



**HAL**  
open science

# A posteriori error estimations and convergence criteria in fast Fourier transform-based computational homogenization

Renaud Ferrier, Cédric Bellis

► **To cite this version:**

Renaud Ferrier, Cédric Bellis. A posteriori error estimations and convergence criteria in fast Fourier transform-based computational homogenization. *International Journal for Numerical Methods in Engineering*, 2023, 124 (4), pp.834-863. 10.1002/nme.7145 . hal-03952103

**HAL Id: hal-03952103**

**<https://hal.science/hal-03952103>**

Submitted on 23 Jan 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.

# A posteriori error estimations and convergence criteria in FFT-based computational homogenization

Renaud Ferrier\*, Cédric Bellis†

January 23, 2023

**Abstract:** A stopping criterion for FFT-based iterative schemes in computational homogenization is proposed and investigated numerically. This criterion is based on the separate evaluation and comparison of the discretization and iteration errors on the computed fields. Some estimators for these errors are proposed and their performances are assessed on a set of 2D problems in the frameworks of both the classical FFT-based methods and these that use a modified version of the featured Green’s operator. In particular, two novel strategies for estimating the discretization error are investigated: either using an image processing approach or transposing to the FFT-based setting the constitutive relation error that is well-established in the context of the finite element method. It is then shown that the resulting stopping criterion leads to a better control of the global error on the computed effective property compared to the classical criterion based on the residual of the iterative scheme alone.

**Key Words:** Error estimators, Fast Fourier Transform, Image Processing, Filtering, Constitutive Relation Error

## 1 Introduction

### 1.1 Context and motivations

Homogenization is commonly employed to model the macroscopic behavior of heterogeneous materials so as to perform efficient numerical computations on extended domains or structures that are made of different constituents at the microscopic scale. It revolves around the calculation or the numerical computation of an equivalent effective homogeneous material from the knowledge of the geometry of the microstructure and of the behavior of the different material phases. In the framework of periodic homogenization, the computation of such an effective model can be achieved via the resolution of an integral convolution equation commonly referred to as the periodic Lippmann-Schwinger equation.

In [20], it has been proposed to discretize the Lippmann-Schwinger equation in the space-based Fourier domain and use a fixed-point algorithm to solve the resulting system in an efficient Fast Fourier Transform (FFT)-based implementation. Since then, numerous variants of this approach have been proposed, either by modifying the discretization method, see e.g. [34, 5, 30, 12], or by using more efficient, i.e. faster, iterative algorithms, see e.g. [6, 36, 11].

Such variants of the FFT-based computational homogenization method amount in an iterative approach associated with a Cartesian grid-based discretization of the Lippmann-Schwinger equation. In this setting, the most commonly used indicator for quantifying the convergence of the solution is the residual of the numerical scheme, which can be interpreted as an indicator of the *iteration error*. If the question of the definition of this indicator has already been investigated in a number of studies, see e.g. [2] and the references therein, the notion of *discretization error estimation* has seldom been considered, see [21]. In this context, the present study aims at filling this gap, by proposing some estimators of the discretization error and a framework to use them in conjunction with an iteration error estimator in order to control the overall error on the computed fields or effective property.

The proposed procedure to estimate the discretization error is inspired by the ideas initially developed for the finite element method (FEM), see the monograph [14], as well as in the context of domain

---

\*Mines Saint-Etienne, Univ Lyon, CNRS, UMR 5307 LGF, Centre SMS, Saint-Etienne, France

†Aix Marseille Univ, CNRS, Centrale Marseille, LMA, Marseille, France

decomposition [24]. Noticeably, the influence of the discretization in a finite volume method and of the incomplete resolution of the associated linear system has also been assessed in [10]. In the present study, to obtain some estimations of the discretization error, we investigate two types of approaches, which are classical for the FEM but still unexploited in the field of FFT-based computational homogenization. The first one draws from the so-called *ZZ1 method* [37] that can be assimilated, in the setting associated with regular computational grids, to smoothing methods used in the field of *image processing*, see also [19] in the present context. The second one is the so-called *constitutive relation error* [13]. To the best of our knowledge, this topic has never been investigated yet, except in a conference paper [4], and the present study constitutes a first step towards the formulation of robust a posteriori error estimators for Fourier-based computational homogenization platforms. Note that, the tools developed in the context of the FEM would also find a natural extension in this context if FE discretizations on regular grids coupled with FFT solvers were used, as in [25, 16].

In the remaining of this section, the homogenization problem for the model conductivity equation is presented, along with the general idea of the proposed approach to distinguish between the discretization and the iteration errors. In section 2, we present the numerical resolution method for the Lippmann-Schwinger equation as well as the test-case geometries considered for the numerical examples to come. Section 3 exposes the relationship between the residual of the discrete problem and the iteration error, but highlights numerically that this residual is not sufficient to provide a reliable indicator for the total error on the effective property. In section 4, three discretization error estimators are introduced and assessed numerically. Finally, Section 5 focuses on the proposition and the numerical evaluation of a stopping criterion based on the comparison between the estimated discretization and iteration errors, in 2D conductivity and 3D elasticity as well.

## 1.2 Problem setting

Consider a periodic medium characterized by the representative unit cell  $\Omega \subset \mathbb{R}^d$  and a spatially-varying isotropic conductivity field  $\gamma \in L_{\text{per}}^\infty(\Omega)$  with  $\gamma > 0$ . The latter constitutes a material parameter prototypical of linear constitutive laws, and a placeholder for possible extensions of the present work to anisotropic or elastic constituents. Our overall objective is to quantify the macroscopic behavior of the medium considered, through the computation of the effective tensor  $\gamma_{\text{eff}}$ , which is the unique tensor defining an effective energy  $W_{\text{eff}}(\bar{\mathbf{e}}) = \frac{1}{2} \bar{\mathbf{e}} \cdot \gamma_{\text{eff}} \cdot \bar{\mathbf{e}}$  for all  $\bar{\mathbf{e}} \in \mathbb{R}^d$ , which satisfies

$$W_{\text{eff}}(\bar{\mathbf{e}}) = \min_{\mathbf{e}^* \in \mathcal{E}_0} W(\bar{\mathbf{e}} + \mathbf{e}^*) \quad \text{with} \quad W(\mathbf{e}) = \frac{1}{2} \langle \mathbf{e}(\mathbf{x}) \gamma(\mathbf{x}) \mathbf{e}(\mathbf{x}) \rangle, \quad (1)$$

where  $\mathbf{e}(\mathbf{x}) = \bar{\mathbf{e}} + \tilde{\mathbf{e}}(\mathbf{x})$  and using the averaging operator  $\langle \cdot \rangle$  defined as

$$\langle \mathbf{f} \rangle = \frac{1}{|\Omega|} \int_{\Omega} \mathbf{f}(\mathbf{x}) \, d\mathbf{x}. \quad (2)$$

In addition, we made use of the following functional space:

$$\mathcal{E}_0 = \{ \mathbf{e}^* \in \mathbf{L}_{\text{per}}^2(\Omega) \mid \exists w^* \in H_{\text{per}}^1(\Omega), \mathbf{e}^* = \nabla w^* \},$$

where  $\mathbf{L}_{\text{per}}^2(\Omega)$  is the subspace of tensor fields with components in  $L_{\text{per}}^2(\Omega)$ . In general  $\gamma_{\text{eff}}$  is a symmetric second-order tensor, i.e. it characterizes an anisotropic effective media. Yet, it can reduce to an isotropic one provided that the microstructure exhibits some symmetries.

Given  $\bar{\mathbf{e}} \in \mathbb{R}^d$  and with  $\tilde{\mathbf{e}}$  being the minimizer of (1), the Euler-Lagrange equations for the energy minimization problem (1) are equivalent to the following system of local equations:

$$\begin{cases} \mathbf{e}(\mathbf{x}) = \bar{\mathbf{e}} + \tilde{\mathbf{e}}(\mathbf{x}), & \tilde{\mathbf{e}} = \nabla u, & u \text{ periodic on } \partial\Omega, \\ \mathbf{j}(\mathbf{x}) = \gamma(\mathbf{x}) \mathbf{e}(\mathbf{x}), \\ \nabla \cdot \mathbf{j}(\mathbf{x}) = 0, & \mathbf{j} \cdot \mathbf{n} \text{ anti-periodic on } \partial\Omega, \end{cases} \quad (3)$$

with  $\mathbf{n}$  being the unit outward normal on  $\partial\Omega$ . Physically,  $\mathbf{e}$  is the electric field with prescribed mean value  $\bar{\mathbf{e}}$  and  $\mathbf{j}$  is the electric current or flux. In addition, the periodic fluctuation  $\tilde{\mathbf{e}}$  is expressed as the

gradient of a scalar potential  $u$  sought in  $H_{\text{per}}^1(\Omega)$ , so that it satisfies the mean-free property  $\langle \tilde{\mathbf{e}} \rangle = \mathbf{0}$ .

We are interested in (1) as it is a prototypical problem in computational homogenization, which has a broad range of applicability as in conductivity, considered here, as well as in elasticity and electromagnetism [18]. The tensor  $\boldsymbol{\gamma}_{\text{eff}}$  can be fully determined component-wise by solving Problem (3) at most  $d$  times using some linearly independent loadings  $\bar{\mathbf{e}}$  and computing the scalar products of the corresponding solutions. Therefore, to be used in the forthcoming analysis, considering  $\alpha \in L_{\text{per}}^\infty(\Omega)$  with  $\alpha > 0$  we introduce the associated energetic scalar product and norm as

$$(\mathbf{f}_1, \mathbf{f}_2)_\alpha = \langle \mathbf{f}_1(\mathbf{x}) \cdot \alpha(\mathbf{x}) \mathbf{f}_2(\mathbf{x}) \rangle \quad \text{and} \quad \|\mathbf{f}\|_\alpha = (\mathbf{f}, \mathbf{f})_\alpha^{1/2}, \quad (4)$$

for all vectors  $\mathbf{f}_1, \mathbf{f}_2 \in \mathbf{L}_{\text{per}}^2(\Omega)$ . As a consequence of this definition, the quadratic energy functional considered in (1) satisfies  $W(\mathbf{f}) = \frac{1}{2} \|\mathbf{f}\|_\gamma^2$ . Lastly, the standard scalar product and norm on  $\mathbf{L}_{\text{per}}^2(\Omega)$  will be denoted as  $(\cdot, \cdot)$  and  $\|\cdot\|$ , respectively.

### 1.3 Objectives

One considers approximating the solution  $\tilde{\mathbf{e}}$  to (1), or equivalently  $u$  to (3), based on (i) the introduction of a subspace  $\mathcal{E}_0^h \subset \mathcal{E}_0$  constructed from a discretization of the domain using a *grid* size  $h$ , suitable basis functions and an integration scheme, and (ii) a given iterative scheme whose iteration number will be denoted by  $k$ , the combination of which yields an approximation  $\tilde{\mathbf{e}}_k^h$ . In this context, the overall objective of this study is to propose some convergence criteria for the effective property  $\boldsymbol{\gamma}_{\text{eff}}$ , which according to (1) amounts in quantifying the error  $\delta_{k,h}^{\text{eff}}$  in the effective energy as

$$\delta_{k,h}^{\text{eff}} \stackrel{\text{def}}{=} |W(\mathbf{e}_k^h) - W(\mathbf{e})|^{1/2} \quad \text{with} \quad \mathbf{e}_k^h = \bar{\mathbf{e}} + \tilde{\mathbf{e}}_k^h, \quad (5)$$

for a given  $\bar{\mathbf{e}} \in \mathbb{R}^d$ . Upon introducing the limit  $\mathbf{e}_\infty^h = \bar{\mathbf{e}} + \tilde{\mathbf{e}}_\infty^h$  of the iterative scheme considered (provided that it converges), the error above will be related in the following to the total error  $\delta_{k,h}^{\text{tot}}$  on the fields, in the energetic norm, which satisfies

$$\delta_{k,h}^{\text{tot}} \stackrel{\text{def}}{=} \|\mathbf{e}_k^h - \mathbf{e}\|_\gamma \leq \underbrace{\|\mathbf{e}_k^h - \mathbf{e}_\infty^h\|_\gamma}_{\delta_k^h} + \underbrace{\|\mathbf{e}_\infty^h - \mathbf{e}\|_\gamma}_{\delta^h}, \quad (6)$$

where  $\delta_k^h$  is the *iteration* error expressed in terms of the converged discrete solution  $\mathbf{e}_\infty^h$  and  $\delta^h$  is the *discretization* error expressed in terms of the continuous (exact) solution to (1). Note that the above errors and estimations are formulated using the *continuous* norm introduced in (4), an issue that we will return to in the ensuing analysis.

As neither the effective property  $\boldsymbol{\gamma}_{\text{eff}}$  nor the continuous solution are known a priori, the errors introduced above are not accessible directly. Therefore, we aim at proposing some a posteriori estimations, denoted as  $\Delta_k^h$  and  $\Delta^h$  for the iteration and discretization errors respectively, which will be readily accessible during computations. Moreover, we also intend to use these error estimators to propose a stopping criterion for the iterative scheme.

**Remark 1.** *In this study, we only consider the iteration and discretization errors, i.e. we do not take into account a possible inaccurate representation of the geometry of the microstructure. In other words, the discrete geometry is somehow considered to be the exact one in the ensuing analysis.*

## 2 Solution methods

### 2.1 Green's operator and linear equation

Considering a homogeneous comparison medium with conductivity  $\gamma_0 > 0$ , we introduce the associated Green's operator  $\boldsymbol{\Gamma}_0$  on  $\mathbf{L}_{\text{per}}^2(\Omega)$  that is defined as

$$\boldsymbol{\Gamma}_0 : \boldsymbol{\tau} \mapsto \mathbf{e}^* = \boldsymbol{\Gamma}_0 \boldsymbol{\tau} \quad \text{with} \quad \mathbf{e}^* \in \mathcal{E}_0 \text{ and } \mathbf{s} = (\gamma_0 \mathbf{e}^* - \boldsymbol{\tau}) \in \mathcal{S}, \quad (7)$$

with the functional space  $\mathcal{S}$  being given by

$$\mathcal{S} = \{ \mathbf{s} \in \mathbf{L}_{\text{per}}^2(\Omega) \mid \nabla \cdot \mathbf{s}(\mathbf{x}) = 0, \mathbf{s} \cdot \mathbf{n} \text{ anti-periodic on } \partial\Omega \}.$$

With a slight abuse of notation, we have  $\mathcal{S} = \mathcal{E}_0^\perp$  for the standard  $\mathbf{L}_{\text{per}}^2$ -scalar product [18, 3]. Moreover, it is known that  $\mathbf{\Gamma}_0\gamma_0$  is an orthogonal projector (self-adjoint) onto  $\mathcal{E}_0$  for the energetic scalar product  $(\cdot, \cdot)_{\gamma_0}$ . Therefore, for the standard inner product,  $\mathbf{\Gamma}_0\gamma_0$  is an oblique projector (non-self-adjoint) onto  $\mathcal{E}_0$  and it is straightforward to show [30] that its adjoint is  $\gamma_0\mathbf{\Gamma}_0$ .

Now, the local problem (3) is equivalent to the following weak formulation: find  $u \in H_{\text{per}}^1(\Omega)$  such that

$$a(u, v) = \ell(v) \quad \forall v \in H_{\text{per}}^1(\Omega), \quad (8)$$

with  $a$  and  $\ell$  being respectively the bilinear and linear forms defined on  $H_{\text{per}}^1(\Omega)$  as:

$$a(u, v) = \int_{\Omega} \gamma \nabla u \cdot \nabla v \, d\mathbf{x} \quad \text{and} \quad \ell(v) = - \int_{\Omega} \gamma \bar{\mathbf{e}} \cdot \nabla v \, d\mathbf{x}.$$

The identity (8) is equivalent to

$$(\gamma(\bar{\mathbf{e}} + \tilde{\mathbf{e}}), \mathbf{e}^*) = 0 \quad \forall \mathbf{e}^* = \nabla v \in \mathcal{E}_0, \quad (9)$$

so that  $\mathbf{j} = \gamma(\bar{\mathbf{e}} + \tilde{\mathbf{e}})$  belongs to the subspace orthogonal to  $\mathcal{E}_0 = \text{Im}(\mathbf{\Gamma}_0\gamma_0)$  for the standard  $\mathbf{L}_{\text{per}}^2$ -scalar product. Therefore, owing to the properties of the adjoint operator, one has  $\mathbf{j} \in \text{Im}(\mathbf{\Gamma}_0\gamma_0)^\perp = \text{Ker}(\gamma_0\mathbf{\Gamma}_0)$ , which according to the above entails  $\mathbf{\Gamma}_0\mathbf{j} = \mathbf{0}$ . The latter identity can finally be recast as the following linear equation for  $\tilde{\mathbf{e}} \in \mathcal{E}_0$ :

$$\mathbf{A}\tilde{\mathbf{e}} = \mathbf{b}, \quad \text{with} \quad \mathbf{A} = \mathbf{\Gamma}_0\gamma : \mathcal{E}_0 \rightarrow \mathcal{E}_0 \quad \text{and} \quad \mathbf{b} = -\mathbf{\Gamma}_0\gamma\bar{\mathbf{e}} \in \mathcal{E}_0. \quad (10)$$

Note that this equation can equivalently be obtained from the first-order optimality condition for the minimization problem (1) by recognizing that the gradient of  $W$  in  $\mathcal{E}_0$  endowed with the energetic scalar product  $(\cdot, \cdot)_{\gamma_0}$  writes as  $\nabla W(\mathbf{e}^*) = \mathbf{\Gamma}_0\gamma(\bar{\mathbf{e}} + \mathbf{e}^*)$  [11, 3].

