{-13} - 6.8 \cdot 10^{-16}i$ |
| 1/75  | 2.046%               | $1.2 \cdot 10^{-14} + 4.0 \cdot 10^{-18}i$ | 0.373%               | $1.2 \cdot 10^{-14} - 5.1 \cdot 10^{-17}i$ |
| 1/100 | 9.266%               | $2.1 \cdot 10^{-15} + 4.5 \cdot 10^{-18}i$ | 1.390%               | $2.1 \cdot 10^{-15} + 2.5 \cdot 10^{-18}i$ |
| 1/150 | 243.2%               | $1.9 \cdot 10^{-16} + 1.5 \cdot 10^{-19}i$ | 67.75%               | $1.8 \cdot 10^{-16} - 8.3 \cdot 10^{-18}i$ |

- A quick comparison of the errors reported in Tables I and II (see also Figure 9 for mDGM) shows that in both formulations a better conditioning of the local matrices leads to more accurate approximations. Moreover, in both methods there is a reduction of factor 2 on the mesh size, while maintaining the stability. These two important points are related to the fact that the seven shape functions remain linearly independent as we refine the mesh.
- As it was shown in Figure 10 and reported in Table II, the numerical instabilities appear earlier in DGM than in mDGM. We believe that this is due to the local systems which become nearly singular with the mesh refining. In mDGM the observed numerical instabilities are due to the non-negativeness nature of the global matrix corresponding to Approach 1.
- Last, we must point out the fact that the seven shape functions are becoming linearly dependent with the mesh refinement, but more slowly than the eight plane waves corresponding to the R-8-2b element. This behavior of the shape functions is predictable when observing the dependence of the smallest eigenvalue of the local matrix with respect to the mesh size. Consequently, this element will be ultimately unstable.

### 5.2.2 Performance assessment in the case of Approach 2

We have performed several numerical experiments to assess the performance of

mDGM when at Step 2 the linear system is given by Approach 2 (see Eqs. (21)-(22)). This system is positive definite and therefore, it is expected that with