Finally, in the periodic setting considered, the Green's operator (7) can be expressed in closed-form using the Fourier transform  $\mathcal{F}$ , see Appendix A, as

$$\mathbf{\Gamma}_0\boldsymbol{\tau}(\mathbf{x}) = \mathcal{F}^{-1} \left[ \hat{\mathbf{\Gamma}}_0(\boldsymbol{\xi}) \cdot \mathcal{F}[\boldsymbol{\tau}](\boldsymbol{\xi}) \right](\mathbf{x}) \quad \forall \mathbf{x} \in \Omega, \quad (11)$$

with the symmetric second-order tensor  $\hat{\mathbf{\Gamma}}_0(\boldsymbol{\xi})$  being defined in the Fourier space by

$$\hat{\mathbf{\Gamma}}_0(\mathbf{0}) = \mathbf{0} \quad \text{and} \quad \hat{\mathbf{\Gamma}}_0(\boldsymbol{\xi}) = \frac{\boldsymbol{\xi} \otimes \boldsymbol{\xi}}{\gamma_0|\boldsymbol{\xi}|^2} \quad \forall \boldsymbol{\xi} \in \mathcal{R}^* \setminus \{\mathbf{0}\}. \quad (12)$$

## 2.2 Discretization and numerical methods

### 2.2.1 Discrete problem

To transpose the continuous linear problem (10) into a discrete setting, we consider a regular grid of size  $h$ , with the corresponding set of interpolation points being  $\{\mathbf{x}_i\}$  where  $i$  a multi-index in dimension  $d$ , and the approximation subspace  $\mathcal{T}^h \subset \mathbf{L}_{\text{per}}^2(\Omega)$  being generated by the trigonometric polynomials associated with the Discrete Fourier Transform (DFT) on the grid considered. Correspondingly, consider the discrete subspace  $\mathcal{E}_0^h = \mathcal{E}_0 \cap \mathcal{T}^h$  of mean-free gradient fields, relatively to the discrete version  $\langle \cdot \rangle_h$  of the averaging operator (2) defined as

$$\langle \mathbf{f} \rangle_h = \sum_i h^d \mathbf{f}(\mathbf{x}_i), \quad (13)$$

a quadrature scheme that coincides with the trapezoidal rule [2]. The definition (13) also yields a discrete version  $(\cdot, \cdot)_{h,\alpha}$  of the scalar product (4), as

$$(\mathbf{f}_1, \mathbf{f}_2)_{h,\alpha} = \sum_i h^d \mathbf{f}_1(\mathbf{x}_i) \cdot \alpha(\mathbf{x}_i) \mathbf{f}_2(\mathbf{x}_i)$$

whose associated norm  $\|\cdot\|_{h,\alpha}$  defines the discrete energy functional

$$W_h(\mathbf{f}) = \frac{1}{2} \|\mathbf{f}\|_{h,\gamma}^2. \quad (14)$$

**Remark 2.** *In the ensuing analysis, we will have to deal with the continuous norms of some discrete fields. To do so, given a discrete field, say  $\mathbf{f}^h$  in a discrete subspace such as  $\mathcal{T}^h$ , then in the evaluation of its continuous norm  $\|\cdot\|_\alpha$  we consider that  $\mathbf{f}^h$  is extended as a piecewise-constant field on each element or pixel of the grid considered. By doing so it holds*

$$\|\mathbf{f}^h\|_\alpha = \|\mathbf{f}^h\|_{h,\alpha}.$$

*This convention will allow us to evaluate quantities such as  $\|\mathbf{f}^h - \mathbf{f}\|_\alpha$ , where  $\mathbf{f}$  is a continuous field. Note that we do not interpolate  $\mathbf{f}^h$  in the associated basis of continuous functions, which in the case of the Fourier interpolation would give continuous fields oscillating between sampling points that would in turn yield an artificially high value of  $\|\mathbf{f}^h - \mathbf{f}\|_\alpha$ .*

In addition, a DFT-based version  $\mathbf{\Gamma}_0^h$  of the Green's operator in (10) is commonly defined by the straightforward transposition of (11–12) to the DFT setting. Let us underline that it has also been proposed [32, 5, 34, 33, 26] to modify the definition of the discrete Green's operator  $\mathbf{\Gamma}_0^h$  by making use of staggered grids or introducing in (12) a filtering of the set of discrete spatial frequencies  $\{\boldsymbol{\xi}_i\}$ . The latter strategy is equivalent to approximating the continuous problem (3) using a finite-differences scheme, whose degree of approximation can be directly related to the chosen frequency filter [7]. In the present study, the considered modification of the Green's operator corresponds to the first degree central finite-differences scheme.

In this setting, one considers the corresponding discrete counterpart of the weak formulation (9) with the inexact integration scheme (13), which leads to the following equation for the sought discrete solution  $\bar{\mathbf{e}}^h \in \mathcal{E}_0^h$ :

$$(\gamma^h(\bar{\mathbf{e}} + \tilde{\mathbf{e}}^h), \mathbf{e}^{*h})_h = 0 \quad \forall \mathbf{e}^{*h} \in \mathcal{E}_0^h, \quad (15)$$

where  $\gamma^h$  denotes the constitutive conductivity field evaluated at the DFT interpolation points, i.e.  $\gamma^h = \{\gamma(\mathbf{x}_i)\}$ .

**Remark 3.** *As a consequence, the discrete conductivity  $\gamma^h$  is considered to be equal to the exact one  $\gamma$  at all the discretization points.*

Now, the equation (15) implies that  $\mathbf{j}^h = \gamma^h(\bar{\mathbf{e}} + \tilde{\mathbf{e}}^h)$  belongs to the orthogonal subspace of  $\mathcal{E}_0^h$  for the (standard) discrete scalar product on  $\mathcal{T}^h$ . Upon noting that the Helmholtz decomposition of  $\mathbf{L}_{\text{per}}^2(\Omega)$  (as the direct sum of the mutually orthogonal subspaces of uniform, mean-free gradient and mean-free divergence-free fields) can be prolonged to  $\mathcal{T}^h$ , it has been shown [28, 30] that  $\mathbf{\Gamma}_0^h \gamma_0$  defines a projector onto  $\mathcal{E}_0^h$ . Again, this operator is non-orthogonal for the standard scalar product while it is self-adjoint for the induced energetic one  $(\cdot, \cdot)_{h,\gamma_0}$ . Therefore, (15) leads to the following linear system:

$$\mathbf{A}^h \bar{\mathbf{e}}^h = \mathbf{b}^h, \quad \text{with} \quad \mathbf{A}^h = \mathbf{\Gamma}_0^h \gamma^h : \mathcal{E}_0^h \rightarrow \mathcal{E}_0^h \quad \text{and} \quad \mathbf{b}^h = -\mathbf{\Gamma}_0^h \gamma^h \bar{\mathbf{e}} \in \mathcal{E}_0^h. \quad (16)$$

This is the system under investigation and it should be noted that it will be handled as is, even if a modified version of the Green's operator is used.

**Remark 4.** *In the computational treatment of composite materials, it can be expected that the geometry of the microstructure cannot be represented exactly on a regular grid. The associated error on the geometry is not explicitly taken into account in this study. However, this error appears implicitly in the difference between the continuous and discrete integral forms, such as  $(\cdot, \cdot)_\gamma$  and  $(\cdot, \cdot)_{h,\gamma}$  respectively.*

### 2.2.2 Iterative schemes

To solve the linear system (16) numerically, we intend to use an iterative scheme that will yield an approximation  $\tilde{\mathbf{e}}_k^h = \nabla u_k^h$  of the discrete solution  $\tilde{\mathbf{e}}^h$  at a given iterate  $k$ . Provided that the scheme considered converges, we formally get

$$\tilde{\mathbf{e}}_k^h \xrightarrow[k \rightarrow \infty]{} \tilde{\mathbf{e}}_\infty^h \equiv \tilde{\mathbf{e}}^h. \quad (17)$$

Here, two of the main approaches used in practice are adopted, namely a fixed-point scheme and the conjugate gradient method:

- (i) Owing to the orthogonal decomposition of the discrete approximation space  $\mathcal{T}^h$  and using again the interpretation of  $\mathbf{\Gamma}_0^h \gamma_0$  as a projector, we have  $\mathbf{\Gamma}_0^h \gamma_0(\bar{\mathbf{e}} + \tilde{\mathbf{e}}^h) = \tilde{\mathbf{e}}^h$  so that (16) can be recast as

$$(\mathbf{I} + \mathbf{\Gamma}_0^h \delta \gamma^h) \mathbf{e}^h = \bar{\mathbf{e}}, \quad \text{with} \quad \mathbf{e}^h = \bar{\mathbf{e}} + \tilde{\mathbf{e}}^h$$

and where  $\delta \gamma^h = \gamma^h - \gamma_0$ . The equation above is a Lippmann-Schwinger integral equation for the total field  $\mathbf{e}^h$ , which can be inverted through a Neumann series expansion. This coincides with the fixed-point algorithm originally proposed [20]. The corresponding results shown hereafter will be labeled as FP or FPM provided that the Green's operator is used in its original or a modified version (as described in Section 2.2.1), respectively.

- (ii) Alternatively, the linear system (16) can be solved using a conjugate gradient solver [36]. Note that, in general, care must be taken to use an energetic scalar product for which the operator  $\mathbf{A}^h$  to be inverted is symmetric. Note that in the case of isotropic conductivity considered here some simplifications occurs. Accordingly, for any discrete field  $\check{\mathbf{e}}^h \in \mathcal{E}_0^h$ , one introduces the corresponding residual:

$$\check{\mathbf{r}}^h(\check{\mathbf{e}}^h) = \mathbf{A}^h \check{\mathbf{e}}^h - \mathbf{b}^h, \quad (18)$$

which is minimized in a norm weighted by  $(\mathbf{A}^h)^{-1}$ . Indeed, the quadratic minimization problem associated with the linear system (16) satisfies

$$\tilde{\mathbf{e}}^h = \arg \min_{\check{\mathbf{e}}^h \in \mathcal{E}_0^h} \left( \frac{1}{2} (\check{\mathbf{e}}^h, \mathbf{A}^h \check{\mathbf{e}}^h)_h - (\check{\mathbf{e}}^h, \mathbf{b}^h)_h \right) = \arg \min_{\check{\mathbf{e}}^h \in \mathcal{E}_0^h} \|\check{\mathbf{r}}^h(\check{\mathbf{e}}^h)\|_{h, (\mathbf{A}^h)^{-1}}^2.$$

The associated results below will be labelled as CG or CGM depending on the version (original or modified) of the Green's operator employed.

**Remark 5.** *In both approaches, the computed iterates  $\tilde{\mathbf{e}}_k^h$  are compatible fields in the sense of the Fourier-based discretization considered, i.e.  $\tilde{\mathbf{e}}_k^h \in \mathcal{E}_0^h = \mathcal{E}_0 \cap \mathcal{T}^h$ , with  $\mathcal{E}_0$  being the subspace of mean-free gradient fields and  $\mathcal{T}^h$  the approximation subspace generated by trigonometric polynomials. The schemes are initialized by setting  $\tilde{\mathbf{e}}_0^h = \mathbf{0}$ .*

## 2.3 Test-cases

In this study, two distinct piecewise-homogeneous 2D material distributions are considered to illustrate the proposed approaches, see Figure 1. Different conductivity ratio  $\rho$  will be considered, with the definition  $\rho = \gamma_i/\gamma_m$  where  $\gamma_i$  is the conductivity of the inclusion (or the set thereof) in red and  $\gamma_m$  this of of the matrix phase in blue.

The first geometry presents a regular pattern of square inclusions, with a surface fraction of 1/4, see Fig. 1a. In [22], the corresponding effective conductivity has been calculated analytically, which allows a reliable computation of the corresponding error by (5). The effective conductivity  $\gamma_{\text{eff}}^{\text{square}}$  is found to be isotropic and writes as

$$\gamma_{\text{eff}}^{\text{square}} = \gamma_m \sqrt{\frac{\gamma_m + 3\gamma_i}{3\gamma_m + \gamma_i}}.$$

This test-case presents the further advantage that the discretized geometry conforms strictly to the exact geometry. In addition, the continuous solution of the conductivity problem exhibits integrable

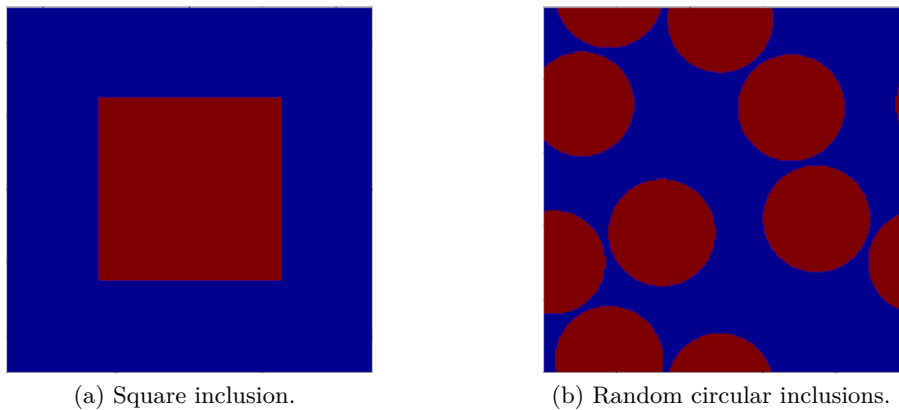


Figure 1: Geometries defining the test-cases considered.

singularities at the corners [22], which poses some difficulties in the computation of an approximation. Such singularities are typical of geometries with corners, which makes this problem interesting in the study of errors estimators.

The second geometry consists in a random pattern of circular inclusions, see Fig. 1b. The inclusions have a random radius satisfying  $0.14 \leq r \leq 0.15$ , the minimal distance between two inclusions is 0.01, and the phase fraction of the inclusions is equal to  $1/4$ . The geometry in this case is not exactly rendered by the discretization, but the solution is regular (provided, as it is ensured here, a minimal exclusion distance is imposed between two inclusions). A reference numerical solution for this problem is computed using the finite element method on a mesh with P1 triangular elements and approximately  $2^8$  nodes on one side of the unit square. The mesh has 88337 nodes in total. The associated conductivity will be considered as the reference value.

To refer to the different configurations, we adopt a compact notation that indicates, in that order: the shape of the geometry (`squ` or `cir`), the phase ratio  $\rho$ , the number of pixels, and the solver used. For example, a computation labeled `[squ;rho=1E3;npix=2**6;FP]` is done on the square inclusion case, with a conductivity ratio equal to  $10^3$ , a discretization of  $2^6 \times 2^6$  pixels, and the fixed point method applied to the non-modified operator.

Finally, while the computation of the complete anisotropic effective conductivity  $\gamma_{\text{eff}}$  requires to solve Problem (3) twice with  $\bar{e}$  describing an orthonormal basis of  $\mathbb{R}^2$ , in the numerical experiments shown hereafter, only the loading  $\bar{e} = (1, 0)$  is considered. Accordingly, Figure 2 presents the maps of the components of the field  $\tilde{e}_k^h = \nabla u_k^h$  computed for the two test-cases considered and obtained at convergence with the CG variant of the FFT-based method on a discretization of  $2^{10} \times 2^{10}$  pixels.

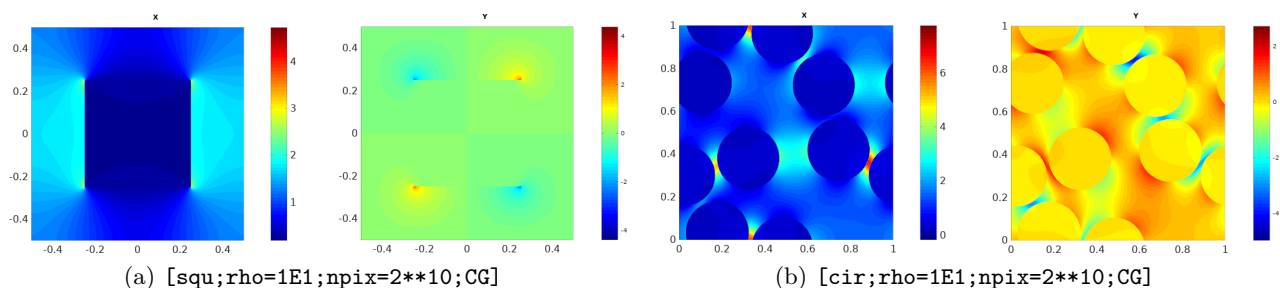


Figure 2: Components of the numerical solution  $\tilde{e}_k^h = \nabla u_k^h$  for the two geometries considered and  $\bar{e} = (1, 0)$ .



### 3 Towards a convergence criterion

#### 3.1 Estimation of the iteration error

##### 3.1.1 Relationship between the residual and the iteration error

As announced, since the converged solution  $\tilde{\mathbf{e}}_\infty^h$  is not accessible a priori, an estimator  $\Delta_k^h$  of the exact iteration error  $\delta_k^h$  is needed. In the case of a gradient descent algorithm, such as the conjugate gradient (CG) considered here, it makes sense to directly use a norm of the residual  $\mathbf{r}_k^h$  in (18) as a measure of convergence. Here, we consider the following energetic norm:

$$\Delta_k^h = \|\mathbf{r}_k^h\|_{h,\gamma} = \|\mathbf{\Gamma}_0^h \gamma^h (\bar{\mathbf{e}} + \tilde{\mathbf{e}}_k^h)\|_{h,\gamma}. \quad (19)$$

This proposition is slightly different from the one in [3], where the norm  $\|\mathbf{r}_k^h\|_{h,\gamma_0} = \|\nabla W_h(\tilde{\mathbf{e}}_k^h)\|_{h,\gamma_0}$  is considered, with the featured gradient of  $W_h$  being computed in the approximation subspace  $\mathcal{E}_0^h$  endowed with the energetic scalar product  $(\cdot, \cdot)_{h,\gamma_0}$ . Yet, it should be noted that, provided that  $\min \gamma^h \leq \gamma_0 \leq \max \gamma^h$  as it will be ensured later on, then the different energetic norms considered are equivalent as it holds

$$\frac{\min \gamma^h}{\gamma_0} \|\mathbf{f}\|_{h,\gamma_0}^2 \leq \|\mathbf{f}\|_{h,\gamma}^2 \leq \frac{\max \gamma^h}{\gamma_0} \|\mathbf{f}\|_{h,\gamma_0}^2.$$

Here, (19) is introduced to be consistent with the approximation of the discretization error in Section 4, and we generalize its use to all of the numerical schemes considered in the present study.

Considering the definitions (18) of the residual  $\mathbf{r}_k^h$  and (17) of the discrete solution  $\tilde{\mathbf{e}}^h$  obtained at convergence when  $k \rightarrow \infty$ , it holds

$$\mathbf{r}_k^h = \mathbf{\Gamma}_0^h \gamma^h (\tilde{\mathbf{e}}_k^h - \tilde{\mathbf{e}}^h) = \mathbf{A}^h (\tilde{\mathbf{e}}_k^h - \tilde{\mathbf{e}}^h).$$

Now, the featured discrete operator  $\mathbf{A}^h$  is symmetric, and hence diagonalizable, in the energetic norm  $(\cdot, \cdot)_{h,\gamma}$ . Indeed, for all  $\mathbf{f}_1, \mathbf{f}_2$ , it holds

$$(\mathbf{A}^h \mathbf{f}_1, \mathbf{f}_2)_{h,\gamma} = (\mathbf{\Gamma}_0^h \gamma^h \mathbf{f}_1, \gamma^h \mathbf{f}_2)_h = (\gamma^h \mathbf{f}_1, \mathbf{\Gamma}_0^h \gamma^h \mathbf{f}_2)_h = (\mathbf{f}_1, \mathbf{A}^h \mathbf{f}_2)_{h,\gamma},$$

where the second equality makes use of the reciprocity identity satisfied by  $\mathbf{\Gamma}_0^h$  [3, Lemma 3]. As a consequence, upon introducing the lowest and largest eigenvalues of  $\mathbf{A}^h$ , denoted as  $\lambda_{\min}$  and  $\lambda_{\max}$  respectively, one gets

$$\lambda_{\min} \|\tilde{\mathbf{e}}_k^h - \tilde{\mathbf{e}}^h\|_{h,\gamma} \leq \|\mathbf{A}^h (\tilde{\mathbf{e}}_k^h - \tilde{\mathbf{e}}^h)\|_{h,\gamma} \leq \lambda_{\max} \|\tilde{\mathbf{e}}_k^h - \tilde{\mathbf{e}}^h\|_{h,\gamma},$$

which, owing to (19) and Remark 2, can be rewritten as

$$\lambda_{\min} \delta_k^h \leq \Delta_k^h \leq \lambda_{\max} \delta_k^h.$$

Finally, as is well known [17], the eigenvalues considered are bounded as

$$\frac{\min \gamma^h}{\gamma_0} \leq \lambda_{\min} \quad \text{and} \quad \lambda_{\max} \leq \frac{\max \gamma^h}{\gamma_0}, \quad (20)$$

and, a common choice for the reference medium that ensures convergence of the fixed-point scheme described in Section 2.2.2 is  $\gamma_0 = \frac{1}{2}(\min \gamma^h + \max \gamma^h)$ . Therefore, upon introducing the contrast  $c = \max \gamma^h / \min \gamma^h \geq 1$  then one arrives at the following bounds:

$$\frac{2}{1+c} \delta_k^h \leq \Delta_k^h \leq \frac{2c}{1+c} \delta_k^h. \quad (21)$$

As a consequence, for very large contrasts then the estimated iteration error  $\Delta_k^h$  can only be bounded as

$$0 \leq \Delta_k^h \leq 2\delta_k^h \quad \text{when} \quad c \gg 1,$$

whereas, conversely, for very small contrasts one has

$$\Delta_k^h \sim \delta_k^h \quad \text{when} \quad c \sim 1.$$

### 3.1.2 Numerical examples

In this section, we consider the case of random circular inclusions, for different values of the conductivity ratio, and a spacial discretization of  $2^6 \times 2^6$  pixels. We illustrate numerically the obtained estimates (21) by considering the scaled ratio  $(1+c)\Delta_k^h/\delta_k^h$ .

This quantity is plotted against iterations  $k$  in Figure 3, for the different schemes considered and for various contrasts with  $\rho = c$ . Note that the contrast is inverted in the case of voids with the conductivity ratio between phases being set in this particular case as  $\rho = 10^{-6} = c^{-1}$ . To help distinguish between the different cases considered, the multiplying factor  $(1+c)$  has been introduced in the representation of the ratio  $\Delta_k^h/\delta_k^h$ . On the downside, this choice artificially and excessively amplifies the relative variations of these quantities. In addition, both axes are in log-scale. Owing to (21), the lower bound is 2 independently of the contrast while the upper bound is equal to  $2c$ . Yet, on most of the numerical examples reported in Fig. 3 one observes that the quantity  $(1+c)\Delta_k^h/\delta_k^h$  considered is about of the order of the contrast, i.e. the error ratio satisfies  $\Delta_k^h/\delta_k^h = \mathcal{O}(1)$ . Note that deviations from this behavior can be observed in some high-contrast cases. In addition, an asymptotic regime is met in all cases, i.e. the residual  $\Delta_k^h$  stabilizes around a value proportional to the exact iteration error  $\delta_k^h$ , typically in a few dozen iterations, as will be considered in the ensuing numerical experiments.

Clearly, the residual  $\Delta_k^h$  is a reliable estimate of the exact iteration error  $\delta_k^h$  when the contrast is low. This appears to also hold in the high contrast cases in examples considered, with the residual and error being of the same order of magnitude despite the loose bound (21). The residual will consequently be used as an estimate of the iteration error in the remainder of this article.

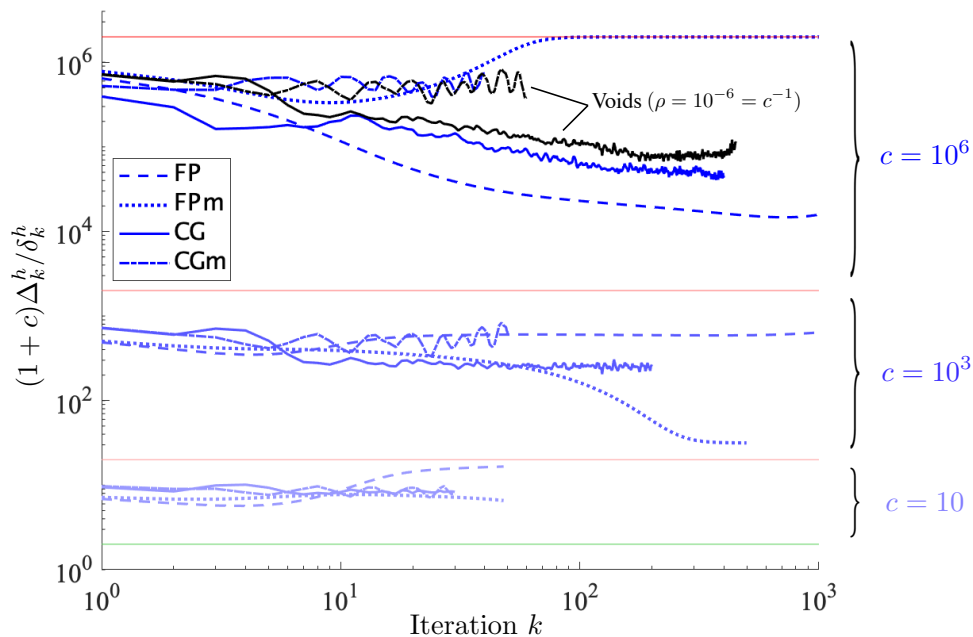


Figure 3: Ratio of the residual  $\Delta_k^h$  by the exact iteration error  $\delta_k^h$ , scaled by the contrast, in the random circular inclusions case for various contrasts and for the different schemes considered. The green line indicates the lower bound equal to 2 and the red lines corresponds to the upper bound  $2c$  in (21).

In Equation (21), the lower bound for  $\Delta_k^h$  would be attained if the residual  $(\tilde{\mathbf{e}}_k^h - \tilde{\mathbf{e}}^h)$  were only spanned by the eigenvectors of  $\mathbf{A}^h$  associated with its lowest eigenvalue  $\lambda_{\min}$ . In order to visualize the projection of this residual onto the eigenvectors of  $\mathbf{A}^h$ , we build the matrix  $\mathbf{M}$  that represents  $\mathbf{A}^h$  in an energy-orthogonal basis  $\{\mathbf{t}_i\}_{i=1,\dots,N}$  of  $\mathcal{E}_0^h = \mathcal{E}_0 \cap \mathcal{T}^h$ , i.e. we have:

$$\begin{cases} (\mathbf{t}_i, \mathbf{t}_j)_{h,\gamma} = \delta_{ij} \\ \mathbf{M}_{ij} = (\mathbf{t}_i, \mathbf{A}^h \mathbf{t}_j)_{h,\gamma} \end{cases} \quad \forall i, j \in \{1, \dots, N\}.$$

In this setting, we display on Figure 4 the numerically-computed eigenvalues  $\lambda_i$  of  $M$  in the case of a conductivity ratio  $\rho = 10^3$ . The vector  $\chi$  is superposed on the same figure; it corresponds to the decomposition on the eigenbasis considered of the residual  $(\tilde{e}_{30}^h - \tilde{e}^h)$  at the 30th fixed-point iteration, i.e. for all  $i = 1, \dots, n$  we have

$$\chi_i = \sum_{j=1}^n U_{ji} \left( t_j, \tilde{e}_{30}^h - \tilde{e}^h \right)_{h,\gamma} \quad \text{with} \quad \begin{cases} U^T M U = \Lambda, \\ U^T U = I. \end{cases}$$

where  $I$  is the identity matrix and  $\Lambda = \text{diag}(\lambda_i)_{i=1, \dots, n}$ . This is illustrated on Figure 4 for the fixed-point method using the standard or the modified operator, FP and FPM respectively. One can

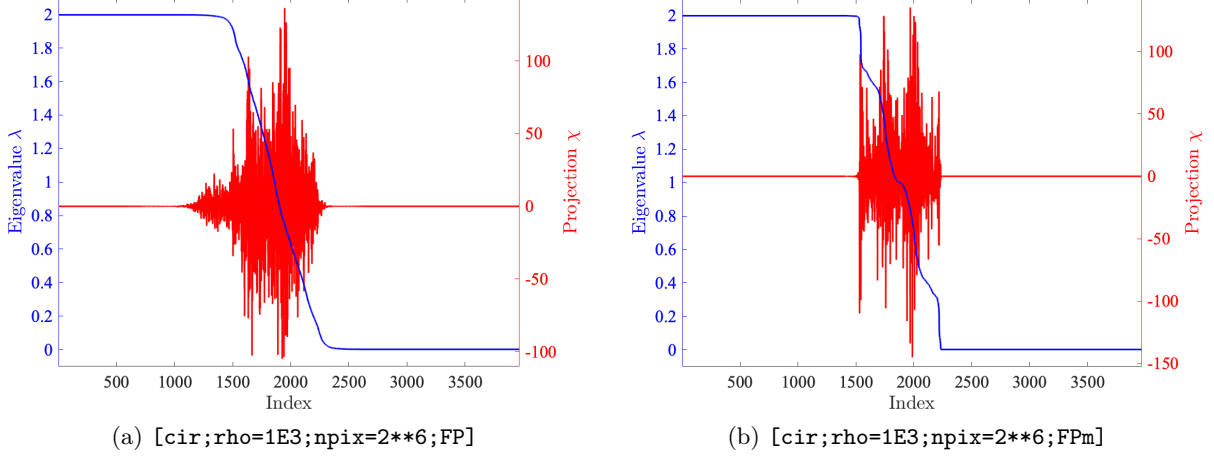


Figure 4: Decomposition of the residual  $(\tilde{e}_{30}^h - \tilde{e}^h)$  at the 30th fixed-point iteration on the discrete eigenvectors of  $A^h$ .

observe a behavior that has already been highlighted [2]: there are two massively multiple eigenvalues (corresponding the bounds (20) and associated to the two distinct values of the local conductivity  $\gamma$ ), that fill the majority of the spectrum, and on which the projection of the residual, and in fact of the solution itself, is zero. Rather, the solution can be decomposed on a smaller number of eigenvectors corresponding to intermediate eigenvalues. For both the non-modified and the modified operator, the overall contribution of each term  $\chi_i$  to the residual are rather equally distributed among these eigenvalues. Accordingly, the ratio  $\Delta_k^h / \delta_k^h$  is expected to be weighted in between  $\lambda_{\min}$  and  $\lambda_{\max}$ .

## 3.2 Examples of convergence behaviors in terms of error and residual

### 3.2.1 Error on the energy

Now that we have assessed the use of the residual  $\Delta_k^h$  in (19) as a reliable indicator of the iteration error, one investigates numerically here its correlation with the error  $\delta_{k,h}^{\text{eff}}$  in the effective energy (5) or with the total error  $\delta_{k,h}^{\text{tot}}$  on the fields (6). It should be noted that the theoretical errors in (5) and (6) are defined relatively to the continuous energetic norm  $\|\cdot\|_\gamma$ . In this context, since  $\tilde{e} \in \mathcal{E}_0$  is the solution of the weak formulation (9) of the problem, and introducing the solution  $\tilde{e}_\star^h$  to (9) in the approximation subspace  $\mathcal{E}_0^h \subset \mathcal{E}_0$  and with exact integration, we have the following Galerkin orthogonality relations:

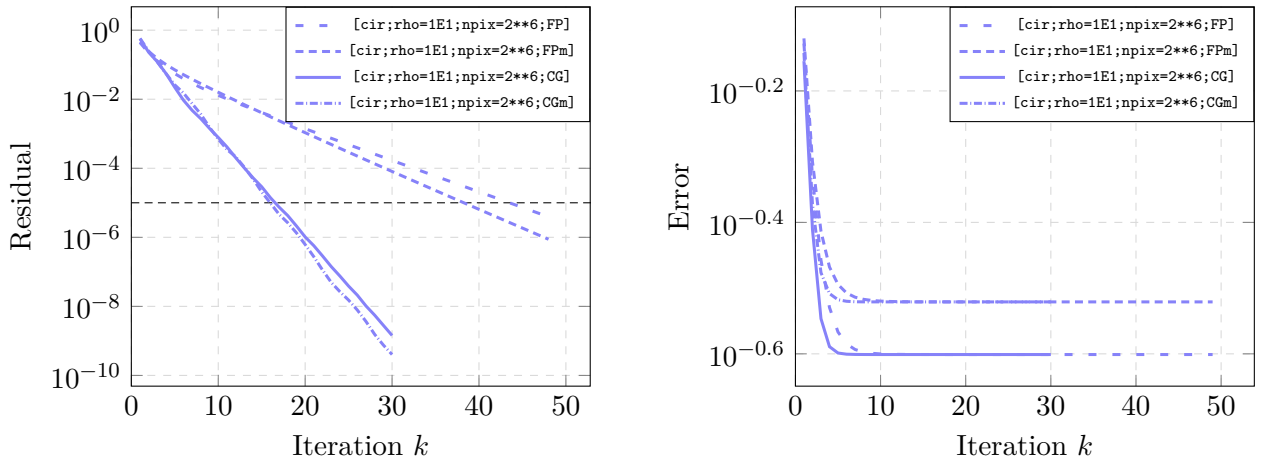
$$(\gamma(\bar{e} + \tilde{e}), \tilde{e}) = 0 \quad \text{and} \quad (\gamma(\bar{e} + \tilde{e}), \tilde{e}_\star^h) = 0, \quad (22)$$

which could formally be rewritten using the energetic scalar product  $(\cdot, \cdot)_\gamma$ . Moreover, considering the following identity

$$\|\tilde{e} - \tilde{e}_\star^h\|_\gamma^2 = \|\bar{e} + \tilde{e}\|_\gamma^2 + \|\bar{e} + \tilde{e}_\star^h\|_\gamma^2 - 2(\bar{e} + \tilde{e}, \bar{e} + \tilde{e}_\star^h)_\gamma,$$

then the last right-hand side term simplifies as  $(\bar{e} + \tilde{e}, \bar{e} + \tilde{e}_\star^h)_\gamma = \|\bar{e} + \tilde{e}\|_\gamma^2$  according to (22). Therefore, we finally get:

$$\frac{1}{2} \|\tilde{e} - \tilde{e}_\star^h\|_\gamma^2 = W(e_\star^h) - W(e^h) \geq 0, \quad (23)$$



(a) Residual  $\Delta_k^h$  and threshold at  $10^{-5}$  (black dashed line)

(b) Error  $d_{k,h}^{\text{eff}}$

Figure 5: Residual and error for the four methods considered.

where we defined  $e_\star^h = \bar{e} + \tilde{e}_\star^h$  and used that  $W(e_\star^h) = \frac{1}{2}\|\bar{e} + \tilde{e}_\star^h\|_\gamma^2$  and  $W(e^h) = \frac{1}{2}\|\bar{e} + \tilde{e}^h\|_\gamma^2$ . In other words, the square error on the fields in energetic norm is equal to the error on the energy. Note that this result has already been discussed extensively [28, 30, 31, 29], in particular to construct guaranteed numerical bounds on the effective properties, but the argument is reproduced here for the reader's convenience.

Noticeably, the above developments are relative to the solution  $\tilde{e}_\star^h$  that satisfies the weak formulation (9) in the approximation subspace  $\mathcal{E}_0^h$  and with the integration being performed exactly, which is a crucial argument in the previous developments. The approximate solution  $\tilde{e}^h \in \mathcal{E}_0^h$  we deal with however, is solution of the weak formulation (15) where the integration is performed numerically according to the inexact trapezoidal rule (13). As a consequence, an inequality such as (23) is not guaranteed anymore for the numerical solution  $\tilde{e}^h$  considered, so that according to (14):

$$\text{the inequality } W_h(e^h) \geq W(e) \text{ cannot be ensured.} \quad (24)$$

Note that in the previous expression, the discrete and exact solutions are each associated with the corresponding definition of the energy. Owing to (24), the effective property estimated using the numerical schemes considered can become smaller than the exact one along iterations, which has already been observed in numerous studies. The impact on the development of error estimators is the impossibility to guarantee strict bounds in such a setting, as discussed later on.

In addition, we show in the remainder of this section that the iteration error estimate  $\Delta_k^h$  is not sufficient in itself and that a reliable discretization error estimate is also needed. To illustrate this and based on (24), we consider next a numerically consistent version  $d_{k,h}^{\text{eff}}$  of  $\delta_{k,h}^{\text{eff}}$  in (5) as

$$d_{k,h}^{\text{eff}} = |W_h(e_k^h) - W(e)|^{1/2}. \quad (25)$$

Therefore, the convergence of the four variants of the iterative FFT-based homogenization method are investigated numerically next on the set of test-cases considered and with the reference energy term  $W(e)$  in (25) being either known analytically or computed using the finite elements method on a reference fine grid, see Section 2.3.