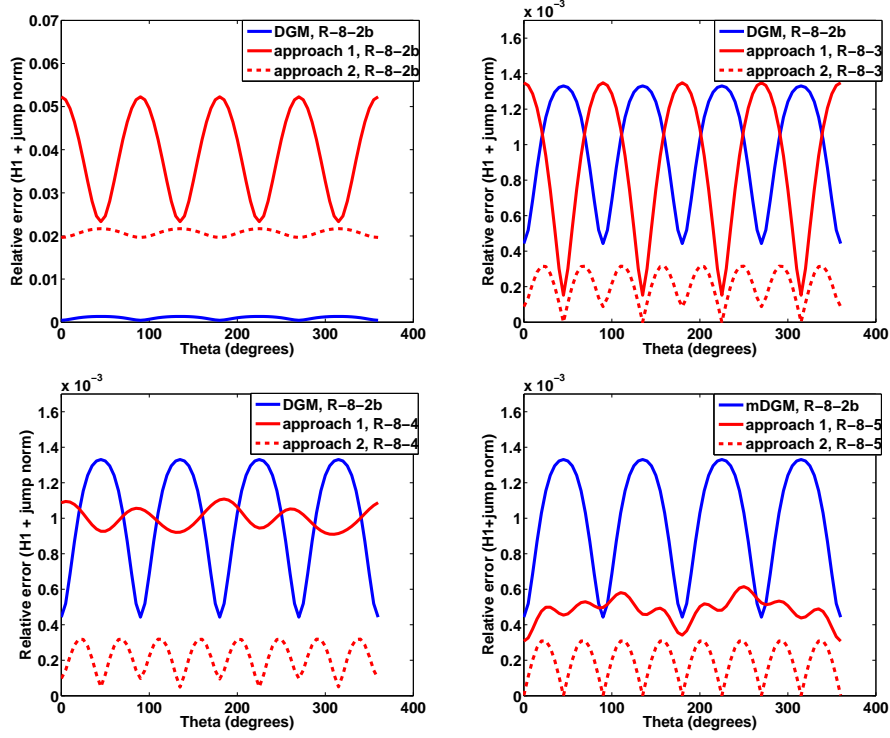


Figure 11: Performance of the three methods,  $kh=1/2$

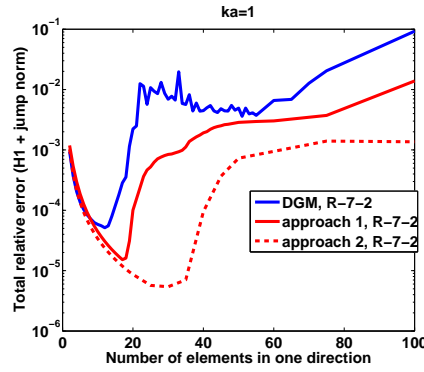


Figure 12: Convergence rates for the R-7-2 element

the well-posedness of the local boundary value problem in Step 1, the method leads to more stable and thus, more accurate numerical results. The obtained numerical results are depicted in Figures 11, 12 and 13. They clearly indicate that Approach 2 of mDGM not only outperforms both Approach 1 of mDGM and DGM, but also delivers more accurate results. The error is reduced by one

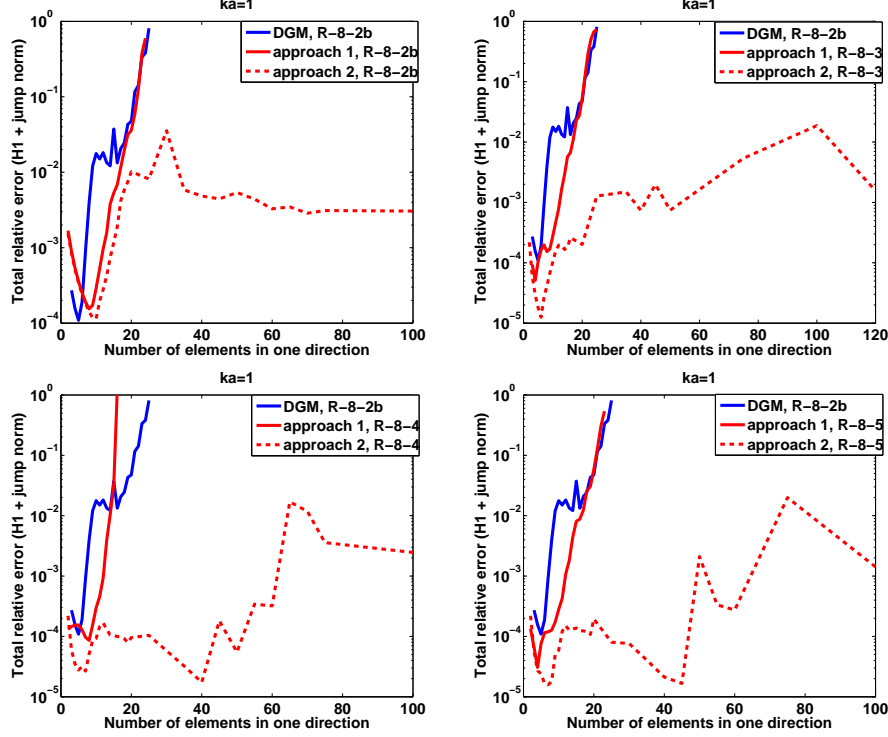


Figure 13: Convergence of the three methods - eight plane waves

to two orders of magnitude depending on the element and the mesh size.

Note that the numerical stability occurs only in the case of mDGM when equipped with Approach 2. Indeed, one can observe that the errors (see dashed curves in Figure 13) are oscillating as  $a/h$  is increasing, but their magnitude remains steadily about 1%, whereas the magnitude of the errors delivered by the two other methods increase to over 100%.