### 3.2.2 Numerical examples

Let us first consider a discretization of  $2^6 \times 2^6$  pixels and a conductivity ratio  $\rho = 10$ . The error and residual convergence history for the random circular inclusions case are displayed on Figure 5. One observes that, for such a moderate contrast, the modified Green operator leads to a slightly faster convergence of the residual  $\Delta_k^h$ , but to a higher asymptotic error  $d_{k,h}^{\text{eff}}$ . This was already noticed in

[21]. As expected, the Conjugate Gradient, when associated with the modified operator or not, leads to a faster convergence, a fact that does not affect in the end the asymptotic error compared to the fixed-point scheme. This last result is due to the fact that the underlying discrete linear problem, and thus its solution, does not depend on the resolution method. With these results at hand, one considers the convergence criterion  $\Delta_k^h \leq 10^{-\beta} \Delta_0^h$ , with  $\beta = 5$  being chosen here as a typical value, and report the iteration number and the error reached when the threshold is met. These values are then reported on Figure 6. The same is done for the four schemes considered, on the two test-cases as well as for a high conductivity ratio  $\rho = 10^3$ .

At all points displayed on the plot, the residual  $\Delta_k^h$  is the same (except for the non-converged points). However, the error  $d_{k,h}^{\text{eff}}$  is different for each of them. This shows that the accuracy of the current iterate cannot be deduced from the residual alone.

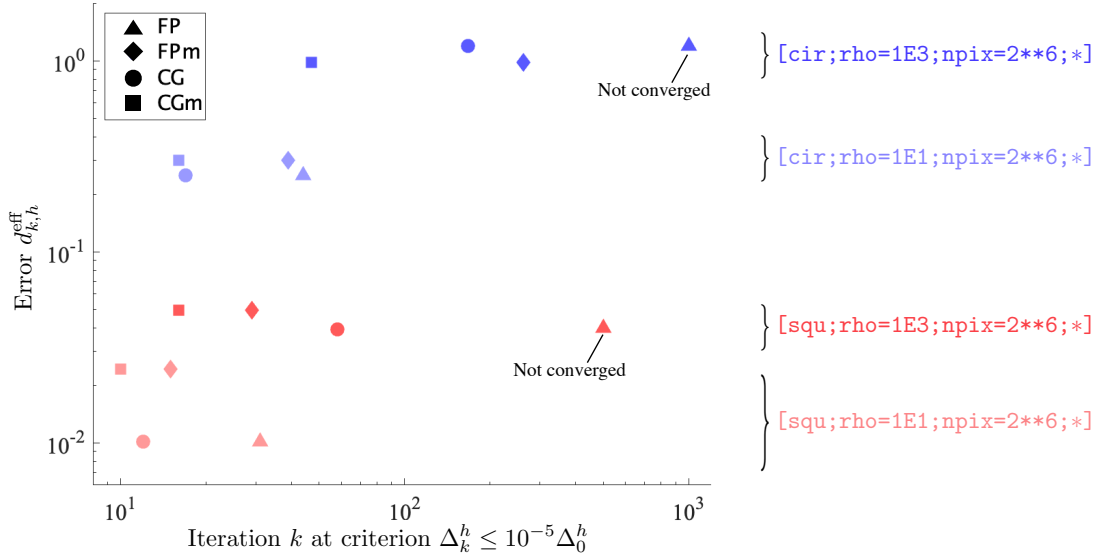


Figure 6: Error  $d_{k,h}^{\text{eff}}$  on the energy vs iteration number when the criterion  $\Delta_k^h \leq 10^{-5} \Delta_0^h$  is met, for the different schemes and test cases considered.

**Remark 6.** *One finds here a result compatible with what was obtained in [34] with a slightly different modification. The convergence of the residual is much faster when using the modified operator. The asymptotic error is better with the modified operator in the case of circular inclusions and high contrast, and worse in the square inclusion cases. This can be explained because the modification tends to filter out the singularities, which is favorable in the circular case, and not in the square case.*

### 3.3 Principle of a stopping criterion

Based on the previous discussion, we now expose the principle of a stopping criterion that would be based on estimations of both the iteration and the convergence errors. Considering (6), the triangle inequality entails  $\delta^h - \delta_k^h \leq \delta_{k,h}^{\text{tot}} \leq \delta_k^h + \delta^h$ , which in turn implies:

$$|\delta_{k,h}^{\text{tot}} - \delta^h| \leq \delta_k^h. \quad (26)$$

As a consequence, if the iteration error  $\delta_k^h$  becomes much smaller than the discretization error  $\delta^h$  then the total error  $\delta_{k,h}^{\text{tot}}$  cannot be expected to decrease noticeably any more. The bound (26) actually traduces our implicit assumption that, since the numerical schemes considered are expected to converge, the total error would tend to be mostly governed by the discretization error. This invites us to consider a stopping criterion relying on a comparison of  $\delta_k^h$  and  $\delta^h$  (or rather their estimations  $\Delta_k^h$  and  $\Delta^h$ ). Note that a similar idea has been used in the context of domain decomposition [24].

Let us introduce a parameter  $\beta > 0$  such that the iterations are to be stopped when

$$\delta_k^h \leq 10^{-\beta} \delta^h. \quad (27)$$

Let us now suppose that the error estimators  $\Delta_k^h$  and  $\Delta^h$  are such that there exist some positive parameters  $m_d^h$ ,  $M_d^h$ ,  $m_i^h$  and  $M_i^h$  such that the following bounds hold:

$$m_i^h \Delta_k^h \leq \delta_k^h \leq M_i^h \Delta_k^h \quad \text{and} \quad m_d^h \Delta^h \leq \delta^h \leq M_d^h \Delta^h.$$

Note that, owing to (21), one could choose  $m_i^h = (1+c)/2c$  and  $M_i^h = (1+c)/2$ . In this context, a sufficient condition ensuring (27) writes as the following uniform bound on  $k$ :

$$M_i^h \Delta_k^h \leq 10^{-\beta} m_d^h \Delta^h,$$

which readily provides a convergence criterion that will be discussed and illustrated numerically in Section 5.

## 4 Estimation of the discretization error

### 4.1 Preliminaries

At convergence when  $k \rightarrow \infty$  the obtained solution  $\tilde{e}^h = \nabla u^h$ , see (17), satisfies the discrete weak formulation (15) and our objective is now to compute an estimation  $\Delta^h$  to the corresponding discretization error  $\delta^h$  in (6) relatively to the exact solution  $\tilde{e} = \nabla u$ . As  $\delta^h$  involves both a discrete quantity and a continuous one, the evaluation of the featured continuous norm has to be consistent. As already discussed in Remark 2, computing this norm from the interpolation of  $\tilde{e}^h$  as a continuous Fourier series would yield an exaggeratedly high error due to the spurious oscillations of the former between sampling points. Therefore, we consider a piecewise-constant extension of the sampled version of  $\tilde{e}^h$  when appropriate, while keeping the same notation for simplicity.

The discretization error estimators developed in this section are inspired by methods that have proven their reliability in the framework of the Finite Element (FE) method. In this context and focusing on 2D problems without loss of generality, the cornerstone of the proposed approach is the interpretation of the DFT grid as a periodic and structured mesh of square Q1-Lagrange finite elements with reduced integration [9], i.e. elements having 4 nodes and 1 Gauss point at the center, see Figure 7a. In addition, we consider that the DFT interpolation points are at the centre of the pixels, which is only a display convention (note for example that it is not the one used by default in Matlab), and we identify these interpolation points as the Gauss points of the periodic FE mesh. The nodes of the FE mesh are thus at the corners of the pixels. This will allow a relatively simple transport of the discrete flux  $\mathbf{j}^h = \gamma^h(\bar{\mathbf{e}} + \tilde{\mathbf{e}}^h)$  computed at the Gauss points, and piecewise-constant over each Q1 element, to a FE version  $\mathbf{j}^h$  available at the nodes.

For a given Gauss point  $\mathbf{g}$ , we consider the surrounding element  $\mathbf{E}_g$  whose associated set of nodes is defined as  $\mathcal{N}(\mathbf{E}_g) = \{\mathbf{n}_j^g, j = 1, \dots, 4\}$ , see Fig. 7b. A point  $\mathbf{a} = (a_1, a_2)$  in the square reference unit element of Fig. 7a is mapped to a point  $\mathbf{x} \in \mathbf{E}_g$  using translations and dilatations of the former as

$$\mathbf{x}(\mathbf{a}) = \sum_{k=1}^4 N_k(\mathbf{a}) \mathbf{x}_k^g \quad \text{with} \quad \begin{cases} N_1(\mathbf{a}) = \frac{1}{4}(1-a_1)(1-a_2) \\ N_2(\mathbf{a}) = \frac{1}{4}(1+a_1)(1-a_2) \\ N_3(\mathbf{a}) = \frac{1}{4}(1+a_1)(1+a_2) \\ N_4(\mathbf{a}) = \frac{1}{4}(1-a_1)(1+a_2), \end{cases} \quad (28)$$

with the bilinear shape functions  $N_k$  in (28) being also used to define an isoparametric FE interpolation of a vectorial field  $\mathbf{f}^h$  at a point  $\mathbf{x} = \mathbf{x}(\mathbf{a}) \in \mathbf{E}_g$  as

$$\mathbf{f}^h(\mathbf{x}) = \sum_{k=1}^4 N_k(\mathbf{a}) \mathbf{f}^h(\mathbf{x}_k^g). \quad (29)$$

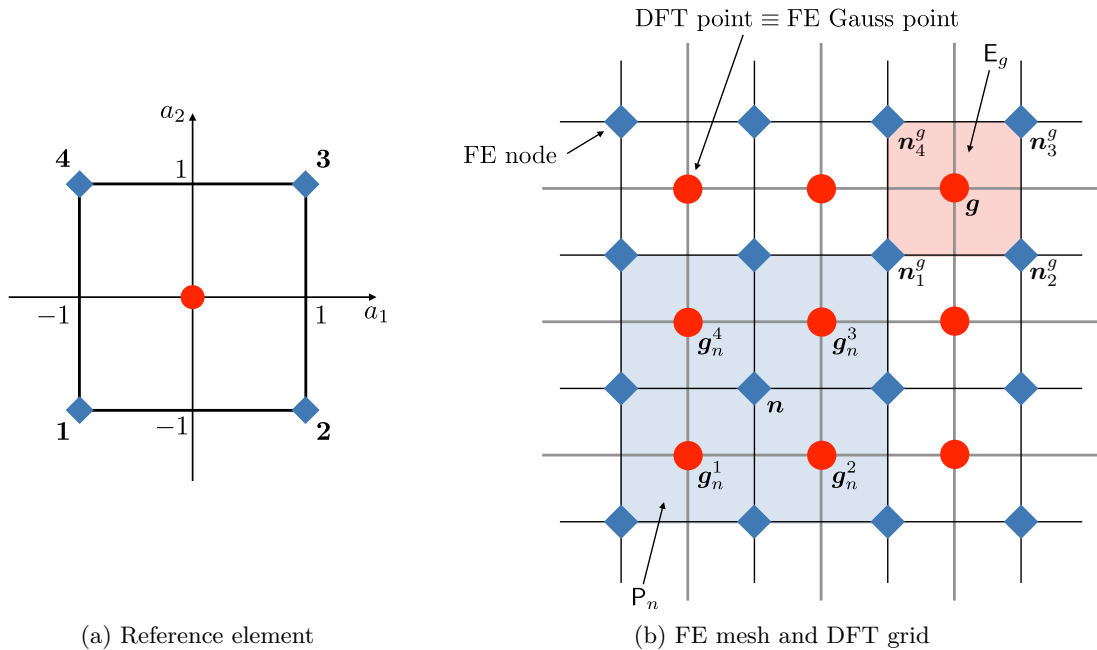


Figure 7: Nodes of the FE mesh (blue diamonds) and interpolation points of the DFT grid (red dots). A given DFT interpolation point  $\mathbf{g}$  at the center of a pixel is seen as the FE Gauss point of the element  $E_g$  that coincides geometrically with that pixel.

For further use, for each node  $\mathbf{n}$  we define the surrounding patch  $P_n$  as the set of 4 elements that share this node. The set of Gauss points that belong to this patch is given by  $\mathcal{G}(P_n) = \{\mathbf{g}_n^j, j = 1, \dots, 4\}$ , see Fig. 7b. Moreover, for each Gauss point  $\mathbf{g}$ , we also define  $A_g$  as the set of 9 adjacent elements (in 2D) that share at least one node with  $E_g$  (including  $E_g$  itself). Finally, the set of Gauss points that belong to these elements are denoted by  $\mathcal{G}(A_g)$ .

The common idea for the three proposed error estimation approaches is that from the available flux  $\mathbf{j}^h$ , that results from a Fourier-based computation of the solution to the discrete problem (16), one constructs a suitably modified flux  $\mathbf{j}_*^h$  and use the energetic distance  $\Delta^h = \|\mathbf{j}^h - \mathbf{j}_*^h\|_{h,\gamma^{-1}}$  as an estimation of the discretization error  $\delta^h$ . In all approaches we will make use of an intermediate FE flux quantity  $\mathbf{j}^h$  computed at the nodes following some principles that are described below. Note that the justification that makes the field  $\mathbf{j}_*^h$  relevant is not the same for the third method and the two others.

**Remark 7.** For the methods described hereafter, it has been chosen to construct the modified flux  $\mathbf{j}_*^h$  at a given Gauss point from the values of  $\mathbf{j}^h$  at the nearest neighbors only. In some cases, for example for highly refined discretizations, this window can be made larger and the proposed methods be easily extended to such configurations. Note also that periodicity must be taken into account for the computational treatment involving the Gauss points of the pixels at the boundary of the image.

**Remark 8.** The mean and median filters are commonly used in image processing. These two error estimators can be interpreted as follows: each component of the field  $\mathbf{j}^h$  is considered as a noisy image (due to the non-physical oscillations), and this image is smoothed out to obtain  $\mathbf{j}_*^h$ . Note that an image-filtering based approach has also been proposed in [19] to improve FFT-based computations in homogenization.

## 4.2 Weighted mean filter

In this method, we aim at constructing a modified flux  $\mathbf{j}_*^h$  that is aimed at being closer to the exact solution  $\mathbf{j} = \gamma(\bar{\mathbf{e}} + \tilde{\mathbf{e}})$  in an appropriate norm. As the Fourier-based discrete solution  $\mathbf{j}^h$  may typically exhibit unphysical oscillations, often seen as aliasing effects or Gibbs phenomena (see for example Fig. 13), we propose here to construct  $\mathbf{j}_*^h$  as a smooth version of  $\mathbf{j}^h$ . To do so, we draw from the

rather simple but seminal approach proposed in [37], where an improved version of the flux  $\mathbf{j}_*^h$  is computed by a nodal averaging process based on the FE shape functions (28–29).

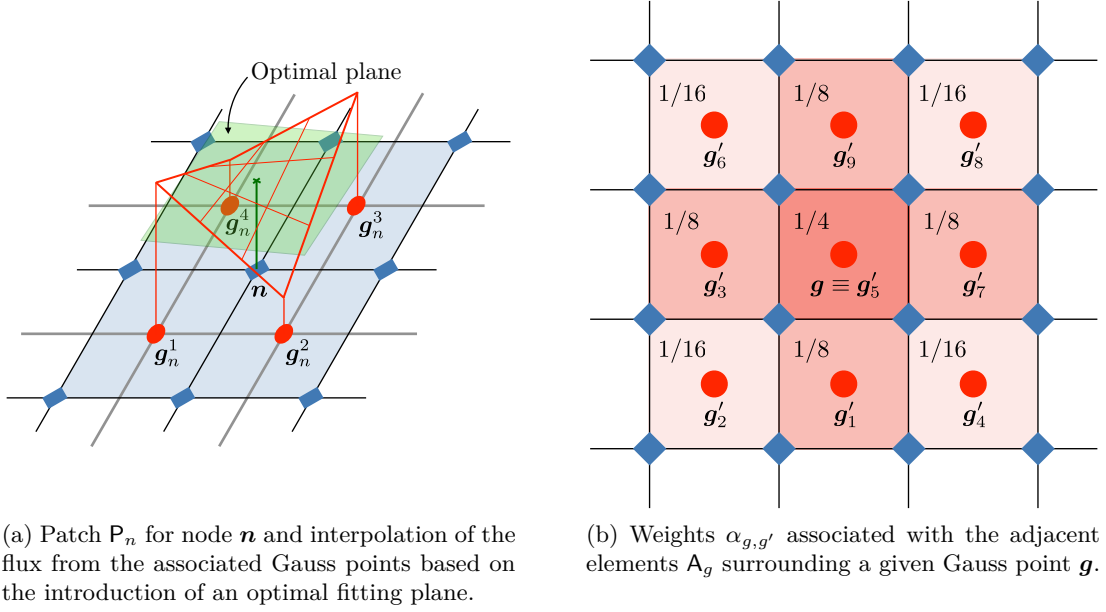


Figure 8: Weighted mean filter process.

This method has two steps. The first one consists in computing  $\mathbf{j}^h(\mathbf{n})$  for each FE node  $\mathbf{n}$  from the field values at the neighboring Gauss points, i.e.  $\mathbf{j}^h(\mathbf{g}_n^j)$  for  $\mathbf{g}_n^j \in \mathcal{G}(P_n)$ . In order to do this, we define  $\mathbf{j}^h(\mathbf{n})$  as a point on an optimal hyper-surface that minimizes the least-squares distances to the values at the chosen Gauss points. Here, this hyper-surface is simply defined as a plane, see Fig. 8a, which amounts to use the interpolating bilinear Lagrange polynomials (28) with the reference square element mapped to the Gauss points  $\mathbf{g}_n^j \in \mathcal{G}(P_n)$ . This polynomial is then evaluated at  $\mathbf{n}$ , the central point, so that we get

$$\mathbf{j}^h(\mathbf{n}) = \frac{1}{4} \sum_{\mathbf{g}_n^j \in \mathcal{G}(P_n)} \mathbf{j}^h(\mathbf{g}_n^j). \quad (30)$$

In the second step, the obtained values are then transported from the nodes back to the Gauss points by evaluation of the shape functions (28). Once again, the use of square elements implies that, for each element  $E_g$  the interpolated values at the central Gauss point  $\mathbf{g}$  is the mean among the four nodal values, i.e.

$$\mathbf{j}_*^h(\mathbf{g}) = \frac{1}{4} \sum_{\mathbf{n}_j^g \in \mathcal{N}(E_{g'})} \mathbf{j}^h(\mathbf{n}_j^g). \quad (31)$$

All in all, combining (30) and (31), leads to a resulting smoothed field computed at a given  $\mathbf{g}$  from a weighted mean over the values at the Gauss points  $\mathbf{g}'$  of the adjacent elements, see Fig. 8b, as

$$\mathbf{j}_*^h(\mathbf{g}) = \sum_{\mathbf{g}' \in \mathcal{G}(A_g)} \alpha_{g,g'} \mathbf{j}^h(\mathbf{g}') \quad (32)$$

where for any two Gauss point  $\mathbf{g}, \mathbf{g}'$ , we denote by  $\mu_{g,g'}$  the number of nodes shared by  $E_g$  and  $E_{g'}$  and define

$$\alpha_{g,g'} = \frac{\mu_{g,g'}}{16} \quad \text{with} \quad \sum_{\mathbf{g}' \in \mathcal{G}(A_g)} \alpha_{g,g'} = 1$$

### 4.3 Weighted median filter

Computing the weighted average (32) can be interpreted as solving for each Gauss point  $\mathbf{g}$ , i.e. for each pixel, the following local weighted quadratic minimization problem:

$$\mathbf{j}_*^h(\mathbf{g}) = \arg \min_{\mathbf{j}_* \in \mathbb{R}^2} \frac{1}{2} \sum_{\mathbf{g}' \in \mathcal{G}(A_g)} \alpha_{g,g'} \|\mathbf{j}^h(\mathbf{g}') - \mathbf{j}_*\|_2^2 \quad (33)$$



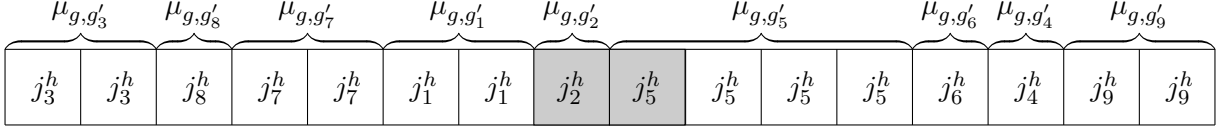


Figure 9: Sorting process for the computation of the weighted median filter (arbitrary example). Each component of  $j_\ell^h = j^h(\mathbf{g}'_\ell)_\zeta$  for  $\mathbf{g}'_\ell \in \mathcal{G}(A_g)$ , with  $\zeta = x$  and  $y$  and using the notations of Fig. 8b, are sorted in ascending order and one of the median values (greyed) is chosen as a minimizer of (34).

where the featured norm  $\|\cdot\|_2$  is the standard Euclidean norm. As it is known that the resulting field  $\mathbf{j}_*^h$  will not respect the possible discontinuities of the original flux, this process may be prone to overestimation of the error. This has been underlined in [27], where an alternative error estimation strategy is proposed, which however cannot be used in our context as it does not directly provides an estimation of  $\|\tilde{e} - \tilde{e}^h\|_\gamma$ . We rather consider here an alternate weighted median filter that is meant to allow for such possible discontinuities. As mentioned in Remark 8, this type of filtering is commonly employed in signal processing, see for example [8, 35], but it has never been used to our knowledge in the present context. To do so, the minimization problem (33) is modified using the  $L^1$ -norm and rewritten component-wise as

$$j_*^h(\mathbf{g})_\zeta = \arg \min_{j_* \in \mathbb{R}} \sum_{\mathbf{g}' \in \mathcal{G}(A_g)} \alpha_{g,g'} \left| j^h(\mathbf{g}')_\zeta - j_* \right| \quad \text{with} \quad \zeta = x \text{ or } y. \quad (34)$$

These minimization problems are then solved independently at each point  $\mathbf{g}$  and for each sought component ( $x$  and  $y$  in 2D) of the vector  $\mathbf{j}_*^h(\mathbf{g})$ . Each one can then be reformulated as two (or three in 3D) unidimensional minimization problems, see (34), for cost functionals that are linear by parts and whose minimum might not be unique. In that case, one of the minimizers is arbitrarily chosen.

Two methods can be used in practice to find a minimizer. The first one consists in evaluating each cost functional in (34) for each of the  $j^h(\mathbf{g}')_\zeta$ , with  $\zeta = x$  and  $y$ , i.e. at the *corners* of such piecewise linear functions, and storing the value for which each of them is minimal. The second one uses the fact that the minimizer in (34) is the weighted median of the set of vector values  $\{j^h(\mathbf{g}')_\zeta, \mathbf{g}' \in \mathcal{G}(A_g)\}$ . To compute the latter, one can sort the elements of this set in ascending order while duplicating them  $\mu_{g,g'}$  times, and pick the median value (here the 8th or 9th as  $\sum_{g'} \mu_{g,g'}$  is even). For the sake of the example an arbitrary sorted set is represented in Fig. 9 with reference to the notations of Fig. 8b. Note that, in the case where the filter window is chosen to be larger, see Remark 7 then fast algorithms can be used instead of the above sorting method [8].

## 4.4 Constitutive Relation Error

The third and last method that will be tested on our problem is the Constitutive Relation Error method (CRE) [13]. This method builds an estimation of the error that, provided a few hypotheses are verified, is proven to be higher than the true error. Its usual limitation is that it is more computationally costly than competing methods. However, as it is shown in this section, in the case where all the elements of the mesh are identical, the cost of the method is similar as the one of the ZZ1 method.

### 4.4.1 Principle of the method

Let us use the field  $\tilde{\mathbf{e}}_*^h$  defined in Section 3.2.1 as the solution to the discrete Lippmann-Schwinger problem with exact integration, with associated total field  $\mathbf{e}_*^h = \bar{\mathbf{e}} + \tilde{\mathbf{e}}_*^h$ . We now aim at estimating the discretization error  $\delta^h = \|\bar{\mathbf{e}} - \tilde{\mathbf{e}}_*^h\|_\gamma = \|\gamma(\mathbf{e} - \mathbf{e}_*^h)\|_{\gamma^{-1}}$ .

For all  $\mathbf{j}_*^h \in \mathcal{S}$  we have:

$$\begin{aligned} \|\mathbf{j}_*^h - \gamma \mathbf{e}_*^h\|_{\gamma^{-1}}^2 &= \|\mathbf{j}_*^h - \gamma \mathbf{e} + \gamma(\mathbf{e} - \mathbf{e}_*^h)\|_{\gamma^{-1}}^2 \\ &= \|\mathbf{j}_*^h - \gamma \mathbf{e}\|_{\gamma^{-1}}^2 + \|\gamma(\mathbf{e} - \mathbf{e}_*^h)\|_{\gamma^{-1}}^2 + 2(\mathbf{j}_*^h - \gamma \mathbf{e}, \gamma(\mathbf{e} - \mathbf{e}_*^h))_{\gamma^{-1}}. \end{aligned}$$

Note that it is the exact conductivity  $\gamma$  that enters the above identities consistently with Remark 3. As  $(\mathbf{e} - \mathbf{e}_\star^h) \in \mathcal{E}_0$  and  $(\mathbf{j}_\star^h - \gamma \mathbf{e}) \in \mathcal{S}$ , and since these functional spaces are orthogonal in the sense of the standard scalar product on  $L^2_{\text{per}}(\Omega)$  (see Section 2.1), then one has:

$$(\mathbf{j}_\star^h - \gamma \mathbf{e}, \gamma(\mathbf{e} - \mathbf{e}_\star^h))_{\gamma^{-1}} = (\mathbf{j}_\star^h - \gamma \mathbf{e}, \mathbf{e} - \mathbf{e}_\star^h) = 0, \quad (35)$$

and then

$$\|\mathbf{j}_\star^h - \gamma \mathbf{e}_\star^h\|_{\gamma^{-1}}^2 = \|\mathbf{j}_\star^h - \gamma \mathbf{e}\|_{\gamma^{-1}}^2 + \|\gamma(\mathbf{e} - \mathbf{e}_\star^h)\|_{\gamma^{-1}}^2. \quad (36)$$

As a conclusion, for any  $\mathbf{j}_\star^h \in \mathcal{S}$ , the quantity  $\|\mathbf{j}_\star^h - \gamma \mathbf{e}_\star^h\|_{\gamma^{-1}}$  is an estimator of the discretization error  $\delta^h = \|\gamma(\mathbf{e} - \mathbf{e}_\star^h)\|_{\gamma^{-1}}$ , that is all the more accurate as  $\|\mathbf{j}_\star^h - \gamma \mathbf{e}\|_{\gamma^{-1}}$  is small, i.e. when  $\mathbf{j}_\star^h$  gets close to the exact flux  $\mathbf{j} = \gamma \mathbf{e}$ . What is more, this estimator is guaranteed to be an upper bound of the true discretization error.

We recall however that, due to inexact integration, the orthogonality relations (35) are not satisfied despite the fact that the fields  $\mathbf{e}_k^h$  computed are compatible, see Remark 5. In other words, the inexact integration amounts to prolongating  $\mathbf{e}^h$  by a piecewise-constant field that will not coincide with  $\mathbf{e}_\star^h$  (see Remark 2). As a result, we reach a conclusion similar to (24), in that the above bound on the discretization error cannot be ensured for  $\mathbf{e}^h$ . Yet, we will make use of the derivation above as a guideline to devise an estimator that we discuss next, i.e. we will apply the CRE procedure, designed to estimate  $\|\tilde{\mathbf{e}} - \tilde{\mathbf{e}}^h\|_\gamma$ , to construct an estimate of  $\|\tilde{\mathbf{e}} - \tilde{\mathbf{e}}^h\|_\gamma$ . The relevance of this approach will be assessed on a number of numerical examples.

The flux-equilibration method [15] will be used to compute a modified field  $\mathbf{j}_\star^h$  from the available numerical solution. It consists in two steps:

- Firstly, determine an equilibrated normal flux  $F_\star^h$  at the boundary of each element  $\mathbf{E}_g$  of the F.E. mesh from the original flux  $\mathbf{j}^h = \gamma^h \mathbf{e}^h$  available at the Gauss points.
- Secondly, compute  $\mathbf{j}_\star^h$  from a F.E. approximation on a refined mesh of the solution  $\mathbf{j}_\star$  to the following continuous subproblem:

$$\begin{cases} \mathbf{j}_\star(\mathbf{x}) = \gamma^g \nabla \mathbf{u}_\star(\mathbf{x}) & \mathbf{x} \in \mathbf{E}_g \\ \nabla \cdot \mathbf{j}_\star(\mathbf{x}) = 0 & \mathbf{x} \in \mathbf{E}_g \\ \mathbf{j}_\star(\mathbf{x}) \cdot \boldsymbol{\nu}_g(\mathbf{x}) = F_\star^h(\mathbf{x}) & \mathbf{x} \in \partial \mathbf{E}_g \end{cases} \quad (37)$$

where  $\gamma^g$  is the *homogeneous* conductivity value sampled at the Gauss point  $\mathbf{g}$  and  $\boldsymbol{\nu}_g$  is the unit outward normal on  $\partial \mathbf{E}_g$ . The fact that the constant value  $\gamma^g$  of  $\gamma^h$  in  $\mathbf{E}_g$  is used in (37) pertains to the assumption that the latter constitutes the exact conductivity, see remarks 1 and 3 (the geometrical error is not accounted for in the present work). In addition, this choice allows the factorization by the material property, which in turn makes the resolution of (37) numerically efficient. If an exact conductivity field  $\gamma(\mathbf{x})$ , different from its discrete counterpart  $\gamma^h$ , were known then it should be used in (37) to get a consistent *non-homogeneous* material distribution in the sub-discretization of  $\mathbf{E}_g$ .

Doing so, the field  $\mathbf{j}_\star^h$  respects the discrete equilibrium equation on the whole domain, and its proximity to  $\gamma \mathbf{e}$  depends on the relevance of the boundary fluxes.

#### 4.4.2 Computation of the boundary fluxes

Given a node  $\mathbf{n}$ , together with the associated patch  $\mathbf{P}_n$ , let  $\phi_n$  denote the associated global shape function constructed from (28). Considering the inter-element interface  $I$  connecting  $\mathbf{n}$  to another node  $\mathbf{n}'$ , see Fig. 10a, the objective is then to compute an equilibrated boundary flux  $F_\star^h(\mathbf{x})$ , spatially varying along  $I$  and which must be consistent with the available values of the flux  $\mathbf{j}^h(\mathbf{g}_n^j)$  at the Gauss points  $\mathbf{g}_n^j \in \mathcal{G}(\mathbf{P}_n)$ , see Fig. 10b.

To do so, one considers the prolongation equation [13]. It amounts in equating the virtual works associated with  $\mathbf{j}^h$  and  $\mathbf{j}_\star$  in any element  $\mathbf{E} \in \mathbf{P}_n$  and for all test function  $\phi_m$  relative to a node  $\mathbf{m}$  of this element, i.e. we set

$$\begin{aligned} \int_{\mathbf{E}} \mathbf{j}^h(\mathbf{x}) \cdot \nabla \phi_m(\mathbf{x}) \, \mathrm{d}\mathbf{x} &= \int_{\mathbf{E}} \mathbf{j}_\star(\mathbf{x}) \cdot \nabla \phi_m(\mathbf{x}) \, \mathrm{d}\mathbf{x} \\ &= \int_{\partial \mathbf{E}} F_\star^h(\mathbf{x}) \phi_m(\mathbf{x}) \, \mathrm{d}\mathbf{x} \end{aligned} \quad (38)$$

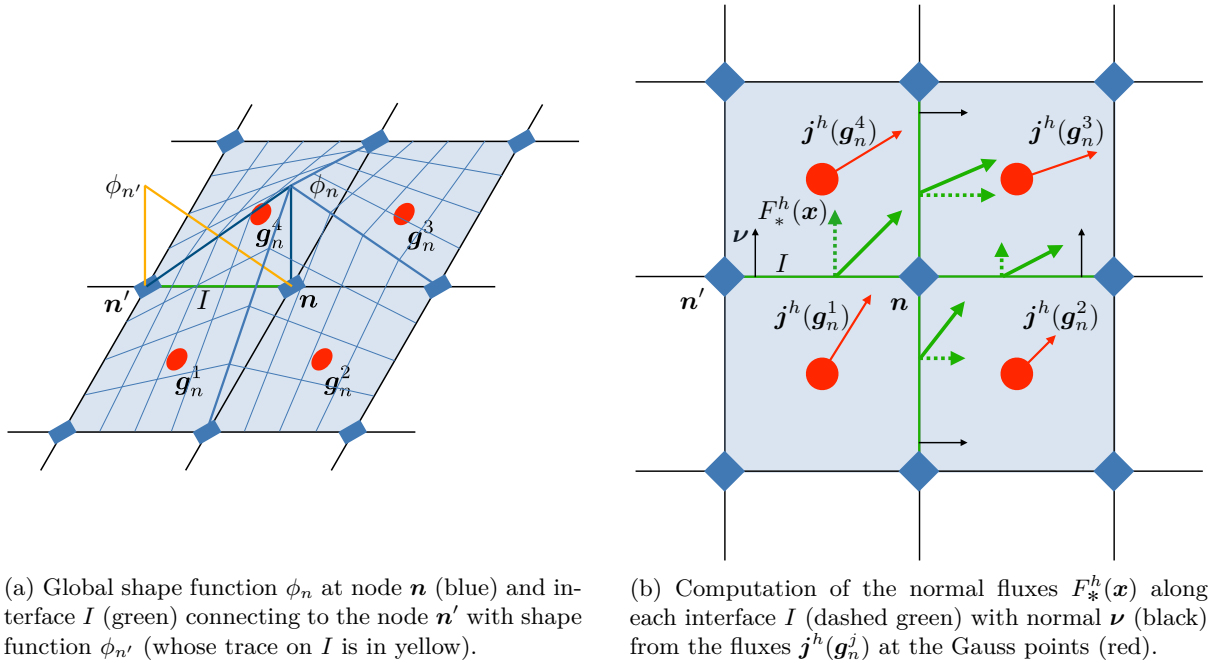


Figure 10: Representation and computation of the boundary fluxes.

where the second equality makes use of (37). Since the set of shape functions  $\{\phi_n\}_n$  satisfies the property of partition of unity  $\sum_n \phi_n(\mathbf{x}) = 1$ , then (38) ensures that  $F_*^h(\mathbf{x})$  is equilibrated over  $\mathbf{E}$ , i.e. it holds

$$\int_{\partial\mathbf{E}} F_*^h(\mathbf{x}) d\mathbf{x} = 0.$$

Now, at a given boundary  $I$  of an element  $\mathbf{E}$ , to enforce the continuity of the flux across elements we write  $F_*^h(\mathbf{x}) = \eta_{\mathbf{E}}^I F_I^h(\mathbf{x})$  with the factor  $\eta_{\mathbf{E}}^I = \pm 1$  being chosen so as to ensure that an outgoing flux is the opposite of an entering flux on each adjacent elements. This is done consistently with the orientation of the normal  $\boldsymbol{\nu}$  in Figure 10b.