## 6 Summary and conclusions

We have designed a new solution methodology, called mDGM, for Helmholtz problems which is easy to understand and implement. At the element level, we approximate the solution by a superposition of plane waves. Consequently, the obtained solution is discontinuous and Lagrange multipliers are introduced to ensure the continuity in a weak sense. Unlike the DGM, the Lagrange multiplier is also discontinuous, which allows us to consider well-posed local problems. The algebraic approach requires solving local linear systems with multiple right-hand side: the system's size is given by the number of plane waves considered in the local basis. These problems are independent from one element to another and therefore can be solved in parallel. The global system, whose size is the number of total dofs used for approximating the Lagrange multiplier, can be either

positive semi-definite or definite depending on the approach adopted at the continuous level. The numerical results we have presented show that the proposed method is more stable than the DGM. When using Approach 2, mDGM is not only more stable than DGM, but also exhibits a better level of accuracy. More specifically, as indicated by the reported numerical results, mDGM reduces the level of errors by one to two orders of magnitude depending on the mesh size and on the element.

## Acknowledgements

The authors acknowledge the support by TOTAL and INRIA/CSUN Associate Team Magic, INRIA Bordeaux Sud-Ouest Center. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of TOTAL, INRIA or CSUN.

## References

- [1] Amara M., Djellouli R., Farhat C. Convergence analysis of a discontinuous Galerkin method with plane waves and Lagrange multipliers for the solution of Helmholtz problems *SIAM J. Numer. Anal.* 2009; **47**(2):1038–1066.
- [2] Babuška I., Melenk I.J.M. The partition of unity method *Internat. J. Numer. Methods Eng.* 1997; **40**:727–758
- [3] Babuška I., Sauter S. Is the Pollution Effect of the FEM Avoidable for the Helmholtz Equation Considering High Wave Numbers? *SIAM J. Numer. Anal.* 1997; **34**:2392–2423
- [4] Cessenat O., Despres B. Application of an ultra-weak variational formulation of elliptic PDEs to the two-dimensional Helmholtz problems *SIAM J. Numer. Anal.* 1998; **35**:255–299.
- [5] Farhat C., Harari I., Hetmaniuk U. A discontinuous Galerkin method with Lagrange multipliers for the solution of Helmholtz problems in the mid-frequency regime *Comput. Methods Appl. Mech. Eng.* 2003; **192**:1389–1419.
- [6] Farhat C., Wiedemann-Goiran P., Tezaur R. A discontinuous Galerkin method with plane waves and Lagrange multipliers for the solution of short wave exterior Helmholtz problems on unstructured meshes *Wave Motion* 2004; **39**:307–317.
- [7] Farhat C., Tezaur R., Wiedemann-Goiran P. Higher-order extensions of a discontinuous Galerkin method for mid-frequency Helmholtz problems *Internat. J. Numer. Methods Eng.* 2004; **61**:1938–1956.
- [8] Franca L.P., Farhat C., Macedo A.P., Lesoinne M. Residual-free bubbles for the Helmholtz equation *Internat. J. Numer. Methods Eng.* 1997; **40**:4003–4009.
- [9] Hadamard J. Lectures on Cauchy's Problem in Linear Partial Differential Equations *Yale University Press, New Haven* 1923;