In addition, consistently with the discretization (28) considered, we define  $F_I^h(\mathbf{x})$  as a linear function of  $\mathbf{x} \in I$ , so that it can be associated with two scalar degrees of freedom denoted as  $F_{I,n}^h$  and  $F_{I,n'}^h$ , and which we defined as the following projection

$$F_{I,n}^h = \int_I F_I^h(\mathbf{x}) \phi_n(\mathbf{x}) d\mathbf{x}, \quad (39)$$

and likewise for  $F_{I,n'}^h$  using  $\phi_{n'}$ . As a consequence for the mesh considered, introducing the set  $\mathbf{B}_n$  of four interfaces between the elements of  $\mathbf{P}_n$ , there are four degrees of freedom associated to each node  $\mathbf{n}$ , which we gather into the following vector:

$$\mathbf{F}_n^h = \left\{ F_{I,n}^h, I \in \mathbf{B}_n \right\} \in \mathbb{R}^4.$$

Therefore, owing to (38) and upon choosing the test function relative to the node  $\mathbf{m} = \mathbf{n}$  we get

$$\sum_{I \in \mathbf{B}_n} \eta_{\mathbf{E}}^I F_{I,n}^h = \int_{\mathbf{E}} \mathbf{j}^h(\mathbf{x}) \cdot \nabla \phi_n(\mathbf{x}) d\mathbf{x}, \quad (40)$$

Given that the available flux  $\mathbf{j}^h$  is piecewise constant in each of the four elements  $\mathbf{E} \in \mathbf{P}_n$  (see Section 4.1), then (40) leads to the following linear system:

$$\mathbf{A} \mathbf{F}_n^h = \mathbf{b}_n$$

with 
$$\mathbf{A} = \begin{pmatrix} -1 & 0 & 0 & -1 \\ 1 & -1 & 0 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 0 & -1 & 1 \end{pmatrix} \quad \text{and} \quad \mathbf{b}_n = \begin{pmatrix} -j^h(\mathbf{g}_n^1)_x - j^h(\mathbf{g}_n^1)_y \\ j^h(\mathbf{g}_n^2)_x - j^h(\mathbf{g}_n^2)_y \\ j^h(\mathbf{g}_n^3)_x + j^h(\mathbf{g}_n^3)_y \\ -j^h(\mathbf{g}_n^4)_x + j^h(\mathbf{g}_n^4)_y \end{pmatrix}. \quad (41)$$

However, the matrix  $\mathbf{A}$  is singular and its null space is of dimension 1. To overcome this, we use its Moore-Penrose pseudo-inverse [1], denoted by  $\mathbf{A}^+$ . Moreover, as proposed in [15], the extra degree of freedom (in the null space of  $\mathbf{A}$ ) is determined by minimizing the distance between the components of  $\mathbf{F}_n^h$  and these of an average vector  $\langle \mathbf{F}_n^h \rangle$  defined on each interface  $I$  from the mean value of the projections onto the normal  $\boldsymbol{\nu}$  of the adjacent fluxes, i.e.  $\frac{1}{2}(\mathbf{j}^h(\mathbf{g}_n^1) + \mathbf{j}^h(\mathbf{g}_n^4)) \cdot \boldsymbol{\nu}$  for example. All in all,  $\mathbf{F}_n^h$  is computed from the available components of the flux by solving the following equation:

$$\mathbf{F}_n^h = (\mathbf{A}^+ \mathbf{B}_1 + \mathbf{k} \mathbf{k}^T \mathbf{B}_2) \mathbf{j}_n^h$$

with  $\mathbf{j}_n^h = (j^h(\mathbf{g}_n^1)_x \ j^h(\mathbf{g}_n^1)_y \ j^h(\mathbf{g}_n^2)_x \ j^h(\mathbf{g}_n^2)_y \ j^h(\mathbf{g}_n^3)_x \ j^h(\mathbf{g}_n^3)_y \ j^h(\mathbf{g}_n^4)_x \ j^h(\mathbf{g}_n^4)_y)^T$ ,

$$\text{span}(\mathbf{k}) = \text{null}(\mathbf{A}), \quad \mathbf{B}_1 = \begin{pmatrix} -1 & -1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & -1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & -1 & 1 \end{pmatrix}, \quad \mathbf{B}_2 = \frac{1}{2} \begin{pmatrix} 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}.$$

The matrices  $\mathbf{B}_1$  and  $\mathbf{B}_2$  build respectively  $\mathbf{b}_n$  and the components of  $\langle \mathbf{F}_n^h \rangle$  from  $\mathbf{j}_n^h$ . As the geometries of the element patches are all the same, the matrices  $\mathbf{A}$ ,  $\mathbf{B}_1$ ,  $\mathbf{B}_2$  and vector  $\mathbf{k}$  do not depend on the specific node  $\mathbf{n}$ . As a consequence, the matrix operator  $(\mathbf{A}^+ \mathbf{B}_1 + \mathbf{k} \mathbf{k}^T \mathbf{B}_2)$  is assembled once for all in a pre-processing step and it is then applied to a multi-vector constituted from the set  $\{\mathbf{j}_n^h\}_n$ , which is indexed by the nodes  $\mathbf{n}$  of the given mesh. The boundary fluxes  $\{\mathbf{F}_n^h\}_n$  are then available after such a computation.

**Remark 9.** *In usual applications of this approach in the context of finite elements [13, 15], one can show that  $\mathbf{b}_n$  is orthogonal to  $\mathbf{k}$ . As a result, one can write the problem as the minimization of  $\|\mathbf{F}_n^h - \langle \mathbf{F}_n^h \rangle\|$  under the constraint of Equation (41). It can be shown that the idea proposed here is mathematically equivalent to solving this minimization problem (by remarking that (41) only imposes the projection of  $\mathbf{F}_n^h$  in the subspace orthogonal to  $\mathbf{k}$ ), with the advantage of being also usable when there is no solution to (41).*

#### 4.4.3 Computation of a self-equilibrated flux

For each Gauss point  $\mathbf{g}$ , we now aim at computing a self-equilibrated flux  $\mathbf{j}_*^h$  for the associated element (or pixel)  $E_g$  that respects the boundary fluxes computed at the previous step. We introduce the vector  $\mathbf{F}_g^h = \{(F_g^h)_{jk}, j = 1, \dots, 4, k = 1, 2\} \in \mathbb{R}^8$  that gathers the relevant degrees of freedom in  $\mathbf{F}_{n_j^g}^h$  for all nodes  $\mathbf{n}_j^g \in \mathcal{N}(E_g)$ , i.e. two associated with each element edge, see Fig. 11.

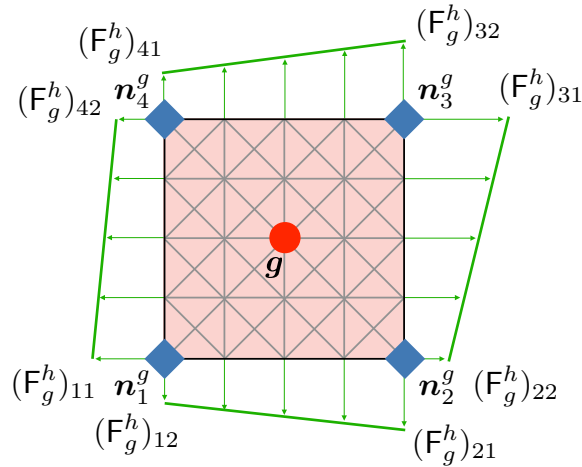


Figure 11: Finite element-based computation of a self-equilibrated flux  $\mathbf{j}_*^h$  on a given element  $E_g$  from the associated boundary fluxes computed previously and using a refined mesh.

Considering the continuous interpolation operator  $\mathcal{M}$  that yields piecewise linear normal fluxes  $\mathbf{F}_*^h(\mathbf{x})$  from its projected components  $\mathbf{F}_g^h$ , see (39), we consider the solution  $\mathbf{j}_*$  to (37) with the

boundary condition

$$\mathbf{j}_*(\mathbf{x}) \cdot \boldsymbol{\nu}_g(\mathbf{x}) = \mathcal{M}(\mathbf{F}_g^h)(\mathbf{x}) \quad \forall \mathbf{x} \in \partial E_g.$$

The Neumann problem (37) is then discretized and solved using the finite element method on a mesh that is refined compared to the original pixel-based discretization, and using a primal formulation. Noticeably, a  $p$ -refinement strategy could be used at this step but  $h$ -refinement is simple and efficient here since one only needs to build a single mesh as all the elements we are dealing with are identical. In addition,  $h$ -refinement is preferred here as it would allow a direct mapping to refined FFT grids if a multi-grid approach were used to improve the computations. This leads to an approximated flux that is considered to be *sufficiently self-equilibrated* provided the mesh of this subproblem is fine enough [15]. Given  $h' < h$  the associated mesh size, we denote by  $\mathbf{K}^{h'}$  the corresponding *stiffness* matrix and by  $\mathbf{M}^{h'}$  the discretized version of the operator  $\mathcal{M}$  and define  $\mathbf{u}_*^{h'}$  as the solution to

$$\mathbf{K}^{h'} \mathbf{u}_*^{h'} = \mathbf{M}^{h'} \mathbf{F}_g^h. \quad (42)$$

Due to the absence of Dirichlet boundary conditions, the matrix  $\mathbf{K}^{h'}$  is singular and its null space corresponds to the fields that are uniform in  $E_g$ . Therefore, we make use of its Moore-Penrose pseudo-inverse  $\mathbf{K}^{h'+}$ . At this stage a self-equilibrated flux is available on a refined mesh.

Finally, in order to compute the energy of the difference between  $\mathbf{j}_*^{h'}$  and the original field  $\mathbf{j}^h$ , we consider a discrete operator  $\mathbf{G}$  that computes a final modified flux  $\mathbf{j}_*^h$  at the original Gauss point  $\mathbf{g}$  from the refined finite element solution  $\mathbf{u}_*^{h'}$ , i.e.  $\mathbf{j}_*^h(\mathbf{g}) = \mathbf{G} \mathbf{u}_*^{h'}$ . Different strategies can be adopted to do so and the operator chosen here simply amounts in averaging locally the fluxes computed at the Gauss points of the finite elements of  $E_g$  that are adjacent to  $\mathbf{g}$ .

As the number of degrees of freedom involved in (42) is relatively small and since all pixels can be discretized using the same refined mesh, then the finite element system (42) can be inverted in a pre-processing step and its inverse stored in matrix form. Therefore, the operator that gives the value  $\mathbf{j}_*^h(\mathbf{g})$  at the Gauss point from the components of the boundary flux  $\mathbf{F}_g^h$  reads

$$\mathbf{j}_*^h(\mathbf{g}) = \mathbf{G} \mathbf{K}^{h'+} \mathbf{M}^{h'} \mathbf{F}_g^h.$$

Again, the key-point of the proposed approach is that, all pixels having the same geometry, the matrices  $\mathbf{G}$  and  $\mathbf{K}^{h'}$  only depend on  $\mathbf{g}$  through a pre-factor  $\gamma^g$ . Therefore,  $(\mathbf{G} \mathbf{K}^{h'+})$  and  $\mathbf{M}^{h'}$  do not depend on  $\mathbf{g}$  and the computation above can be directly performed on a multi-vector  $\{\mathbf{F}_g^h\}_g$  that is a mere reordering of the nodal degrees of freedom  $\{\mathbf{F}_n^h\}_n$  of the boundary fluxes. This is slightly different in the case of elasticity, see Remark 12.

## 4.5 Error maps

The three proposed smoothing methods are now illustrated numerically on the two test cases. First, the Obnosov square inclusion geometry with a conductivity ratio  $\rho = 10$  is considered and a numerical solution  $\mathbf{j}^h$  is computed on a coarse discretization of  $2^6 \times 2^6$  pixels using  $k = 40$  iterations of the method FP. Then we show on Figure 12 the different modified fields  $\mathbf{j}_*^h$  computed using the mean, median and CRE methods, which can thus be compared to the field computed originally and to a reference field. The latter corresponds to a computation on a finer grid of  $2^{10} \times 2^{10}$  pixels and averaged locally to be shown on the coarse grid.

One can observe that the mean smoothing method does not respect the discontinuities of the flux at the top and bottom boundaries of the inclusion. The CRE and the median methods seem to allow for these discontinuities. In addition the median smoothing does not preserve the corner singularities on the square inclusion. Figure 14a presents the corresponding maps of the local estimated error  $(\mathbf{j}^h - \mathbf{j}_*^h)/\sqrt{\gamma}$ . It appears that the over-smoothing done by the mean method leads to a large over-estimation of the error in the discontinuity regions. The median method over-estimates the error nearby the singularities while, on this test-case, the CRE method performs at best. On Figure 13, we present the reconstructed fluxes for the random circles test-case with the local errors being displayed on Figure 14b.

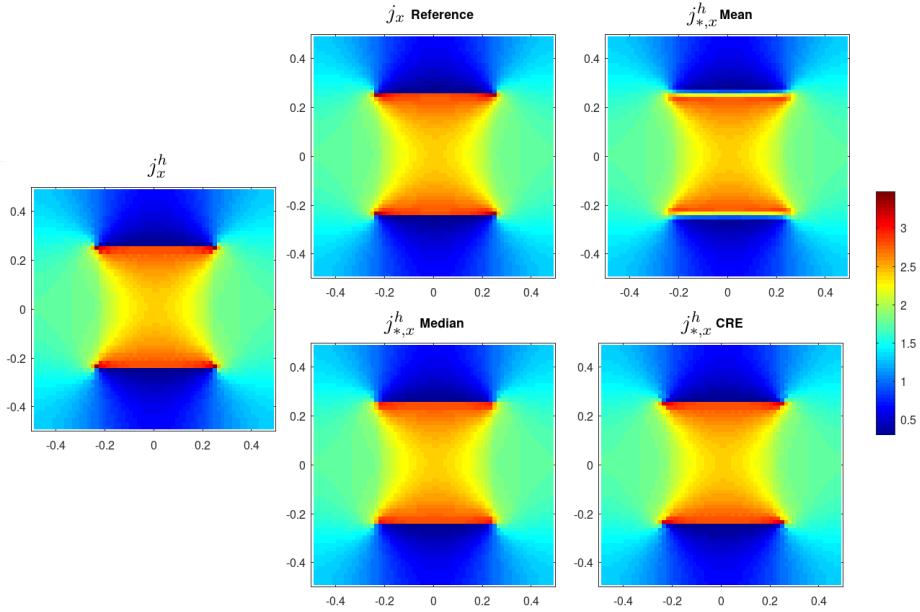


Figure 12: [squ;rho=1E1;npix=2\*\*6;FP]: Discrete flux  $j_x^h$ , reference flux  $j_x$  and reconstructed fluxes  $j_{*,x}^h$  with the different methods.

#### 4.6 Performance of the global error estimate

In this section, we investigate the accuracy of the three different error estimators on 16 different settings. They are tested on both the square and random circles inclusions geometries, with both  $2^6 \times 2^6$  and  $2^8 \times 2^8$  pixels, with both the modified and original operators, and with both  $\rho = 10$  and  $\rho = 10^3$ . The fixed-point algorithm was used but the specific choice of solver does not matter here as only the error at convergence is studied (a number  $k = 40$  of iterations were used when  $\rho = 10$ , and  $k = 500$  when  $\rho = 10^3$ ). The results are displayed on Figure 15 for all combinations of the numerical parameter that lead to a true error  $\delta^h \leq 1$ . Therefore, in these figures we have discarded the points for which the true error is above 1, which are here obtained in some of the high-contrast cases. In such cases our estimations are not relevant.

The mean filter based error estimator gives mostly always the worst estimation. It must be noticed however that this estimator is the easiest to develop, and is marginally less CPU costly than the CRE one. When the original Green's operator is used (Figure 15a), the CRE method appears to lead to the most reliable estimator in the cases of low contrast (although not a guaranteed bound, see Section 4.4.1). However, in the cases of high contrast, this estimator does not perform much better than the mean filter. Finally, the median filter-based method appears to lead to the best estimator in the high contrasts cases, while being acceptable in the low contrast cases considered.

When the modified operator is used (Figure 15b), the median filter and CRE-based error estimators turn out to have quite equivalent accuracies, and perform reasonably well on the test cases considered. It must be noticed that, for the investigated configurations, these estimators tend to yield values lower than the exact one, which is often considered to be undesirable in error estimation. Yet, the theoretical bound of Section 4.4.1 has no reason to be satisfied when using the modified operator. In addition and as previously discussed, because of inexact integration, the schemes considered violate the identity (36), which conventionally yields a guaranteed upper bound with the CRE method. This explains that we can get estimated errors below the true ones even with this method. The use of exact numerical integration would allow to recover guaranteed bounds.

## 5 Convergence criterion: numerical results

In this section, we assess the relevance of a stopping criterion based on the comparison between the estimations of the iteration error and of the discretization one, as discussed in Section 3.3. Here the

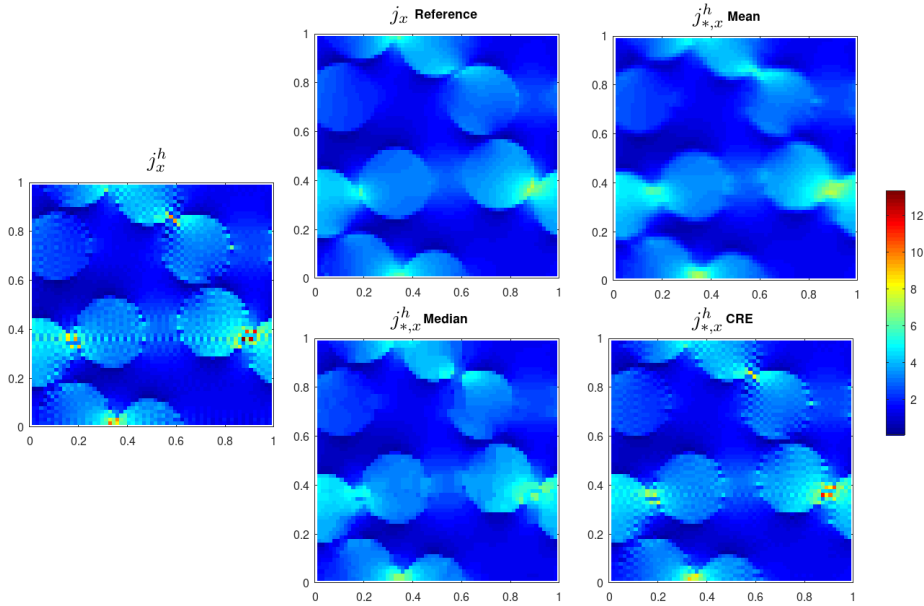


Figure 13: [cir;rho=1E1;npix=2\*\*6;FP]: Discrete flux  $j_x^h$ , reference flux  $j_x$  and reconstructed fluxes  $j_{*,x}^h$  with the different methods.

discretization is considered to be given and fixed.

Two convergence criteria, denoted by  $R_\beta$  and  $E_\beta^{\text{method}}$  are considered:

- (i)  $R_\beta$  consists in stopping the algorithm at the first iteration for which  $\Delta_k^h \leq 10^{-\beta} \Delta_0^h$  where  $\Delta_0^h$  is the initial residual. This criterion corresponds to the one most employed in the literature.
- (ii)  $E_\beta^{\text{method}}$  consists in stopping the algorithm at the first iteration for which  $\Delta_k^h \leq 10^{-\beta} \Delta^h$ . In other words, it consists in waiting until the estimated iteration error has no more significant effect on the global error by being a threshold lower than the estimated discretization error. This is similar to the criterion employed in [10], which couples estimations of the discretization error and of the algebraic error associated with an inexact solution of the linear system in a finite volume method. Here, criterion  $E_\beta^{\text{med}}$  uses the median filter to estimate  $\Delta^h$ , while  $E_\beta^{\text{cre}}$  uses the CRE method. Noticeably, we will here make use of the approaches described in Section 4 to compute some estimations  $\Delta^h$  of the discretization error, not in the limit  $e^h = \lim_{k \rightarrow \infty} e_k^h$  but for the specific values of  $k$  considered. In practice, it was noticed that at the very first iterations, the estimator  $\Delta^h$  was much smaller than the residual  $\Delta_k^h$ , which prevented the stopping criterion to be met accidentally due to the use of  $e_k^h$  instead of  $e^h$  in the computation of  $\Delta^h$ .

Once the chosen stopping criterion is met then the computation is simply stopped. At this point, the numerical results could be further improved by performing a new computation on a refined grid and using the last computation as an initial guess, see e.g. [6]. This is however beyond the scope of this work.

## 5.1 2D conductivity test-cases

In the case of the square inclusion, Figure 16 plots the *true* error  $\delta_{k,h}^{\text{eff}}$  on the effective conductivity and the associated number of iterations. Here,  $\delta_{k,h}^{\text{eff}}$  is computed from the knowledge of the analytical effective conductivity [22]. The conductivity ratio is  $\rho = 10^3$ , and we use the variant of the method denoted by CG because of its faster convergence. Different discretizations are investigated.

On this test-case, the criteria  $E_2^{\text{med}}$  and  $E_2^{\text{cre}}$  yields similar numbers of iterations. However, as for a high contrast the CRE method over-estimated the discretization error (see Section 4.6), the corresponding stopping criterion requires less iterations in that case. The proposed criteria give errors that are relatively close to the standard criterion  $R_5$  but they require different numbers of iterations.

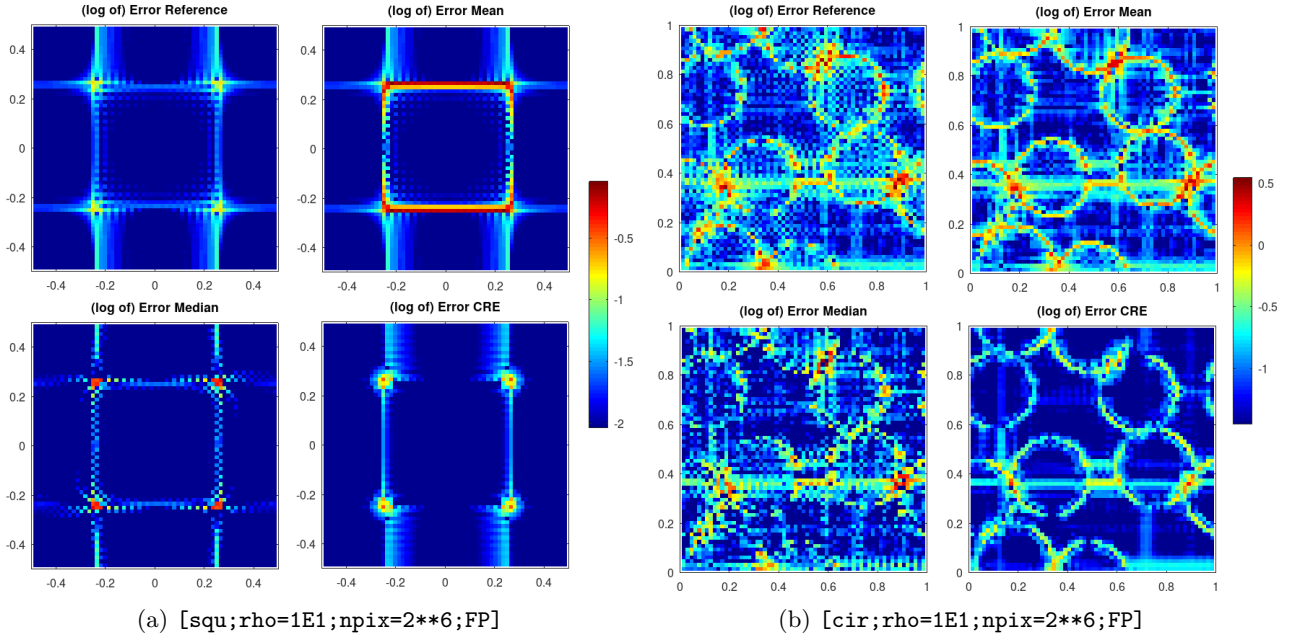


Figure 14: Reference and error estimators (logarithmic plot with scaling by the conductivity).

As a consequence, we introduce a new indicator, that we will refer to as *indicator of improvement*, in order to understand which stopping criterion is the most efficient. This indicator, denoted by  $\mathcal{I}_k^h$ , aims at measuring how much the error on the effective property could have been reduced if the iterative procedure had not been stopped. It is defined as the ratio between the iteration error and discretization error but uses the reference effective property and a numerical solution  $\mathbf{e}_\infty^h$  obtained after a sufficiently large number of iterations as

$$\mathcal{I}_k^h = \frac{|W_h(\mathbf{e}_k^h) - W_h(\mathbf{e}_\infty^h)|^{1/2}}{|W(\mathbf{e}) - W_h(\mathbf{e}_\infty^h)|^{1/2}}.$$

A *good* criterion would stop at an iteration  $k$  that ensures that the error on the effective property (ie. on the energy) cannot be decreased noticeably anymore, disregarding the discretization, geometry or conductivity ratio. This can be ensured when  $\mathcal{I}_k^h = \alpha$ , with  $\alpha \ll 1$  being a parameter that tunes how small  $\mathcal{I}_k^h$  is wanted to be, and that does not vary with the discretization, geometry nor conductivity ratio. Typically, we choose  $\alpha = 0.1$

**Remark 10.** Criterion  $E_\beta^{\text{method}}$  is based on the discretization error estimators developed in Section 4, that estimate the energetic error on the field  $\mathbf{e}$ , defined in Equation (6), rather than the error on the effective property, or equivalently on the energy which is used to evaluate  $\mathcal{I}_k^h$ . According to (23), if exact integration were used, both errors would be equal.

On Figure 17a, we plot the evolution of the indicator of improvement  $\mathcal{I}_k^h$  for the different stopping criteria considered, as functions of the discretization. One can notice that the criteria of type  $E_\beta^{\text{method}}$ , which are based on the proposed error estimators, tend to make  $\mathcal{I}_k^h$  independent to the discretization, while the conventional residual-based criterion  $R_5$  makes it increase with the discretization. This means that, for a coarse grid, the latter requires too many iterations, while for a finer grid, it does not allow to benefit from the full precision offered by the discretization.

In the case of random circular inclusions, the trends on the error and number of iterations are similar to the case of a square inclusion (see Figure 16) but they are not shown here. However, on Figure 17b, we display the evolution of the indicator of improvement  $\mathcal{I}_k^h$  with the discretization for the former test-case. One notices once again that this quantity is stable when the stopping criteria  $E_\beta^{\text{method}}$  are used, while it increases with the number of degrees of freedom with the conventional criterion  $R_5$ . In addition, the comparison between Figures 17a and 17b shows that the levels of  $\mathcal{I}_k^h$  have the same magnitude for both test-cases, and according to our computations on the test-cases



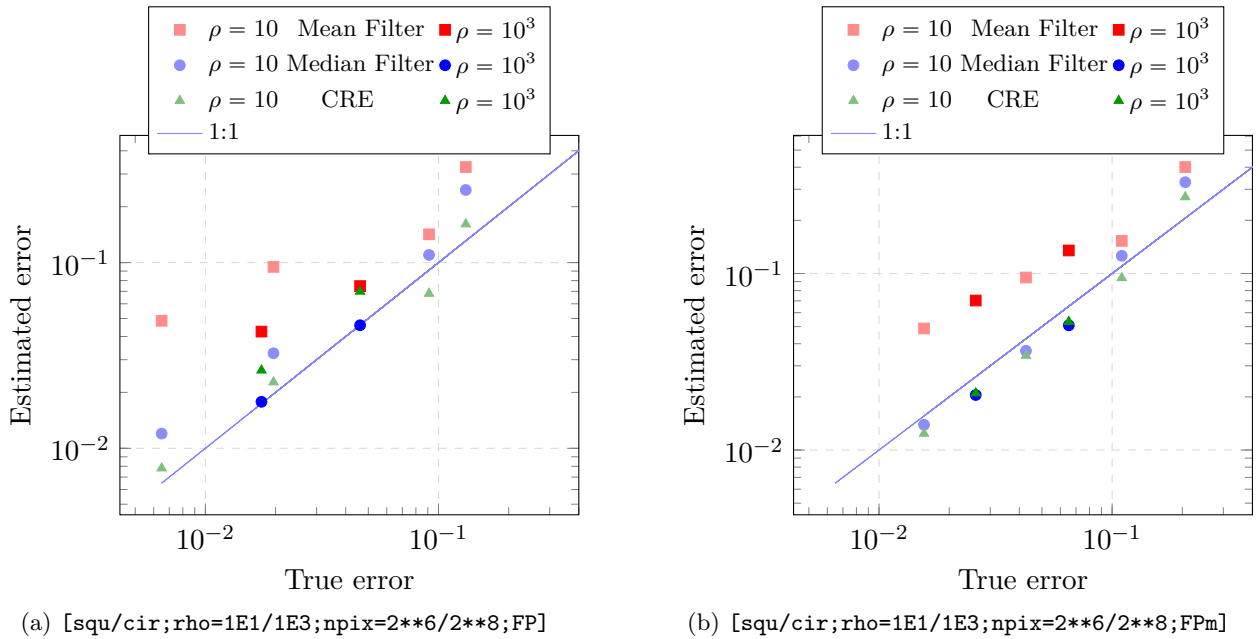


Figure 15: Estimated error  $\Delta^h$  vs true error  $\delta^h$  for different configurations and choice of numerical parameters.

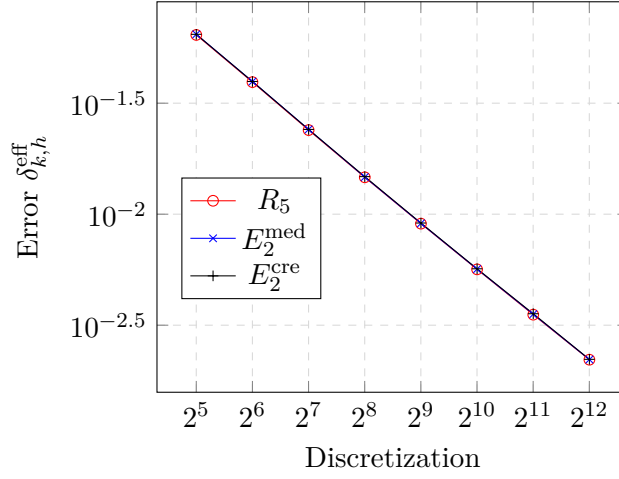
considered, this magnitude seems to be correlated with  $\beta$ . This implies that the latter parameter can be used to control the significance of the discretization in the total error on the effective property.

Figure 18 displays the evolution of the indicator  $\mathcal{I}_k^h$  against the number of pixels and for the same two geometries but at a lower conductivity ratio  $\rho = 10$ . For the criteria  $E_\beta^{\text{method}}$ , the behavior of this indicator in the case of random circular inclusions is satisfactory, with all stopping criteria leading to rather stable values of  $\mathcal{I}_k^h$ , while the criterion  $R_5$  leads to an unstable indicator, which increases with the number of pixels. The square inclusion test-case leads to slightly different behaviors: the residual-based stopping criterion  $R_5$  makes  $\mathcal{I}_k^h$  increase for more than two orders of magnitude when the number of pixels increases; the proposed criteria  $E_\beta^{\text{method}}$  also lead to an increase of  $\mathcal{I}_k^h$  but it is however smaller. This shows that, even when the proposed stopping criteria do not exhibit an optimal behavior, they however lead to a less unstable  $\mathcal{I}_k^h$ .

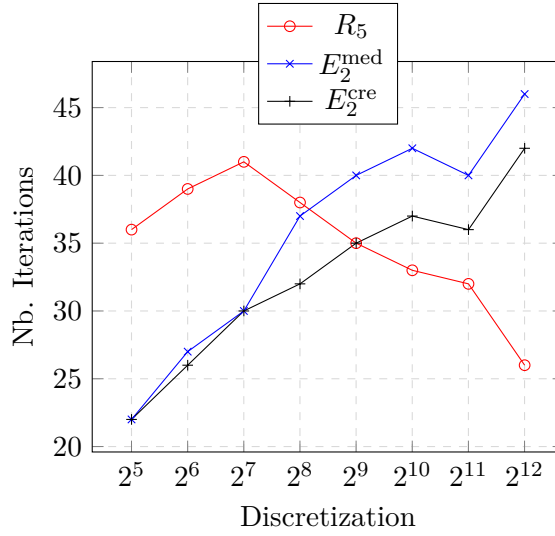
**Remark 11.** *The relatively sub-optimal behavior of the proposed stopping criteria on the test-case [sqr;rho=1E1;npix=\*;CG] of Figure 18a can be explained as follows: the discretization error estimators do not match the error on the effective property (see Remark 10) and the discrepancy seems to slightly increase with the discretization. In addition, on this test-case, all stopping criteria require only a few iterations to be met (typically 8 or 9), so that a variation of a few iterations leads to potentially significantly different values of  $\mathcal{I}_k^h$ .*

## 5.2 3D elasticity test-cases

The proposed procedure can be applied to the case of 3D linear elasticity. Adapting the proposed errors estimators is straightforward although more demanding both in terms of implementation and hardware. We consider a periodic microstructure of hard ellipsoids (Young modulus  $E = 200$  GPa) in a soft matrix ( $E = 20$  GPa for  $\rho = 10$  or  $E = 200$  MPa for  $\rho = 10^3$ ). The ellipsoids are allowed to interpenetrate. The Poisson coefficient is homogeneous and set to  $\nu = 0.3$ . The geometry of the inclusions is displayed on Figure 19a for an example discretization with  $(2^6)^3 = 262\,144$  voxels. The number of voxels used in the computations will vary between  $(2^4)^3 = 4\,096$  and  $(2^9)^3 = 134\,217\,728$ . For this computation, no Finite Element solution was computed, and the reference solution is obtained with an overkill FFT computation with  $(2^{10})^3 = 1\,073\,741\,824$  voxels. Note that the use of this overkill solution as a reference for the finest computation being evaluated is arguable. The stress  $\sigma_{xx}^h$  at



(a) Relative errors



(b) Nb. Iterations

Figure 16: Error and required number of iterations for different convergence criteria as functions of the discretization and for the different stopping criteria. [sqr;rho=1E3;npix=\*;CG]

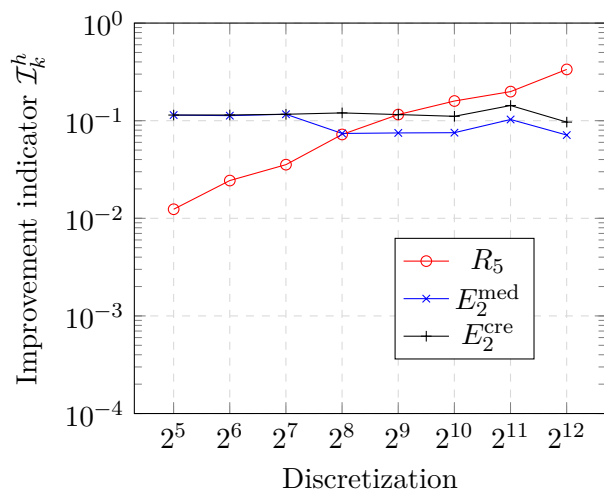
convergence (with  $k = 40$  iterations) and its transformed versions with the median and CRE methods are displayed on Figures 19b, 19c and 19d.

**Remark 12.** For the CRE estimator, a heterogeneous Poisson coefficient would require to invert several rigidity matrices  $\mathbf{K}_\nu^{h'}$  at the step of Section 4.4.3 (one for each different value of  $\nu$ ), and to apply the operator  $\mathbf{G}_\nu \mathbf{K}_\nu^{h'} + \mathbf{M}^{h'}$  independently to the voxels of each inclusion. However, this should have only a limited impact on the computational cost and implementation difficulty of the method provided there is only a small number of different homogeneous inclusions.

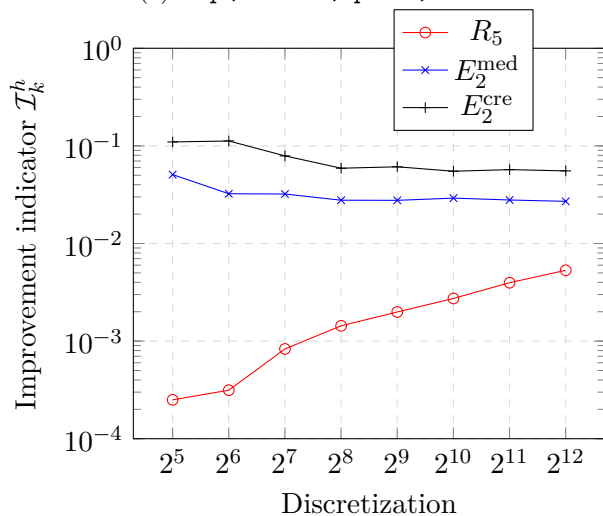
We display on Figure 20 the evolution of the improvement indicator  $\mathcal{I}_k^h$  (ratio between the iteration and total errors on the energy) with the number of voxels for the different stopping criteria. The results obtained for these numerical experiments in 3D elasticity are similar to those of Section 5.1 in that both stopping criteria  $E_\beta^{\text{med}}$  and  $E_\beta^{\text{cre}}$  lead to a ratio between iteration error and total error that is more stable than  $R_\beta$  and thus are more reliable.

## 6 Conclusion

The present study addresses the issue of the computation of a global convergence criterion for iterative solvers in the FFT-based computational homogenization of periodic materials. The main idea is to separate and estimate the contributions to the error of the iterative scheme and of the discretization.



(a) [squ;rho=1E3;npix=\*;CG]



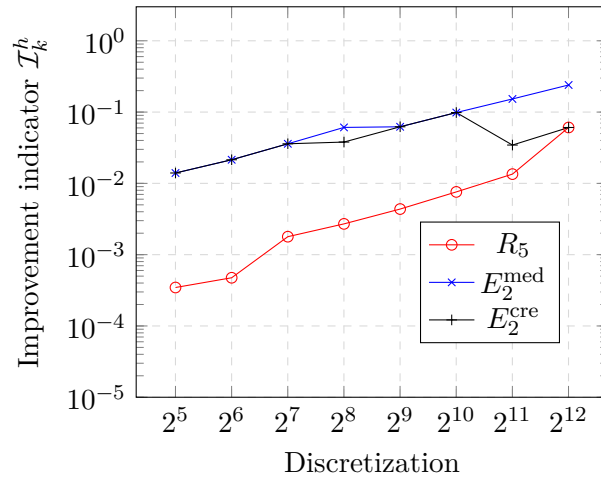
(b) [cir;rho=1E3;npix=\*;CG]

Figure 17: Indicator of improvement  $\mathcal{I}_k^h$  as a function of the discretization and for the different stopping criteria.