- [10] Harari I., Hughes T.J.R. Galerkin/least-squares finite element methods for the reduced wave equation with non-reflecting boundary conditions in unbounded domains *Comput. Methods Appl. Mech. Eng.* 1992; **98**:411–454.
- [11] Harari I., Hetmaniuk U. Private communication
- [12] Hörmander L. The Analysis of Linear Partial Differential Operator *Springer-Verlag, New York* 1985;
- [13] Ihlenburg F. Finite Element Analysis of Acoustic Scattering *Appl. Math. Sci* 132, *Springer-Verlag, New York* 1998;
- [14] Karypis G., Kumar V. A fast and high quality multilevel scheme for partitioning irregular graph *SIAM Journal on Scientific Computing* 1998; **20**: 359–392.
- [15] Magoulès F. Computational Methods for Acoustics Problems *Saxe-Coburg Publications* 2008;
- [16] Monk P., Wang D.Q. A least-squares method for the Helmholtz equation *Comput. Methods Appl. Mech. Eng.* 1999; **175**:411–454.
- [17] Rose M.E. Weak element approximations to elliptic differential equations *Numer. Math.* 1975; **24**:185–204.
- [18] Schenk O., Gärtner K. Solving unsymmetric sparse systems of linear equations with PARDISO *Journal of Future Generation Computer Systems* 2004; **20**: 475–487.
- [19] Schenk O., Gärtner K. On fast factorization pivoting methods for symmetric indefinite systems *Elec. Trans. Numer. Anal.* 2006; **23**: 158–179.
- [20] Taylor M. E. Partial Differential Equations I: Basic Theory *Springer-Verlag, New York* 1997;

## Contents

|          |  |           |
|----------|--|-----------|
| <b>1</b> | <b>Introduction</b>  | <b>3</b>  |
| <b>2</b> | <b>Preliminaries</b>                                       | <b>4</b>  |
| 2.1      | The mathematical model . . . . .                           | 4         |
| 2.2      | Nomenclature and assumptions . . . . .                     | 5         |
| <b>3</b> | <b>The continuous approach</b>                             | <b>5</b>  |
| 3.1      | Step 1: The restriction procedure . . . . .                | 6         |
| 3.2      | Step 2: The optimization procedure . . . . .               | 8         |
| <b>4</b> | <b>The algebraic approach</b>                              | <b>9</b>  |
| 4.1      | Step 1: the restriction procedure . . . . .                | 9         |
| 4.2      | Step 2: The optimization procedure . . . . .               | 10        |
| <b>5</b> | <b>Numerical investigation</b>                             | <b>13</b> |
| 5.1      | Four plane waves per element . . . . .                     | 14        |
| 5.2      | Eight plane waves per element . . . . .                    | 16        |
| 5.2.1    | Performance assessment in the case of Approach 1 . . . . . | 19        |
| 5.2.2    | Performance assessment in the case of Approach 2 . . . . . | 26        |
| <b>6</b> | <b>Summary and conclusions</b>                             | <b>28</b> |



---

Centre de recherche INRIA Bordeaux – Sud Ouest  
Domaine Universitaire - 351, cours de la Libération - 33405 Talence Cedex (France)

Centre de recherche INRIA Grenoble – Rhône-Alpes : 655, avenue de l'Europe - 38334 Montbonnot Saint-Ismier  
Centre de recherche INRIA Lille – Nord Europe : Parc Scientifique de la Haute Borne - 40, avenue Halley - 59650 Villeneuve d'Ascq  
Centre de recherche INRIA Nancy – Grand Est : LORIA, Technopôle de Nancy-Brabois - Campus scientifique  
615, rue du Jardin Botanique - BP 101 - 54602 Villers-lès-Nancy Cedex  
Centre de recherche INRIA Paris – Rocquencourt : Domaine de Voluceau - Rocquencourt - BP 105 - 78153 Le Chesnay Cedex  
Centre de recherche INRIA Rennes – Bretagne Atlantique : IRISA, Campus universitaire de Beaulieu - 35042 Rennes Cedex  
Centre de recherche INRIA Saclay – Île-de-France : Parc Orsay Université - ZAC des Vignes : 4, rue Jacques Monod - 91893 Orsay Cedex  
Centre de recherche INRIA Sophia Antipolis – Méditerranée : 2004, route des Lucioles - BP 93 - 06902 Sophia Antipolis Cedex

---

Éditeur  
INRIA - Domaine de Voluceau - Rocquencourt, BP 105 - 78153 Le Chesnay Cedex (France)  
<http://www.inria.fr>  
ISSN 0249-6399