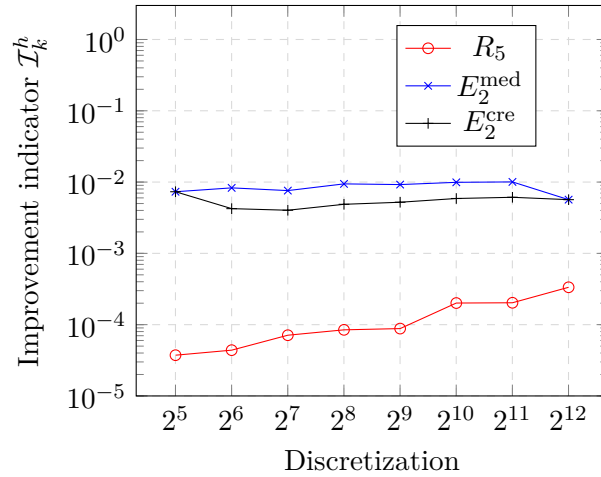
Some associated estimators have been proposed and their performances have been assessed on two prototypical 2D conductivity test-cases. From this study, we can conclude that:

- The residual of the discrete problem is not sufficient for a reliable estimation of the total error on the effective property.
- Provided that the conductivity ratio is not too high (up to  $10^3$  for the configurations considered), this residual is nonetheless a suitable estimator of the iteration error.
- Two methods have proven to be useful for estimating the discretization error. The first one is the *median filter*, which works best with rather high conductivity ratios ( $\sim 10^3$ ), while the second one, the *constitutive relation error*, appears in our computations to be the best choice for moderate conductivity ratios ( $\sim 10$ ).
- A parameter-controlled stopping criterion has been introduced and tested: it consists in interrupting the iterative scheme when the estimated iteration error is *much* smaller than the estimated discretization error, the sense of much being given by the featured parameter.

Finally, we based the evaluation of the stopping criteria on an *improvement indicator* that measures how much the error on the effective property could have been reduced if the iterative scheme had not been stopped. In the test cases considered, the proposed stopping criteria appear to be superior in several ways to the conventional criterion solely based on the residual:



(a) [sqr;rho=1E1;npix=\*;CG]



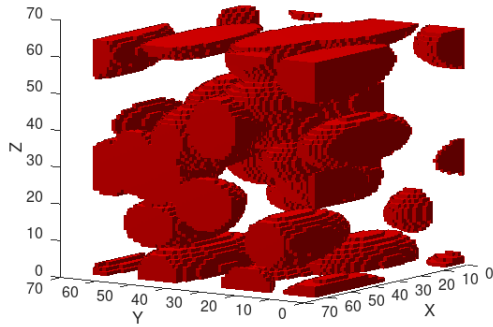
(b) [cir;rho=1E1;npix=\*;CG]

Figure 18: Indicator of improvement  $\mathcal{I}_k^h$  as a function of the discretization for the different stopping criteria.

- The residual-based criterion tends to require less iterations when the discretization becomes finer. As a consequence, the improvement indicator increases with the number of pixels. On the contrary, our stopping criteria tends to require more iterations when the discretization is finer, which leads to a more stable improvement indicator so that full advantage can be taken of the discretization.
- The proposed stopping criteria involve a parameter  $\beta$ , which value seems to control the improvement indicator associated with the computations. If confirmed, this would mean that, when choosing the value of this parameter, one could adjust the share of the iteration error in the total error on the effective property.

At this stage, a number of perspectives emerge for this work:

- The estimation of the discretization error has a non-negligible cost in the computation (despite the fact that the regular grid allows for significant speedups). For this reason, it will be essential to compute it only when necessary, i.e. possibly not at every iterations, which calls for an adapted strategy.
- The *constitutive relation method* appears to be well suited in the cases of small contrast. It could be interesting to improve it in cases of strong contrast, in particular by using a more accurate method of flux reconstruction. One could for example use the approach proposed in [23].
- There are three valuable theoretical results that do not hold due to the fact that the discretization methods we work with are inexactly integrated Galerkin methods. These results are (i) the



(a) Ellipsoidal inclusions with  $(2^6)^3$  voxels

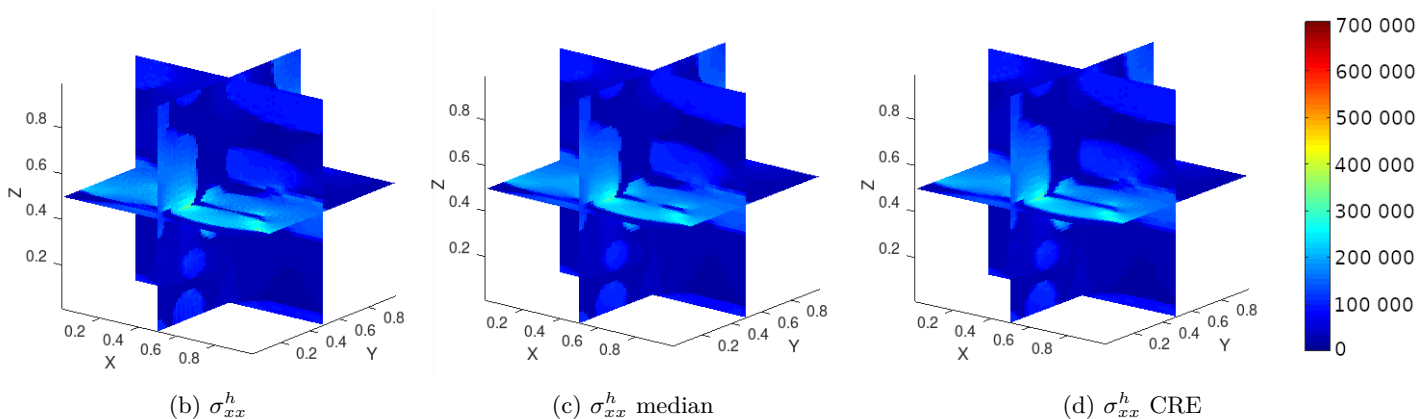


Figure 19: Geometry and stress field  $\sigma_{xx}^h$  after  $k = 40$  iterations.

equality between the energy error on the field and the error on the effective property, (ii) the bound between the computed and the exact effective properties and (iii) the bound between the exact error and the estimator provided by the CRE method. As a consequence, it would be of great interest to extend the methods proposed in the present paper to exactly integrated Galerkin schemes [30, 5] or FEM-based discretizations coupled with FFT solvers [25, 16].

## A Fourier transforms

Consider the unit cell  $\Omega$  filling the space  $\mathbb{R}^d$  by translation along  $d$  vectors  $\mathbf{Y}_1, \dots, \mathbf{Y}_d$ . The lattice  $\mathcal{R}$  generated by these vectors is defined as

$$\mathcal{R} = \left\{ \mathbf{Y} \mid \mathbf{Y} = \sum_{j=1}^d n_j \mathbf{Y}_j, n_j \in \mathbb{Z} \right\}.$$

Let  $\mathcal{R}^*$  denote the reciprocal lattice of  $\mathcal{R}$  generated by the vectors

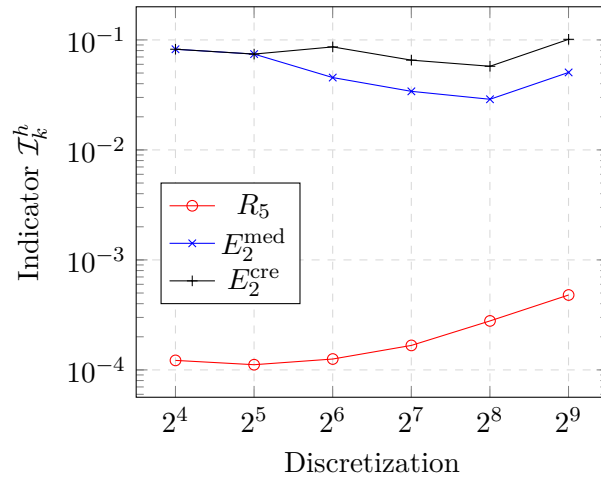
$$\mathbf{Y}_i^* = \frac{2\pi}{|\Omega|} \mathbf{Y}_j \wedge \mathbf{Y}_k,$$

where  $(i, j, k)$  is a direct circular permutation. The Fourier transform  $\hat{f}$  of  $f$  is defined on  $\mathcal{R}^*$  as:

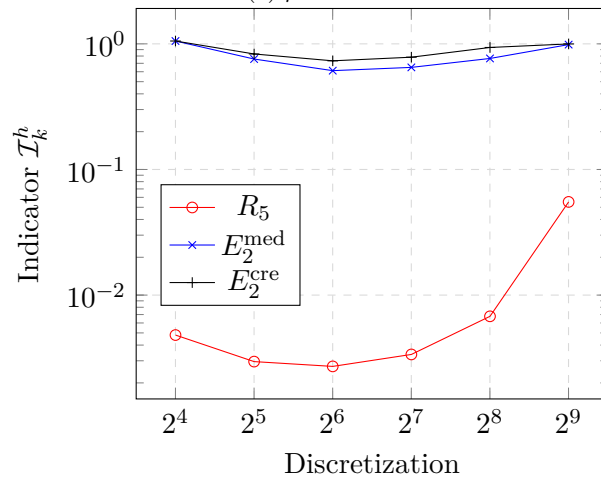
$$\hat{f}(\boldsymbol{\xi}) = \mathcal{F}[f](\boldsymbol{\xi}) = \frac{1}{|\Omega|} \int_{\Omega} f(\mathbf{x}) e^{-i\boldsymbol{\xi} \cdot \mathbf{x}} d\mathbf{x}, \quad \text{where } i = \sqrt{-1}.$$

The periodic function  $f$  in  $L_{\text{per}}^2(\Omega)$  can be reconstructed from its Fourier transform by

$$f(\mathbf{x}) = \mathcal{F}^{-1}[\hat{f}](\mathbf{x}) = \sum_{\boldsymbol{\xi} \in \mathcal{R}^*} \hat{f}(\boldsymbol{\xi}) e^{i\boldsymbol{\xi} \cdot \mathbf{x}}.$$



(a)  $\rho = 10$



(b)  $\rho = 10^3$

Figure 20: Indicator of improvement  $\mathcal{I}_k^h$  as a function of the discretization and for the different stopping criteria considering the 3D elasticity test-cases.

## References

- [1] J. C. A. Barata and M. S. Hussein. The moore–penrose pseudoinverse: A tutorial review of the theory. *Brazilian Journal of Physics*, 42(1):146–165, 2012.
- [2] C. Bellis, H. Moulinec, and P. Suquet. Eigendecomposition-based convergence analysis of the Neumann series for laminated composites and discretization error estimation. *International Journal for Numerical Methods in Engineering*, 121(2):201–232, 2020.
- [3] C. Bellis and P. Suquet. Geometric variational principles for computational homogenization. *Journal of Elasticity*, 137(2):119–149, 2019.
- [4] S. Brisard and L. Chamoin. Constitutive relation error for FFT-based methods. In *ECCOMAS Congress 2016 Proceedings. VII European Congress on Computational Methods in Applied Sciences and Engineering (ECCOMAS Congress2016)*, Greece, 2016.
- [5] S. Brisard and L. Dormieux. Combining Galerkin approximation techniques with the principle of Hashin and Shtrikman to derive a new FFT-based numerical method for the homogenization of composites. *Comput Methods Appl Mech Eng.*, 217–220:197–212, 2012.
- [6] D. J. Eyre and G. W. Milton. A fast numerical scheme for computing the response of composites using grid refinement. *The European Physical Journal Applied Physics*, 6(1):41–47, 1999.

- [7] B. Fornberg. The pseudospectral method: Comparisons with finite differences for the elastic wave equation. *Geophysics*, 52(4):483–501, 1987.
- [8] T. Huang, G. Yang, and G. Tang. A fast two-dimensional median filtering algorithm. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 27(1):13–18, 1979.
- [9] T. J. R. Hughes. *The Finite Element Method: Linear Static and Dynamic Finite Element Analysis*. Prentice-Hall, Inc., 1987.
- [10] P. Jiránek, Z. Strakoš, and M. Vohralík. A posteriori error estimates including algebraic error and stopping criteria for iterative solvers. *SIAM Journal on Scientific Computing*, 32(3):1567–1590, 2010.
- [11] M. Kabel, T. Böhlke, and M. Schneider. Efficient fixed point and Newton–Krylov solvers for FFT-based homogenization of elasticity at large deformations. *Computational Mechanics*, 54(6):1497–1514, 2014.
- [12] J. Kochmann, K. Manjunatha, C. Gierden, S. Wulfinghoff, B. Svendsen, and S. Reese. A simple and flexible model order reduction method for fft-based homogenization problems using a sparse sampling technique. *Computer Methods in Applied Mechanics and Engineering*, 347:622–638, 2019.
- [13] P. Ladeveze and D. Leguillon. Error estimate procedure in the finite element method and applications. *SIAM Journal on Numerical Analysis*, 20(3):485–509, 1983.
- [14] P. Ladevèze and J.-P. Pelle. *Mastering Calculations in Linear and Nonlinear Mechanics*. Springer, New York, 2005.
- [15] P. Ladevèze and P. Rougeot. New advances on a posteriori error on constitutive relation in f.e. analysis. *Computer Methods in Applied Mechanics and Engineering*, 150(1-4):239–249, 1997.
- [16] M. Leuschner and F. Fritzen. Fourier-accelerated nodal solvers (fans) for homogenization problems. *Comput. Mech.*, 62, 2018.
- [17] J.-C. Michel, H. Moulinec, and P. Suquet. A computational scheme for linear and non-linear composites with arbitrary phase contrast. *Int J Numer Methods Eng.*, 52:139–160, 2001.
- [18] G. W. Milton. *The Theory of Composites*. Cambridge university press, 2002.
- [19] L. Morin, R. Brenner, K. Derrien, and K. Dorhmi. Periodic smoothing splines for fft-based solvers. *Computer Methods in Applied Mechanics and Engineering*, 373:113549, 2021.
- [20] H. Moulinec and P. Suquet. A numerical method for computing the overall response of non-linear composites with complex microstructure. *Computer methods in applied mechanics and engineering*, 157(1-2):69–94, 1998.
- [21] H. Moulinec, P. Suquet, and G. W. Milton. Convergence of iterative methods based on neumann series for composite materials: Theory and practice: Convergence of iterative methods. *International Journal for Numerical Methods in Engineering*, 114(10), 2018.
- [22] Y. V. Obnosov. Periodic heterogeneous structures: new explicit solutions and effective characteristics of refraction of an imposed field. *SIAM Journal on Applied Mathematics*, 59(4):1267–1287, 1999.
- [23] A. Parret-Fréaud, V. Rey, P. Gosselet, and C. Rey. Improved recovery of admissible stress in domain decomposition methods—application to heterogeneous structures and new error bounds for feti-dp. *International Journal for Numerical Methods in Engineering*, 111(1):69–87, 2017.

- [24] V. Rey, C. Rey, and P. Gosselet. A strict error bound with separated contributions of the discretization and of the iterative solver in non-overlapping domain decomposition methods. *Computer Methods in Applied Mechanics and Engineering*, 270:293–303, 2014.
- [25] M. Schneider, D. Merkert, and M. Kabel. Fft-based homogenization for microstructures discretized by linear hexahedral elements. *International Journal for Numerical Methods in Engineering*, 109(10), 2017.
- [26] M. Schneider, F. Ospald, and M. Kabel. Computational homogenization of elasticity on a staggered grid. *Int J Numer Methods Eng.*, 105:693–720, 2016.
- [27] R. M. Sydenstricker, A. L. Coutinho, M. A. Martins, L. Landau, and J. L. Alves. A posteriori error estimate for stress analysis of homogeneous and heterogeneous materials: An engineering approach. *Finite elements in analysis and design*, 42(3):171–188, 2005.
- [28] J. Vondřejc. *FFT-based method for homogenization of periodic media: Theory and applications*. PhD thesis, Czech Technical University, 2013.
- [29] J. Vondřejc. Improved guaranteed computable bounds on homogenized properties of periodic media by the Fourier–Galerkin method with exact integration. *International Journal for Numerical Methods in Engineering*, 107(13):1106–1135, 2016.
- [30] J. Vondřejc, J. Zeman, and I. Marek. An FFT-based Galerkin method for homogenization of periodic media. *Computers & Mathematics with Applications*, 68(3):156–173, 2014.
- [31] J. Vondřejc, J. Zeman, and I. Marek. Guaranteed upper–lower bounds on homogenized properties by FFT-based Galerkin method. *Computer Methods in Applied Mechanics and Engineering*, 297:258–291, 2015.
- [32] W. H. Müller. Mathematical versus experimental stress analysis of inhomogeneities in solids. *J Phys IV*, 6:139–148, 1996.
- [33] F. Willot. Fourier-based schemes for computing the mechanical response of composites with accurate local fields. *Comptes Rendus Mécanique*, 343(3):232–245, 2015.
- [34] F. Willot, B. Abdallah, and Y.-P. Pellegrini. Fourier-based schemes with modified Green operator for computing the electrical response of heterogeneous media with accurate local fields. *International Journal for Numerical Methods in Engineering*, 98(7):518–533, 2014.
- [35] L. Yin, R. Yang, M. Gabbouj, and Y. Neuvo. Weighted median filters: a tutorial. *IEEE Transactions on Circuits and Systems II: Analog and Digital Signal Processing*, 43(3):157–192, 1996.
- [36] J. Zeman, J. Vondřejc, J. Novák, and I. Marek. Accelerating a FFT-based solver for numerical homogenization of periodic media by conjugate gradients. *Journal of Computational Physics*, 229(21):8065–8071, 2010.
- [37] O. C. Zienkiewicz and J. Z. Zhu. A simple error estimator and adaptive procedure for practical engineering analysis. *International journal for numerical methods in engineering*, 24(2):337–357, 1987.